

# ECONOMETRICS

BRUCE E. HANSEN

©2000, 2018<sup>1</sup>

**University of Wisconsin**

**Department of Economics**

This Revision: January 2018

Comments Welcome

<sup>1</sup>This manuscript may be printed and reproduced for individual or instructional use, but may not be printed for commercial purposes.

# Contents

Preface . . . . .	x
<b>1 Introduction . . . . .</b>	<b>1</b>
1.1 What is Econometrics? . . . . .	1
1.2 The Probability Approach to Econometrics . . . . .	1
1.3 Econometric Terms and Notation . . . . .	2
1.4 Observational Data . . . . .	3
1.5 Standard Data Structures . . . . .	4
1.6 Sources for Economic Data . . . . .	6
1.7 Econometric Software . . . . .	7
1.8 Data Files for Textbook . . . . .	7
1.9 Reading the Manuscript . . . . .	8
1.10 Common Symbols . . . . .	10
<b>2 Conditional Expectation and Projection . . . . .</b>	<b>11</b>
2.1 Introduction . . . . .	11
2.2 The Distribution of Wages . . . . .	11
2.3 Conditional Expectation . . . . .	13
2.4 Log Differences* . . . . .	15
2.5 Conditional Expectation Function . . . . .	16
2.6 Continuous Variables . . . . .	17
2.7 Law of Iterated Expectations . . . . .	18
2.8 CEF Error . . . . .	20
2.9 Intercept-Only Model . . . . .	21
2.10 Regression Variance . . . . .	22
2.11 Best Predictor . . . . .	22
2.12 Conditional Variance . . . . .	23
2.13 Homoskedasticity and Heteroskedasticity . . . . .	25
2.14 Regression Derivative . . . . .	26
2.15 Linear CEF . . . . .	26
2.16 Linear CEF with Nonlinear Effects . . . . .	28
2.17 Linear CEF with Dummy Variables . . . . .	28
2.18 Best Linear Predictor . . . . .	30
2.19 Linear Predictor Error Variance . . . . .	36
2.20 Regression Coefficients . . . . .	37
2.21 Regression Sub-Vectors . . . . .	37
2.22 Coefficient Decomposition . . . . .	38
2.23 Omitted Variable Bias . . . . .	39
2.24 Best Linear Approximation . . . . .	40
2.25 Regression to the Mean . . . . .	41
2.26 Reverse Regression . . . . .	42
2.27 Limitations of the Best Linear Projection . . . . .	43

2.28	Random Coefficient Model . . . . .	43
2.29	Causal Effects . . . . .	45
2.30	Expectation: Mathematical Details* . . . . .	50
2.31	Moment Generating and Characteristic Functions* . . . . .	52
2.32	Existence and Uniqueness of the Conditional Expectation* . . . . .	52
2.33	Identification* . . . . .	53
2.34	Technical Proofs* . . . . .	54
	Exercises . . . . .	58
<b>3</b>	<b>The Algebra of Least Squares</b>	<b>61</b>
3.1	Introduction . . . . .	61
3.2	Samples . . . . .	61
3.3	Moment Estimators . . . . .	62
3.4	Least Squares Estimator . . . . .	63
3.5	Solving for Least Squares with One Regressor . . . . .	65
3.6	Solving for Least Squares with Multiple Regressors . . . . .	65
3.7	Illustration . . . . .	67
3.8	Least Squares Residuals . . . . .	68
3.9	Demeaned Regressors . . . . .	69
3.10	Model in Matrix Notation . . . . .	69
3.11	Projection Matrix . . . . .	71
3.12	Orthogonal Projection . . . . .	72
3.13	Estimation of Error Variance . . . . .	73
3.14	Analysis of Variance . . . . .	74
3.15	Regression Components . . . . .	74
3.16	Residual Regression . . . . .	76
3.17	Prediction Errors . . . . .	77
3.18	Influential Observations . . . . .	78
3.19	CPS Data Set . . . . .	80
3.20	Programming . . . . .	81
3.21	Technical Proofs* . . . . .	84
	Exercises . . . . .	85
<b>4</b>	<b>Least Squares Regression</b>	<b>88</b>
4.1	Introduction . . . . .	88
4.2	Random Sampling . . . . .	88
4.3	Sample Mean . . . . .	89
4.4	Linear Regression Model . . . . .	90
4.5	Mean of Least-Squares Estimator . . . . .	90
4.6	Variance of Least Squares Estimator . . . . .	92
4.7	Gauss-Markov Theorem . . . . .	94
4.8	Generalized Least Squares . . . . .	95
4.9	Residuals . . . . .	96
4.10	Estimation of Error Variance . . . . .	97
4.11	Mean-Square Forecast Error . . . . .	99
4.12	Covariance Matrix Estimation Under Homoskedasticity . . . . .	100
4.13	Covariance Matrix Estimation Under Heteroskedasticity . . . . .	101
4.14	Standard Errors . . . . .	104
4.15	Covariance Matrix Estimation with Sparse Dummy Variables . . . . .	105
4.16	Computation . . . . .	106
4.17	Measures of Fit . . . . .	107
4.18	Empirical Example . . . . .	108

4.19	Multicollinearity . . . . .	110
4.20	Clustered Sampling . . . . .	113
4.21	Inference with Clustered Samples . . . . .	119
	Exercises . . . . .	121
<b>5</b>	<b>Normal Regression and Maximum Likelihood</b>	<b>126</b>
5.1	Introduction . . . . .	126
5.2	The Normal Distribution . . . . .	126
5.3	Chi-Square Distribution . . . . .	128
5.4	Student t Distribution . . . . .	129
5.5	F Distribution . . . . .	130
5.6	Joint Normality and Linear Regression . . . . .	131
5.7	Normal Regression Model . . . . .	131
5.8	Distribution of OLS Coefficient Vector . . . . .	133
5.9	Distribution of OLS Residual Vector . . . . .	134
5.10	Distribution of Variance Estimate . . . . .	135
5.11	t-statistic . . . . .	135
5.12	Confidence Intervals for Regression Coefficients . . . . .	136
5.13	Confidence Intervals for Error Variance . . . . .	138
5.14	t Test . . . . .	138
5.15	Likelihood Ratio Test . . . . .	140
5.16	Likelihood Properties . . . . .	142
5.17	Information Bound for Normal Regression . . . . .	143
5.18	Gamma Function* . . . . .	144
5.19	Technical Proofs* . . . . .	145
<b>6</b>	<b>An Introduction to Large Sample Asymptotics</b>	<b>154</b>
6.1	Introduction . . . . .	154
6.2	Asymptotic Limits . . . . .	155
6.3	Convergence in Probability . . . . .	156
6.4	Weak Law of Large Numbers . . . . .	157
6.5	Almost Sure Convergence and the Strong Law* . . . . .	158
6.6	Vector-Valued Moments . . . . .	159
6.7	Convergence in Distribution . . . . .	160
6.8	Central Limit Theorem . . . . .	161
6.9	Multivariate Central Limit Theorem . . . . .	164
6.10	Higher Moments . . . . .	165
6.11	Functions of Moments . . . . .	166
6.12	Delta Method . . . . .	168
6.13	Stochastic Order Symbols . . . . .	169
6.14	Uniform Stochastic Bounds* . . . . .	171
6.15	Semiparametric Efficiency . . . . .	171
6.16	Technical Proofs* . . . . .	174
	Exercises . . . . .	181
<b>7</b>	<b>Asymptotic Theory for Least Squares</b>	<b>184</b>
7.1	Introduction . . . . .	184
7.2	Consistency of Least-Squares Estimator . . . . .	185
7.3	Asymptotic Normality . . . . .	186
7.4	Joint Distribution . . . . .	189
7.5	Consistency of Error Variance Estimators . . . . .	193

7.6	Homoskedastic Covariance Matrix Estimation . . . . .	194
7.7	Heteroskedastic Covariance Matrix Estimation . . . . .	194
7.8	Summary of Covariance Matrix Notation . . . . .	196
7.9	Alternative Covariance Matrix Estimators* . . . . .	197
7.10	Functions of Parameters . . . . .	198
7.11	Asymptotic Standard Errors . . . . .	200
7.12	t-statistic . . . . .	202
7.13	Confidence Intervals . . . . .	203
7.14	Regression Intervals . . . . .	205
7.15	Forecast Intervals . . . . .	206
7.16	Wald Statistic . . . . .	208
7.17	Homoskedastic Wald Statistic . . . . .	208
7.18	Confidence Regions . . . . .	209
7.19	Semiparametric Efficiency in the Projection Model . . . . .	210
7.20	Semiparametric Efficiency in the Homoskedastic Regression Model* . . . . .	211
7.21	Uniformly Consistent Residuals* . . . . .	213
7.22	Asymptotic Leverage* . . . . .	214
	Exercises . . . . .	216
<b>8</b>	<b>Restricted Estimation</b>	<b>223</b>
8.1	Introduction . . . . .	223
8.2	Constrained Least Squares . . . . .	224
8.3	Exclusion Restriction . . . . .	225
8.4	Finite Sample Properties . . . . .	225
8.5	Minimum Distance . . . . .	228
8.6	Asymptotic Distribution . . . . .	229
8.7	Variance Estimation and Standard Errors . . . . .	231
8.8	Efficient Minimum Distance Estimator . . . . .	231
8.9	Exclusion Restriction Revisited . . . . .	232
8.10	Variance and Standard Error Estimation . . . . .	234
8.11	Hausman Equality . . . . .	234
8.12	Example: Mankiw, Romer and Weil (1992) . . . . .	235
8.13	Misspecification . . . . .	239
8.14	Nonlinear Constraints . . . . .	241
8.15	Inequality Restrictions . . . . .	242
8.16	Technical Proofs* . . . . .	243
	Exercises . . . . .	245
<b>9</b>	<b>Hypothesis Testing</b>	<b>248</b>
9.1	Hypotheses . . . . .	248
9.2	Acceptance and Rejection . . . . .	249
9.3	Type I Error . . . . .	250
9.4	t tests . . . . .	250
9.5	Type II Error and Power . . . . .	252
9.6	Statistical Significance . . . . .	252
9.7	P-Values . . . . .	253
9.8	t-ratios and the Abuse of Testing . . . . .	255
9.9	Wald Tests . . . . .	256
9.10	Homoskedastic Wald Tests . . . . .	258
9.11	Criterion-Based Tests . . . . .	259
9.12	Minimum Distance Tests . . . . .	259
9.13	Minimum Distance Tests Under Homoskedasticity . . . . .	260

9.14 F Tests . . . . .	261
9.15 Hausman Tests . . . . .	262
9.16 Score Tests . . . . .	264
9.17 Problems with Tests of Nonlinear Hypotheses . . . . .	265
9.18 Monte Carlo Simulation . . . . .	268
9.19 Confidence Intervals by Test Inversion . . . . .	271
9.20 Multiple Tests and Bonferroni Corrections . . . . .	272
9.21 Power and Test Consistency . . . . .	273
9.22 Asymptotic Local Power . . . . .	274
9.23 Asymptotic Local Power, Vector Case . . . . .	277
9.24 Technical Proofs* . . . . .	278
Exercises . . . . .	280
<b>10 Multivariate Regression</b>	<b>287</b>
10.1 Introduction . . . . .	287
10.2 Regression Systems . . . . .	287
10.3 Least-Squares Estimator . . . . .	288
10.4 Mean and Variance of Systems Least-Squares . . . . .	290
10.5 Asymptotic Distribution . . . . .	291
10.6 Covariance Matrix Estimation . . . . .	293
10.7 Seemingly Unrelated Regression . . . . .	293
10.8 Maximum Likelihood Estimator . . . . .	295
10.9 Reduced Rank Regression . . . . .	296
Exercises . . . . .	300
<b>11 Instrumental Variables</b>	<b>302</b>
11.1 Introduction . . . . .	302
11.2 Examples . . . . .	303
11.3 Instrumental Variables . . . . .	304
11.4 Example: College Proximity . . . . .	306
11.5 Reduced Form . . . . .	308
11.6 Reduced Form Estimation . . . . .	309
11.7 Identification . . . . .	310
11.8 Instrumental Variables Estimator . . . . .	311
11.9 Demeaned Representation . . . . .	313
11.10 Wald Estimator . . . . .	314
11.11 Two-Stage Least Squares . . . . .	315
11.12 Limited Information Maximum Likelihood . . . . .	318
11.13 Consistency of 2SLS . . . . .	321
11.14 Asymptotic Distribution of 2SLS . . . . .	322
11.15 Determinants of 2SLS Variance . . . . .	324
11.16 Covariance Matrix Estimation . . . . .	324
11.17 Asymptotic Distribution and Covariance Estimation for LIML . . . . .	326
11.18 Functions of Parameters . . . . .	327
11.19 Hypothesis Tests . . . . .	328
11.20 Finite Sample Theory . . . . .	328
11.21 Clustered Dependence . . . . .	329
11.22 Generated Regressors . . . . .	329
11.23 Regression with Expectation Errors . . . . .	333
11.24 Control Function Regression . . . . .	335
11.25 Endogeneity Tests . . . . .	338
11.26 Subset Endogeneity Tests . . . . .	341

11.27	OverIdentification Tests	343
11.28	Subset OverIdentification Tests	345
11.29	Local Average Treatment Effects	348
11.30	Identification Failure	351
11.31	Weak Instruments	352
11.32	Weak Instruments with $k_2 > 1$	359
11.33	Many Instruments	362
11.34	Example: Acemoglu, Johnson and Robinson (2001)	365
11.35	Example: Angrist and Krueger (1991)	366
11.36	Programming	369
	Exercises	371
<b>12</b>	<b>Generalized Method of Moments</b>	<b>379</b>
12.1	Moment Equation Models	379
12.2	Method of Moments Estimators	379
12.3	Overidentified Moment Equations	381
12.4	Linear Moment Models	382
12.5	GMM Estimator	382
12.6	Distribution of GMM Estimator	383
12.7	Efficient GMM	384
12.8	Efficient GMM versus 2SLS	385
12.9	Estimation of the Efficient Weight Matrix	385
12.10	Iterated GMM	386
12.11	Covariance Matrix Estimation	387
12.12	Clustered Dependence	387
12.13	Wald Test	388
12.14	Restricted GMM	389
12.15	Constrained Regression	390
12.16	Distance Test	391
12.17	Continuously-Updated GMM	392
12.18	OverIdentification Test	393
12.19	Subset OverIdentification Tests	394
12.20	Endogeneity Test	395
12.21	Subset Endogeneity Test	395
12.22	GMM: The General Case	396
12.23	Conditional Moment Equation Models	397
12.24	Technical Proofs*	399
	Exercises	401
<b>13</b>	<b>The Bootstrap</b>	<b>407</b>
13.1	Definition of the Bootstrap	407
13.2	The Empirical Distribution Function	407
13.3	Nonparametric Bootstrap	409
13.4	Bootstrap Estimation of Bias and Variance	409
13.5	Percentile Intervals	410
13.6	Percentile-t Equal-Tailed Interval	412
13.7	Symmetric Percentile-t Intervals	412
13.8	Asymptotic Expansions	413
13.9	One-Sided Tests	415
13.10	Symmetric Two-Sided Tests	415
13.11	Percentile Confidence Intervals	417
13.12	Bootstrap Methods for Regression Models	417

13.13 Bootstrap GMM Inference . . . . .	418
Exercises . . . . .	420
<b>14 Univariate Time Series</b>	<b>424</b>
14.1 Stationarity and Ergodicity . . . . .	424
14.2 Autoregressions . . . . .	426
14.3 Stationarity of AR(1) Process . . . . .	427
14.4 Lag Operator . . . . .	427
14.5 Stationarity of AR(k) . . . . .	428
14.6 Estimation . . . . .	428
14.7 Asymptotic Distribution . . . . .	429
14.8 Bootstrap for Autoregressions . . . . .	430
14.9 Trend Stationarity . . . . .	430
14.10 Testing for Omitted Serial Correlation . . . . .	431
14.11 Model Selection . . . . .	432
14.12 Autoregressive Unit Roots . . . . .	432
<b>15 Multivariate Time Series</b>	<b>434</b>
15.1 Vector Autoregressions (VARs) . . . . .	434
15.2 Estimation . . . . .	435
15.3 Restricted VARs . . . . .	435
15.4 Single Equation from a VAR . . . . .	435
15.5 Testing for Omitted Serial Correlation . . . . .	436
15.6 Selection of Lag Length in an VAR . . . . .	436
15.7 Granger Causality . . . . .	437
15.8 Cointegration . . . . .	437
15.9 Cointegrated VARs . . . . .	438
<b>16 Panel Data</b>	<b>440</b>
16.1 Individual-Effects Model . . . . .	440
16.2 Fixed Effects . . . . .	440
16.3 Dynamic Panel Regression . . . . .	442
Exercises . . . . .	443
<b>17 NonParametric Regression</b>	<b>444</b>
17.1 Introduction . . . . .	444
17.2 Binned Estimator . . . . .	444
17.3 Kernel Regression . . . . .	446
17.4 Local Linear Estimator . . . . .	447
17.5 Nonparametric Residuals and Regression Fit . . . . .	448
17.6 Cross-Validation Bandwidth Selection . . . . .	450
17.7 Asymptotic Distribution . . . . .	453
17.8 Conditional Variance Estimation . . . . .	455
17.9 Standard Errors . . . . .	456
17.10 Multiple Regressors . . . . .	457
<b>18 Series Estimation</b>	<b>459</b>
18.1 Approximation by Series . . . . .	459
18.2 Splines . . . . .	459
18.3 Partially Linear Model . . . . .	461
18.4 Additively Separable Models . . . . .	461
18.5 Uniform Approximations . . . . .	461



18.6	Runge's Phenomenon . . . . .	463
18.7	Approximating Regression . . . . .	463
18.8	Residuals and Regression Fit . . . . .	466
18.9	Cross-Validation Model Selection . . . . .	466
18.10	Convergence in Mean-Square . . . . .	467
18.11	Uniform Convergence . . . . .	468
18.12	Asymptotic Normality . . . . .	469
18.13	Asymptotic Normality with Undersmoothing . . . . .	470
18.14	Regression Estimation . . . . .	471
18.15	Kernel Versus Series Regression . . . . .	472
18.16	Technical Proofs . . . . .	472
	Exercises . . . . .	478
<b>19</b>	<b>Empirical Likelihood</b>	<b>479</b>
19.1	Non-Parametric Likelihood . . . . .	479
19.2	Asymptotic Distribution of EL Estimator . . . . .	481
19.3	Overidentifying Restrictions . . . . .	482
19.4	Testing . . . . .	483
19.5	Numerical Computation . . . . .	484
<b>20</b>	<b>Regression Extensions</b>	<b>486</b>
20.1	Nonlinear Least Squares . . . . .	486
20.2	Generalized Least Squares . . . . .	489
20.3	Testing for Heteroskedasticity . . . . .	492
20.4	Testing for Omitted Nonlinearity . . . . .	492
20.5	Least Absolute Deviations . . . . .	493
20.6	Quantile Regression . . . . .	495
	Exercises . . . . .	498
<b>21</b>	<b>Limited Dependent Variables</b>	<b>500</b>
21.1	Binary Choice . . . . .	500
21.2	Count Data . . . . .	501
21.3	Censored Data . . . . .	502
21.4	Sample Selection . . . . .	503
	Exercises . . . . .	505
<b>22</b>	<b>Nonparametric Density Estimation</b>	<b>506</b>
22.1	Kernel Density Estimation . . . . .	506
22.2	Asymptotic MSE for Kernel Estimates . . . . .	507
<b>A</b>	<b>Matrix Algebra</b>	<b>510</b>
A.1	Notation . . . . .	510
A.2	Complex Matrices* . . . . .	511
A.3	Matrix Addition . . . . .	511
A.4	Matrix Multiplication . . . . .	512
A.5	Trace . . . . .	513
A.6	Rank and Inverse . . . . .	513
A.7	Determinant . . . . .	515
A.8	Eigenvalues . . . . .	516
A.9	Positive Definite Matrices . . . . .	517
A.10	Generalized Eigenvalues . . . . .	517
A.11	Extrema of Quadratic Forms . . . . .	518

A.12 Idempotent Matrices . . . . .	520
A.13 Singular Values . . . . .	521
A.14 Cholesky Decomposition . . . . .	521
A.15 Matrix Calculus . . . . .	522
A.16 Kronecker Products and the Vec Operator . . . . .	523
A.17 Vector Norms . . . . .	524
A.18 Matrix Norms . . . . .	527
A.19 Matrix Inequalities . . . . .	529
 <b>B Probability Inequalities</b>	 <b>532</b>

# Preface

This book is intended to serve as the textbook a first-year graduate course in econometrics.

Students are assumed to have an understanding of multivariate calculus, probability theory, linear algebra, and mathematical statistics. A prior course in undergraduate econometrics would be helpful, but not required. Two excellent undergraduate textbooks are Wooldridge (2015) and Stock and Watson (2014).

For reference, some of the basic tools of matrix algebra and probability inequalities are reviewed in the Appendix.

For students wishing to deepen their knowledge of matrix algebra in relation to their study of econometrics, I recommend *Matrix Algebra* by Abadir and Magnus (2005).

An excellent introduction to probability and statistics is *Statistical Inference* by Casella and Berger (2002). For those wanting a deeper foundation in probability, I recommend Ash (1972) or Billingsley (1995). For more advanced statistical theory, I recommend Lehmann and Casella (1998), van der Vaart (1998), Shao (2003), and Lehmann and Romano (2005).

For further study in econometrics beyond this text, I recommend Davidson (1994) for asymptotic theory, Hamilton (1994) and Kilian and Lütkepohl (2017) for time-series methods, Wooldridge (2010) for panel data and discrete response models, and Li and Racine (2007) for nonparametrics and semiparametric econometrics. Beyond these texts, the *Handbook of Econometrics* series provides advanced summaries of contemporary econometric methods and theory.

The end-of-chapter exercises are important parts of the text and are meant to help teach students of econometrics. Answers are not provided, and this is intentional.

I would like to thank Ying-Ying Lee and Wooyoung Kim for providing research assistance in preparing some of the empirical examples presented in the text.

This is a manuscript in progress. Chapters 1-11 are mostly complete. Chapters 12-18 are incomplete.

# Chapter 1

## Introduction

### 1.1 What is Econometrics?

The term “econometrics” is believed to have been crafted by Ragnar Frisch (1895-1973) of Norway, one of the three principal founders of the Econometric Society, first editor of the journal *Econometrica*, and co-winner of the first Nobel Memorial Prize in Economic Sciences in 1969. It is therefore fitting that we turn to Frisch’s own words in the introduction to the first issue of *Econometrica* to describe the discipline.

A word of explanation regarding the term econometrics may be in order. Its definition is implied in the statement of the scope of the [Econometric] Society, in Section I of the Constitution, which reads: “The Econometric Society is an international society for the advancement of economic theory in its relation to statistics and mathematics.... Its main object shall be to promote studies that aim at a unification of the theoretical-quantitative and the empirical-quantitative approach to economic problems....”

But there are several aspects of the quantitative approach to economics, and no single one of these aspects, taken by itself, should be confounded with econometrics. Thus, econometrics is by no means the same as economic statistics. Nor is it identical with what we call general economic theory, although a considerable portion of this theory has a definitely quantitative character. Nor should econometrics be taken as synonymous with the application of mathematics to economics. Experience has shown that each of these three view-points, that of statistics, economic theory, and mathematics, is a necessary, but not by itself a sufficient, condition for a real understanding of the quantitative relations in modern economic life. It is the *unification* of all three that is powerful. And it is this unification that constitutes econometrics.

Ragnar Frisch, *Econometrica*, (1933), 1, pp. 1-2.

This definition remains valid today, although some terms have evolved somewhat in their usage. Today, we would say that econometrics is the unified study of economic models, mathematical statistics, and economic data.

Within the field of econometrics there are sub-divisions and specializations. **Econometric theory** concerns the development of tools and methods, and the study of the properties of econometric methods. **Applied econometrics** is a term describing the development of quantitative economic models and the application of econometric methods to these models using economic data.

### 1.2 The Probability Approach to Econometrics

The unifying methodology of modern econometrics was articulated by Trygve Haavelmo (1911-1999) of Norway, winner of the 1989 Nobel Memorial Prize in Economic Sciences, in his seminal

paper “The probability approach in econometrics” (1944). Haavelmo argued that quantitative economic models must necessarily be *probability models* (by which today we would mean *stochastic*). Deterministic models are blatantly inconsistent with observed economic quantities, and it is incoherent to apply deterministic models to non-deterministic data. Economic models should be explicitly designed to incorporate randomness; stochastic errors should not be simply added to deterministic models to make them random. Once we acknowledge that an economic model is a probability model, it follows naturally that an appropriate tool way to quantify, estimate, and conduct inferences about the economy is through the powerful theory of mathematical statistics. The appropriate method for a quantitative economic analysis follows from the probabilistic construction of the economic model.

Haavelmo’s probability approach was quickly embraced by the economics profession. Today no quantitative work in economics shuns its fundamental vision.

While all economists embrace the probability approach, there has been some evolution in its implementation.

The **structural approach** is the closest to Haavelmo’s original idea. A probabilistic economic model is specified, and the quantitative analysis performed under the assumption that the economic model is correctly specified. Researchers often describe this as “taking their model seriously.” The structural approach typically leads to likelihood-based analysis, including maximum likelihood and Bayesian estimation.

A criticism of the structural approach is that it is misleading to treat an economic model as correctly specified. Rather, it is more accurate to view a model as a useful abstraction or approximation. In this case, how should we interpret structural econometric analysis? The **quasi-structural approach** to inference views a structural economic model as an approximation rather than the truth. This theory has led to the concepts of the pseudo-true value (the parameter value defined by the estimation problem), the quasi-likelihood function, quasi-MLE, and quasi-likelihood inference.

Closely related is the **semiparametric approach**. A probabilistic economic model is partially specified but some features are left unspecified. This approach typically leads to estimation methods such as least-squares and the Generalized Method of Moments. The semiparametric approach dominates contemporary econometrics, and is the main focus of this textbook.

Another branch of quantitative structural economics is the **calibration approach**. Similar to the quasi-structural approach, the calibration approach interprets structural models as approximations and hence inherently false. The difference is that the calibrationist literature rejects mathematical statistics (deeming classical theory as inappropriate for approximate models) and instead selects parameters by matching model and data moments using non-statistical *ad hoc*<sup>1</sup> methods.

## 1.3 Econometric Terms and Notation

In a typical application, an econometrician has a set of repeated measurements on a set of variables. For example, in a labor application the variables could include weekly earnings, educational attainment, age, and other descriptive characteristics. We call this information the **data**, **dataset**, or **sample**.

We use the term **observations** to refer to the distinct repeated measurements on the variables. An individual observation often corresponds to a specific economic unit, such as a person, household, corporation, firm, organization, country, state, city or other geographical region. An individual observation could also be a measurement at a point in time, such as quarterly GDP or a daily interest rate.

---

<sup>1</sup> *Ad hoc* means “for this purpose” – a method designed for a specific problem – and not based on a generalizable principle.

Economists typically denote variables by the italicized roman characters  $y$ ,  $x$ , and/or  $z$ . The convention in econometrics is to use the character  $y$  to denote the variable to be explained, while the characters  $x$  and  $z$  are used to denote the conditioning (explaining) variables.

Following mathematical convention, real numbers (elements of the real line  $\mathbb{R}$ , also called **scalars**) are written using lower case italics such as  $y$ , and vectors (elements of  $\mathbb{R}^k$ ) by lower case bold italics such as  $\mathbf{x}$ , e.g.

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{pmatrix}.$$

Upper case bold italics such as  $\mathbf{X}$  are used for matrices.

We denote the number of observations by the natural number  $n$ , and subscript the variables by the index  $i$  to denote the individual observation, e.g.  $y_i$ ,  $\mathbf{x}_i$  and  $\mathbf{z}_i$ . In some contexts we use indices other than  $i$ , such as in time-series applications where the index  $t$  is common and  $T$  is used to denote the number of observations. In panel studies we typically use the double index  $it$  to refer to individual  $i$  at a time period  $t$ .

The  $i^{th}$  **observation** is the set  $(y_i, \mathbf{x}_i, \mathbf{z}_i)$ . The **sample** is the set  $\{(y_i, \mathbf{x}_i, \mathbf{z}_i) : i = 1, \dots, n\}$ .

It is proper mathematical practice to use upper case  $X$  for random variables and lower case  $x$  for realizations or specific values. Since we use upper case to denote matrices, the distinction between random variables and their realizations is not rigorously followed in econometric notation. Thus the notation  $y_i$  will in some places refer to a random variable, and in other places a specific realization. This is undesirable but there is little to be done about it without terrifically complicating the notation. Hopefully there will be no confusion as the use should be evident from the context.

We typically use Greek letters such as  $\beta$ ,  $\theta$  and  $\sigma^2$  to denote unknown parameters of an econometric model, and will use boldface, e.g.  $\boldsymbol{\beta}$  or  $\boldsymbol{\theta}$ , when these are vector-valued. Estimates are typically denoted by putting a hat “ $\hat{\cdot}$ ”, tilde “ $\tilde{\cdot}$ ” or bar “ $\bar{\cdot}$ ” over the corresponding letter, e.g.  $\hat{\beta}$  and  $\tilde{\beta}$  are estimates of  $\beta$ .

The covariance matrix of an econometric estimator will typically be written using the capital boldface  $\mathbf{V}$ , often with a subscript to denote the estimator, e.g.  $\mathbf{V}_{\hat{\beta}} = \text{var}(\hat{\beta})$  as the covariance matrix for  $\hat{\beta}$ . Hopefully without causing confusion, we will use the notation  $\mathbf{V}_{\beta} = \text{avar}(\hat{\beta})$  to denote the asymptotic covariance matrix of  $\sqrt{n}(\hat{\beta} - \beta)$  (the variance of the asymptotic distribution). Estimates will be denoted by appending hats or tildes, e.g.  $\hat{\mathbf{V}}_{\beta}$  is an estimate of  $\mathbf{V}_{\beta}$ .

## 1.4 Observational Data

A common econometric question is to quantify the impact of one set of variables on another variable. For example, a concern in labor economics is the returns to schooling – the change in earnings induced by increasing a worker’s education, holding other variables constant. Another issue of interest is the earnings gap between men and women.

Ideally, we would use **experimental** data to answer these questions. To measure the returns to schooling, an experiment might randomly divide children into groups, mandate different levels of education to the different groups, and then follow the children’s wage path after they mature and enter the labor force. The differences between the groups would be direct measurements of the effects of different levels of education. However, experiments such as this would be widely

condemned as immoral! Consequently, in economics non-laboratory experimental data sets are typically narrow in scope.

Instead, most economic data is **observational**. To continue the above example, through data collection we can record the level of a person's education and their wage. With such data we can measure the joint distribution of these variables, and assess the joint dependence. But from observational data it is difficult to infer **causality**, as we are not able to manipulate one variable to see the direct effect on the other. For example, a person's level of education is (at least partially) determined by that person's choices. These factors are likely to be affected by their personal abilities and attitudes towards work. The fact that a person is highly educated suggests a high level of ability, which suggests a high relative wage. This is an alternative explanation for an observed positive correlation between educational levels and wages. High ability individuals do better in school, and therefore choose to attain higher levels of education, and their high ability is the fundamental reason for their high wages. The point is that multiple explanations are consistent with a positive correlation between schooling levels and education. Knowledge of the joint distribution alone may not be able to distinguish between these explanations.

Most economic data sets are observational, not experimental. This means that all variables must be treated as random and possibly jointly determined.

This discussion means that it is difficult to infer causality from observational data alone. Causal inference requires identification, and this is based on strong assumptions. We will discuss these issues on occasion throughout the text.

## 1.5 Standard Data Structures

There are five major types of economic data sets: cross-sectional, time-series, panel, clustered, and spatial. They are distinguished by the dependence structure across observations.

Cross-sectional data sets have one observation per individual. Surveys and administrative records are a typical source for cross-sectional data. In typical applications, the individuals surveyed are persons, households, firms or other economic agents. In many contemporary econometric cross-section studies the sample size  $n$  is quite large. It is conventional to assume that cross-sectional observations are mutually independent. Most of this text is devoted to the study of cross-section data.

Time-series data are indexed by time. Typical examples include macroeconomic aggregates, prices and interest rates. This type of data is characterized by serial dependence. Most aggregate economic data is only available at a low frequency (annual, quarterly or perhaps monthly) so the sample size is typically much smaller than in cross-section studies. An exception is financial data where data are available at a high frequency (weekly, daily, hourly, or by transaction) so sample sizes can be quite large.

Panel data combines elements of cross-section and time-series. These data sets consist of a set of individuals (typically persons, households, or corporations) measured repeatedly over time. The common modeling assumption is that the individuals are mutually independent of one another, but a given individual's observations are mutually dependent. In some panel data contexts, the number of time series observations  $T$  per individual is small while the number of individuals  $n$  is large. In other panel data contexts (for example when countries or states are taken as the unit of measurement) the number of individuals  $n$  can be small while the number of time series observations  $T$  can be moderately large. An important issue in econometric panel data is the treatment of error components.

Clustered samples are increasing popular in applied economics, and is related to panel data. In clustered sampling, the observations are grouped into “clusters” which are treated as mutually independent, yet allowed to be dependent within the cluster. The major difference with panel data is that clustered sampling typically does not explicitly model error component structures, nor the dependence within clusters, but rather is concerned with inference which is robust to arbitrary forms of within-cluster correlation.

Spatial dependence is another model of interdependence. The observations are treated as mutually dependent according to a spatial measure (for example, geographic proximity). Unlike clustering, spatial models allow all observations to be mutually dependent, and typically rely on explicit modeling of the dependence relationships. Spatial dependence can also be viewed as a generalization of time series dependence.

### Data Structures

- Cross-section
- Time-series
- Panel
- Clustered
- Spatial

As we mentioned above, most of this text will be devoted to cross-sectional data under the assumption of mutually independent observations. By mutual independence we mean that the  $i^{th}$  observation  $(y_i, \mathbf{x}_i, \mathbf{z}_i)$  is independent of the  $j^{th}$  observation  $(y_j, \mathbf{x}_j, \mathbf{z}_j)$  for  $i \neq j$ . (Sometimes the label “independent” is misconstrued. It is a statement about the relationship between observations  $i$  and  $j$ , not a statement about the relationship between  $y_i$  and  $\mathbf{x}_i$  and/or  $\mathbf{z}_i$ .) In this case we say that the data are **independently distributed**.

Furthermore, if the data is randomly gathered, it is reasonable to model each observation as a draw from the same probability distribution. In this case we say that the data are **identically distributed**. If the observations are mutually independent and identically distributed, we say that the observations are **independent and identically distributed**, **iid**, or a **random sample**. For most of this text we will assume that our observations come from a random sample.

**Definition 1.5.1** *The observations  $(y_i, \mathbf{x}_i, \mathbf{z}_i)$  are a **sample** from the distribution  $F$  if they are identically distributed across  $i = 1, \dots, n$  with joint distribution  $F$ .*

**Definition 1.5.2** *The observations  $(y_i, \mathbf{x}_i, \mathbf{z}_i)$  are a **random sample** if they are mutually independent and identically distributed (**iid**) across  $i = 1, \dots, n$ .*



In the random sampling framework, we think of an individual observation  $(y_i, \mathbf{x}_i, \mathbf{z}_i)$  as a realization from a joint probability distribution  $F(y, \mathbf{x}, \mathbf{z})$  which we can call the **population**. This “population” is infinitely large. This abstraction can be a source of confusion as it does not correspond to a physical population in the real world. It is an abstraction since the distribution  $F$  is unknown, and the goal of statistical inference is to learn about features of  $F$  from the sample. The *assumption* of random sampling provides the mathematical foundation for treating economic statistics with the tools of mathematical statistics.

The random sampling framework was a major intellectual breakthrough of the late 19th century, allowing the application of mathematical statistics to the social sciences. Before this conceptual development, methods from mathematical statistics had not been applied to economic data as the latter was viewed as non-random. The random sampling framework enabled economic samples to be treated as random, a necessary precondition for the application of statistical methods.

## 1.6 Sources for Economic Data

Fortunately for economists, the internet provides a convenient forum for dissemination of economic data. Many large-scale economic datasets are available without charge from governmental agencies. An excellent starting point is the Resources for Economists Data Links, available at [rfe.org](http://rfe.org). From this site you can find almost every publically available economic data set. Some specific data sources of interest include

- Bureau of Labor Statistics
- US Census
- Current Population Survey
- Survey of Income and Program Participation
- Panel Study of Income Dynamics
- Federal Reserve System (Board of Governors and regional banks)
- National Bureau of Economic Research
- U.S. Bureau of Economic Analysis
- CompuStat
- International Financial Statistics

Another good source of data is from authors of published empirical studies. Most journals in economics require authors of published papers to make their datasets generally available. For example, in its instructions for submission, *Econometrica* states:

*Econometrica* has the policy that all empirical, experimental and simulation results must be replicable. Therefore, authors of accepted papers must submit data sets, programs, and information on empirical analysis, experiments and simulations that are needed for replication and some limited sensitivity analysis.

The *American Economic Review* states:

All data used in analysis must be made available to any researcher for purposes of replication.

The *Journal of Political Economy* states:

It is the policy of the *Journal of Political Economy* to publish papers only if the data used in the analysis are clearly and precisely documented and are readily available to any researcher for purposes of replication.

If you are interested in using the data from a published paper, first check the journal's website, as many journals archive data and replication programs online. Second, check the website(s) of the paper's author(s). Most academic economists maintain webpages, and some make available replication files complete with data and programs. If these investigations fail, email the author(s), politely requesting the data. You may need to be persistent.

As a matter of professional etiquette, all authors absolutely have the obligation to make their data and programs available. Unfortunately, many fail to do so, and typically for poor reasons. The irony of the situation is that it is typically in the best interests of a scholar to make as much of their work (including all data and programs) freely available, as this only increases the likelihood of their work being cited and having an impact.

Keep this in mind as you start your own empirical project. Remember that as part of your end product, you will need (and want) to provide all data and programs to the community of scholars. The greatest form of flattery is to learn that another scholar has read your paper, wants to extend your work, or wants to use your empirical methods. In addition, public openness provides a healthy incentive for transparency and integrity in empirical analysis.

## 1.7 Econometric Software

Economists use a variety of econometric, statistical, and programming software.

Stata ([www.stata.com](http://www.stata.com)) is a powerful statistical program with a broad set of pre-programmed econometric and statistical tools. It is quite popular among economists, and is continuously being updated with new methods. It is an excellent package for most econometric analysis, but is limited when you want to use new or less-common econometric methods which have not yet been programmed.

R ([www.r-project.org](http://www.r-project.org)), GAUSS ([www.aptech.com](http://www.aptech.com)), MATLAB ([www.mathworks.com](http://www.mathworks.com)), and Ox-Metrics ([www.oxmetrics.net](http://www.oxmetrics.net)) are high-level matrix programming languages with a wide variety of built-in statistical functions. Many econometric methods have been programmed in these languages and are available on the web. The advantage of these packages is that you are in complete control of your analysis, and it is easier to program new methods than in Stata. Some disadvantages are that you have to do much of the programming yourself, programming complicated procedures takes significant time, and programming errors are hard to prevent and difficult to detect and eliminate. Of these languages, GAUSS used to be quite popular among econometricians, but currently MATLAB is more popular. A smaller but growing group of econometricians are enthusiastic fans of R, which of these languages is uniquely open-source, user-contributed, and best of all, completely free!

For highly-intensive computational tasks, some economists write their programs in a standard programming language such as Fortran or C. This can lead to major gains in computational speed, at the cost of increased time in programming and debugging.

As these different packages have distinct advantages, many empirical economists end up using more than one package. As a student of econometrics, you will learn at least one of these packages, and probably more than one.

## 1.8 Data Files for Textbook

On the textbook webpage <http://www.ssc.wisc.edu/~bhansen/econometrics/> there are posted a number of files containing data sets which are used in this textbook both for illustration and for end-of-chapter empirical exercises. For each data sets there are four files: (1) Description (pdf format); (2) Excel data file; (3) Text data file; (4) Stata data file. The three data files are identical

in content, the observations and variables are listed in the same order in each, all have variable labels.

For example, the text makes frequent reference to a wage data set extracted from the Current Population Survey. This data set is named `cps09mar`, and is represented by the files `cps09mar_description.pdf`, `cps09mar.xlsx`, `cps09mar.txt`, and `cps09mar.dta`.

The data sets currently included are

- `cps09mar`
  - household survey data extracted from the March 2009 Current Population Survey
- `DDK2011`
  - Data file from Duflo, Dupas and Kremer (2011)
- `invest`
  - Data file from B.E. Hansen (1999), extracted from Hall and Hall (1993)
- `Nerlove1963`
  - Data file from Nerlov (1963)
- `MRW1992`
  - Data file from Mankiw, Romer and Weil (1992)
- `Card1995`
  - Data file from Card (1995)
- `AJR2001`
  - Data file from Acemoglu, Johnson and Robinson (2001)
- `AK1991`
  - Data file from Angrist and Krueger (1991)
- `hprice1`
  - Housing price data. The only files posted are `hprice1.txt` and `hprice1.pdf` which are the data in text format and description, respectively

## 1.9 Reading the Manuscript

I have endeavored to use a unified notation and nomenclature. The development of the material is cumulative, with later chapters building on the earlier ones. Nevertheless, every attempt has been made to make each chapter self-contained, so readers can pick and choose topics according to their interests.

To fully understand econometric methods, it is necessary to have a mathematical understanding of its mechanics, and this includes the mathematical proofs of the main results. Consequently, this text is self-contained, with nearly all results proved with full mathematical rigor. The mathematical development and proofs aim at brevity and conciseness (sometimes described as mathematical

elegance), but also at pedagogy. To understand a mathematical proof, it is not sufficient to simply *read* the proof, you need to follow it, and re-create it for yourself.

Nevertheless, many readers will not be interested in each mathematical detail, explanation, or proof. This is okay. To use a method it may not be necessary to understand the mathematical details. Accordingly I have placed the more technical mathematical proofs and details in chapter appendices. These appendices and other technical sections are marked with an asterisk (\*). These sections can be skipped without any loss in exposition.

## 1.10 Common Symbols

$y$	scalar
$\mathbf{x}$	vector
$\mathbf{X}$	matrix
$\mathbb{R}$	real line
$\mathbb{R}^k$	Euclidean $k$ space
$\mathbb{E}(y)$	mathematical expectation
$\text{var}(y)$	variance
$\text{cov}(x, y)$	covariance
$\text{var}(\mathbf{x})$	covariance matrix
$\text{corr}(x, y)$	correlation
$\text{Pr}$	probability
$\longrightarrow$	limit
$\xrightarrow{p}$	convergence in probability
$\xrightarrow{d}$	convergence in distribution
$\text{plim}_{n \rightarrow \infty}$	probability limit
$N(0, 1)$	standard normal distribution
$N(\mu, \sigma^2)$	normal distribution with mean $\mu$ and variance $\sigma^2$
$\chi_k^2$	chi-square distribution with $k$ degrees of freedom
$\mathbf{I}_n$	$n \times n$ identity matrix
$\text{tr } \mathbf{A}$	trace
$\mathbf{A}'$	matrix transpose
$\mathbf{A}^{-1}$	matrix inverse
$\mathbf{A} > 0$	positive definite
$\mathbf{A} \geq 0$	positive semi-definite
$\ \mathbf{a}\ $	Euclidean norm
$\ \mathbf{A}\ $	matrix (Frobinus or spectral) norm
$\approx$	approximate equality
$\stackrel{\text{def}}{=}$	definitional equality
$\sim$	is distributed as
$\log$	natural logarithm

## Chapter 2

# Conditional Expectation and Projection

### 2.1 Introduction

The most commonly applied econometric tool is least-squares estimation, also known as **regression**. As we will see, least-squares is a tool to estimate an approximate conditional mean of one variable (the **dependent variable**) given another set of variables (the **regressors**, **conditioning variables**, or **covariates**).

In this chapter we abstract from estimation, and focus on the probabilistic foundation of the conditional expectation model and its projection approximation.

### 2.2 The Distribution of Wages

Suppose that we are interested in wage rates in the United States. Since wage rates vary across workers, we cannot describe wage rates by a single number. Instead, we can describe wages using a probability distribution. Formally, we view the wage of an individual worker as a random variable *wage* with the **probability distribution**

$$F(u) = \Pr(\text{wage} \leq u).$$

When we say that a person's wage is random we mean that we do not know their wage before it is measured, and we treat observed wage rates as realizations from the distribution  $F$ . Treating unobserved wages as random variables and observed wages as realizations is a powerful mathematical abstraction which allows us to use the tools of mathematical probability.

A useful thought experiment is to imagine dialing a telephone number selected at random, and then asking the person who responds to tell us their wage rate. (Assume for simplicity that all workers have equal access to telephones, and that the person who answers your call will respond honestly.) In this thought experiment, the wage of the person you have called is a single draw from the distribution  $F$  of wages in the population. By making many such phone calls we can learn the distribution  $F$  of the entire population.

When a distribution function  $F$  is differentiable we define the probability density function

$$f(u) = \frac{d}{du}F(u).$$

The density contains the same information as the distribution function, but the density is typically easier to visually interpret.

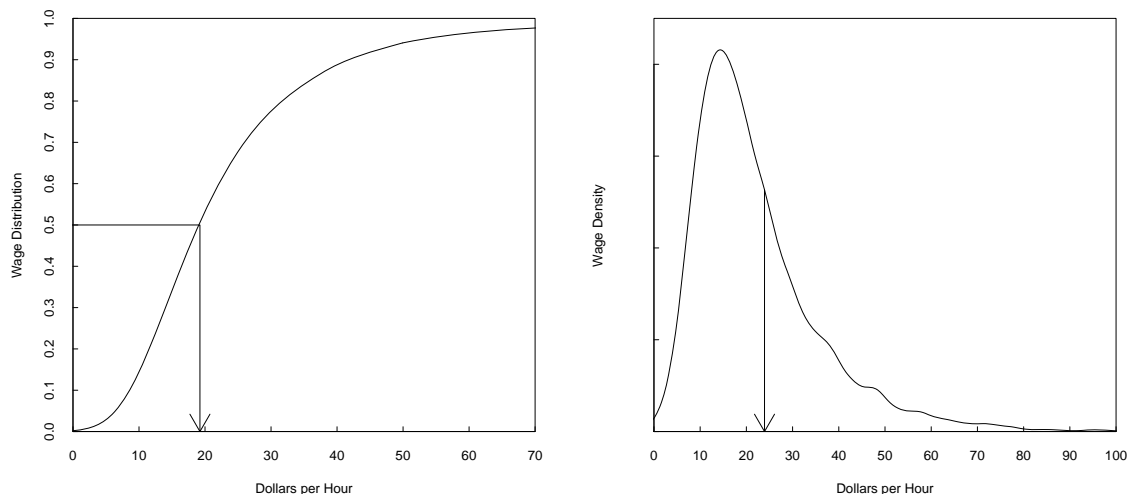


Figure 2.1: Wage Distribution and Density. All full-time U.S. workers

In Figure 2.1 we display estimates<sup>1</sup> of the probability distribution function (on the left) and density function (on the right) of U.S. wage rates in 2009. We see that the density is peaked around \$15, and most of the probability mass appears to lie between \$10 and \$40. These are ranges for typical wage rates in the U.S. population.

Important measures of central tendency are the median and the mean. The **median**  $m$  of a continuous<sup>2</sup> distribution  $F$  is the unique solution to

$$F(m) = \frac{1}{2}.$$

The median U.S. wage (\$19.23) is indicated in the left panel of Figure 2.1 by the arrow. The median is a robust<sup>3</sup> measure of central tendency, but it is tricky to use for many calculations as it is not a linear operator.

The **expectation** or **mean** of a random variable  $y$  with density  $f$  is

$$\mu = \mathbb{E}(y) = \int_{-\infty}^{\infty} uf(u)du.$$

Here we have used the common and convenient convention of using the single character  $y$  to denote a random variable, rather than the more cumbersome label *wage*. A general definition of the mean is presented in Section 2.30. The mean U.S. wage (\$23.90) is indicated in the right panel of Figure 2.1 by the arrow.

We sometimes use the notation  $\mathbb{E}y$  instead of  $\mathbb{E}(y)$  when the variable whose expectation is being taken is clear from the context. There is no distinction in meaning.

The mean is a convenient measure of central tendency because it is a linear operator and arises naturally in many economic models. A disadvantage of the mean is that it is not robust<sup>4</sup> especially in the presence of substantial skewness or thick tails, which are both features of the wage

<sup>1</sup>The distribution and density are estimated nonparametrically from the sample of 50,742 full-time non-military wage-earners reported in the March 2009 Current Population Survey. The wage rate is constructed as annual individual wage and salary earnings divided by hours worked.

<sup>2</sup>If  $F$  is not continuous the definition is  $m = \inf\{u : F(u) \geq \frac{1}{2}\}$

<sup>3</sup>The median is not sensitive to perturbations in the tails of the distribution.

<sup>4</sup>The mean is sensitive to perturbations in the tails of the distribution.

distribution as can be seen easily in the right panel of Figure 2.1. Another way of viewing this is that 64% of workers earn less than the mean wage of \$23.90, suggesting that it is incorrect to describe the mean as a “typical” wage rate.

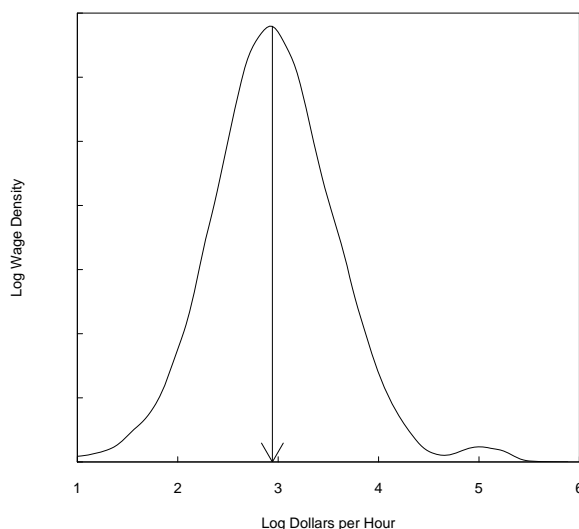


Figure 2.2: Log Wage Density

In this context it is useful to transform the data by taking the natural logarithm<sup>5</sup>. Figure 2.2 shows the density of log hourly wages  $\log(wage)$  for the same population, with its mean 2.95 drawn in with the arrow. The density of log wages is much less skewed and fat-tailed than the density of the level of wages, so its mean

$$\mathbb{E}(\log(wage)) = 2.95$$

is a much better (more robust) measure<sup>6</sup> of central tendency of the distribution. For this reason, wage regressions typically use log wages as a dependent variable rather than the level of wages.

Another useful way to summarize the probability distribution  $F(u)$  is in terms of its quantiles. For any  $\alpha \in (0, 1)$ , the  $\alpha^{th}$  quantile of the continuous<sup>7</sup> distribution  $F$  is the real number  $q_\alpha$  which satisfies

$$F(q_\alpha) = \alpha.$$

The quantile function  $q_\alpha$ , viewed as a function of  $\alpha$ , is the inverse of the distribution function  $F$ . The most commonly used quantile is the median, that is,  $q_{0.5} = m$ . We sometimes refer to quantiles by the percentile representation of  $\alpha$ , and in this case they are often called percentiles, e.g. the median is the 50<sup>th</sup> percentile.

## 2.3 Conditional Expectation

We saw in Figure 2.2 the density of log wages. Is this distribution the same for all workers, or does the wage distribution vary across subpopulations? To answer this question, we can compare wage distributions for different groups – for example, men and women. The plot on the left in Figure 2.3 displays the densities of log wages for U.S. men and women with their means (3.05 and 2.81) indicated by the arrows. We can see that the two wage densities take similar shapes but the density for men is somewhat shifted to the right with a higher mean.

<sup>5</sup>Throughout the text, we will use  $\log(y)$  or  $\log y$  to denote the natural logarithm of  $y$ .

<sup>6</sup>More precisely, the geometric mean  $\exp(\mathbb{E}(\log w)) = \$19.11$  is a robust measure of central tendency.

<sup>7</sup>If  $F$  is not continuous the definition is  $q_\alpha = \inf\{u : F(u) \geq \alpha\}$



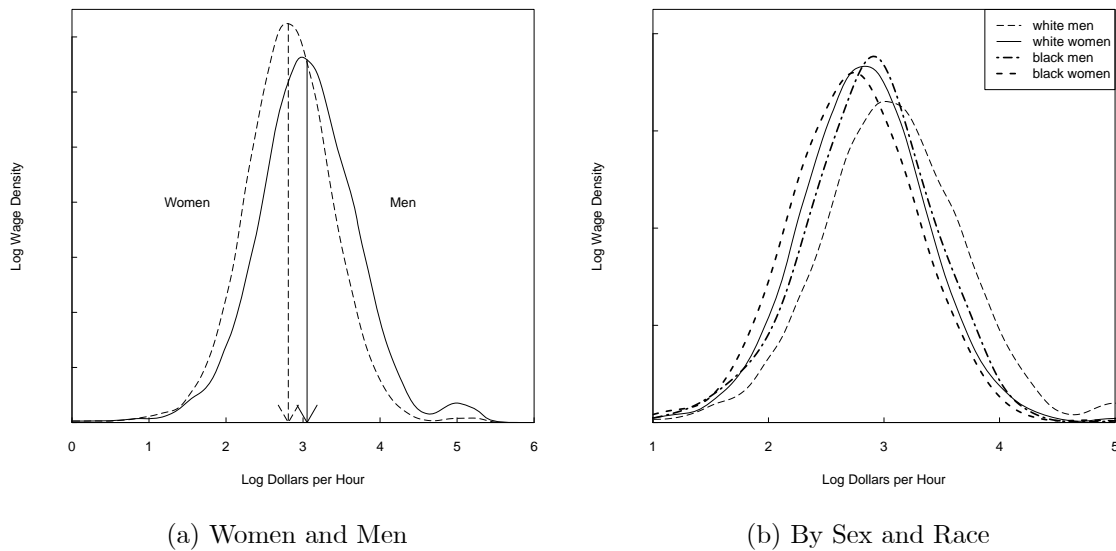


Figure 2.3: Log Wage Density by Sex and Race

The values 3.05 and 2.81 are the mean log wages in the subpopulations of men and women workers. They are called the **conditional means** (or **conditional expectations**) of log wages given sex. We can write their specific values as

$$\mathbb{E}(\log(wage) \mid sex = man) = 3.05 \quad (2.1)$$

$$\mathbb{E}(\log(wage) \mid sex = woman) = 2.81. \quad (2.2)$$

We call these means *conditional* as they are conditioning on a fixed value of the variable *sex*. While you might not think of a person's sex as a random variable, it is random from the viewpoint of econometric analysis. If you randomly select an individual, the sex of the individual is unknown and thus random. (In the population of U.S. workers, the probability that a worker is a woman happens to be 43%.) In observational data, it is most appropriate to view all measurements as random variables, and the means of subpopulations are then conditional means.

As the two densities in Figure 2.3 appear similar, a hasty inference might be that there is not a meaningful difference between the wage distributions of men and women. Before jumping to this conclusion let us examine the differences in the distributions of Figure 2.3 more carefully. As we mentioned above, the primary difference between the two densities appears to be their means. This difference equals

$$\begin{aligned} \mathbb{E}(\log(wage) \mid sex = man) - \mathbb{E}(\log(wage) \mid sex = woman) &= 3.05 - 2.81 \\ &= 0.24. \end{aligned} \quad (2.3)$$

A difference in expected log wages of 0.24 implies an average 24% difference between the wages of men and women, which is quite substantial. (For an explanation of logarithmic and percentage differences see Section 2.4.)

Consider further splitting the men and women subpopulations by race, dividing the population into whites, blacks, and other races. We display the log wage density functions of four of these groups on the right in Figure 2.3. Again we see that the primary difference between the four density functions is their central tendency.

	men	women
white	3.07	2.82
black	2.86	2.73
other	3.03	2.86

Table 2.1: Mean Log Wages by Sex and Race

Focusing on the means of these distributions, Table 2.1 reports the mean log wage for each of the six sub-populations.

The entries in Table 2.1 are the conditional means of  $\log(wage)$  given *sex* and *race*. For example

$$\mathbb{E}(\log(wage) \mid sex = man, race = white) = 3.07$$

and

$$\mathbb{E}(\log(wage) \mid sex = woman, race = black) = 2.73.$$

One benefit of focusing on conditional means is that they reduce complicated distributions to a single summary measure, and thereby facilitate comparisons across groups. Because of this simplifying property, conditional means are the primary interest of regression analysis and are a major focus in econometrics.

Table 2.1 allows us to easily calculate average wage differences between groups. For example, we can see that the wage gap between men and women continues after disaggregation by race, as the average gap between white men and white women is 25%, and that between black men and black women is 13%. We also can see that there is a race gap, as the average wages of blacks are substantially less than the other race categories. In particular, the average wage gap between white men and black men is 21%, and that between white women and black women is 9%.

## 2.4 Log Differences\*

A useful approximation for the natural logarithm for small  $x$  is

$$\log(1+x) \approx x. \quad (2.4)$$

This can be derived from the infinite series expansion of  $\log(1+x)$ :

$$\begin{aligned} \log(1+x) &= x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \cdots \\ &= x + O(x^2). \end{aligned}$$

The symbol  $O(x^2)$  means that the remainder is bounded by  $Ax^2$  as  $x \rightarrow 0$  for some  $A < \infty$ . A plot of  $\log(1+x)$  and the linear approximation  $x$  is shown in Figure 2.4. We can see that  $\log(1+x)$  and the linear approximation  $x$  are very close for  $|x| \leq 0.1$ , and reasonably close for  $|x| \leq 0.2$ , but the difference increases with  $|x|$ .

Now, if  $y^*$  is  $c\%$  greater than  $y$ , then

$$y^* = (1 + c/100)y.$$

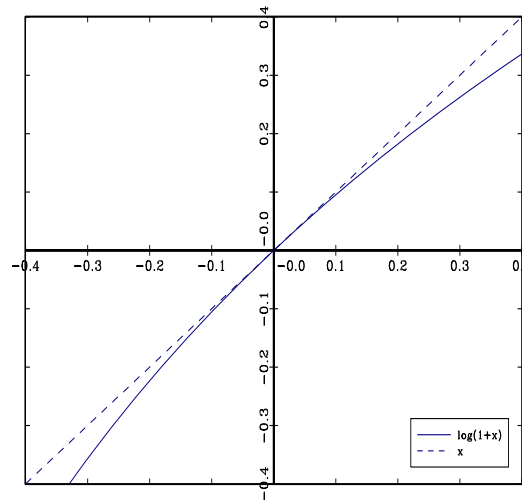
Taking natural logarithms,

$$\log y^* = \log y + \log(1 + c/100)$$

or

$$\log y^* - \log y = \log(1 + c/100) \approx \frac{c}{100}$$

where the approximation is (2.4). This shows that 100 multiplied by the difference in logarithms is approximately the percentage difference between  $y$  and  $y^*$ , and this approximation is quite good for  $|c| \leq 10$ .

Figure 2.4:  $\log(1+x)$ 

## 2.5 Conditional Expectation Function

An important determinant of wage levels is education. In many empirical studies economists measure educational attainment by the number of years<sup>8</sup> of schooling, and we will write this variable as *education*.

The conditional mean of log wages given *sex*, *race*, and *education* is a single number for each category. For example

$$\mathbb{E}(\log(\text{wage}) \mid \text{sex} = \text{man}, \text{race} = \text{white}, \text{education} = 12) = 2.84.$$

We display in Figure 2.5 the conditional means of  $\log(\text{wage})$  for white men and white women as a function of *education*. The plot is quite revealing. We see that the conditional mean is increasing in years of education, but at a different rate for schooling levels above and below nine years. Another striking feature of Figure 2.5 is that the gap between men and women is roughly constant for all education levels. As the variables are measured in logs this implies a constant average percentage gap between men and women regardless of educational attainment.

In many cases it is convenient to simplify the notation by writing variables using single characters, typically  $y$ ,  $x$  and/or  $z$ . It is conventional in econometrics to denote the dependent variable (e.g.  $\log(\text{wage})$ ) by the letter  $y$ , a conditioning variable (such as *sex*) by the letter  $x$ , and multiple conditioning variables (such as *race*, *education* and *sex*) by the subscripted letters  $x_1, x_2, \dots, x_k$ .

Conditional expectations can be written with the generic notation

$$\mathbb{E}(y \mid x_1, x_2, \dots, x_k) = m(x_1, x_2, \dots, x_k).$$

We call this the **conditional expectation function** (CEF). The CEF is a function of  $(x_1, x_2, \dots, x_k)$  as it varies with the variables. For example, the conditional expectation of  $y = \log(\text{wage})$  given  $(x_1, x_2) = (\text{sex}, \text{race})$  is given by the six entries of Table 2.1. The CEF is a function of  $(\text{sex}, \text{race})$  as it varies across the entries.

For greater compactness, we will typically write the conditioning variables as a vector in  $\mathbb{R}^k$  :

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{pmatrix}. \quad (2.5)$$

---

<sup>8</sup>Here, *education* is defined as years of schooling beyond kindergarten. A high school graduate has *education*=12, a college graduate has *education*=16, a Master's degree has *education*=18, and a professional degree (medical, law or PhD) has *education*=20.

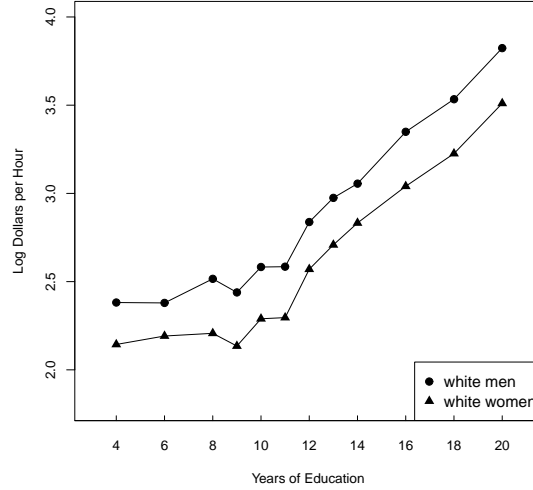


Figure 2.5: Mean Log Wage as a Function of Years of Education

Here we follow the convention of using lower case bold italics  $\mathbf{x}$  to denote a vector. Given this notation, the CEF can be compactly written as

$$\mathbb{E}(y \mid \mathbf{x}) = m(\mathbf{x}).$$

The CEF  $\mathbb{E}(y \mid \mathbf{x})$  is a random variable as it is a function of the random variable  $\mathbf{x}$ . It is also sometimes useful to view the CEF as a function of  $\mathbf{x}$ . In this case we can write  $m(\mathbf{u}) = \mathbb{E}(y \mid \mathbf{x} = \mathbf{u})$ , which is a function of the argument  $\mathbf{u}$ . The expression  $\mathbb{E}(y \mid \mathbf{x} = \mathbf{u})$  is the conditional expectation of  $y$ , given that we know that the random variable  $\mathbf{x}$  equals the specific value  $\mathbf{u}$ . However, sometimes in econometrics we take a notational shortcut and use  $\mathbb{E}(y \mid \mathbf{x})$  to refer to this function. Hopefully, the use of  $\mathbb{E}(y \mid \mathbf{x})$  should be apparent from the context.

## 2.6 Continuous Variables

In the previous sections, we implicitly assumed that the conditioning variables are discrete. However, many conditioning variables are continuous. In this section, we take up this case and assume that the variables  $(y, \mathbf{x})$  are continuously distributed with a joint density function  $f(y, \mathbf{x})$ .

As an example, take  $y = \log(\text{wage})$  and  $x = \text{experience}$ , the number of years of potential labor market experience<sup>9</sup>. The contours of their joint density are plotted on the left side of Figure 2.6 for the population of white men with 12 years of education.

Given the joint density  $f(y, \mathbf{x})$  the variable  $\mathbf{x}$  has the marginal density

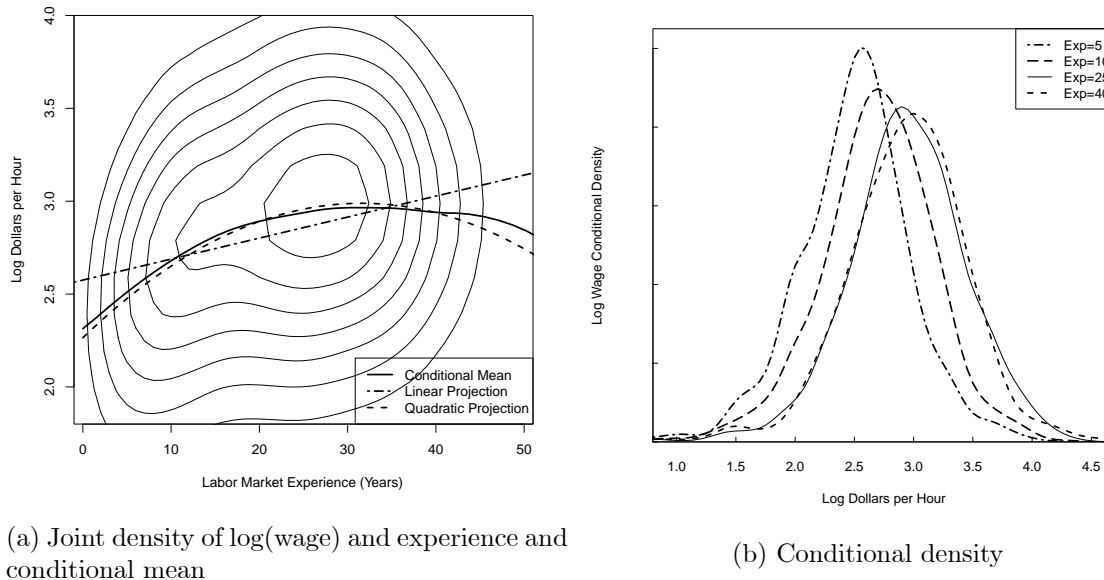
$$f_{\mathbf{x}}(\mathbf{x}) = \int_{-\infty}^{\infty} f(y, \mathbf{x}) dy.$$

For any  $\mathbf{x}$  such that  $f_{\mathbf{x}}(\mathbf{x}) > 0$  the conditional density of  $y$  given  $\mathbf{x}$  is defined as

$$f_{y|\mathbf{x}}(y \mid \mathbf{x}) = \frac{f(y, \mathbf{x})}{f_{\mathbf{x}}(\mathbf{x})}. \quad (2.6)$$

The conditional density is a (renormalized) slice of the joint density  $f(y, \mathbf{x})$  holding  $\mathbf{x}$  fixed. The slice is renormalized (divided by  $f_{\mathbf{x}}(\mathbf{x})$  so that it integrates to one and is thus a density.) We can

<sup>9</sup>Here, *experience* is defined as potential labor market experience, equal to *age* – *education* – 6

Figure 2.6: White men with *education*=12

visualize this by slicing the joint density function at a specific value of  $\mathbf{x}$  parallel with the  $y$ -axis. For example, take the density contours on the left side of Figure 2.6 and slice through the contour plot at a specific value of *experience*, and then renormalize the slice so that it is a proper density. This gives us the conditional density of  $\log(\text{wage})$  for white men with 12 years of education and this level of experience. We do this for four levels of *experience* (5, 10, 25, and 40 years), and plot these densities on the right side of Figure 2.6. We can see that the distribution of wages shifts to the right and becomes more diffuse as experience increases from 5 to 10 years, and from 10 to 25 years, but there is little change from 25 to 40 years experience.

The CEF of  $y$  given  $\mathbf{x}$  is the mean of the conditional density (2.6)

$$m(\mathbf{x}) = \mathbb{E}(y \mid \mathbf{x}) = \int_{-\infty}^{\infty} y f_{y|\mathbf{x}}(y \mid \mathbf{x}) dy. \quad (2.7)$$

Intuitively,  $m(\mathbf{x})$  is the mean of  $y$  for the idealized subpopulation where the conditioning variables are fixed at  $\mathbf{x}$ . This is idealized since  $\mathbf{x}$  is continuously distributed so this subpopulation is infinitely small.

This definition (2.7) is appropriate when the conditional density (2.6) is well defined. However, the conditional mean  $m(\mathbf{x})$  exists quite generally. In Theorem 2.32.1 in Section 2.32 we show that  $m(\mathbf{x})$  exists so long as  $\mathbb{E}|y| < \infty$ .

In Figure 2.6 the CEF of  $\log(\text{wage})$  given *experience* is plotted as the solid line. We can see that the CEF is a smooth but nonlinear function. The CEF is initially increasing in *experience*, flattens out around *experience* = 30, and then decreases for high levels of experience.

## 2.7 Law of Iterated Expectations

An extremely useful tool from probability theory is the **law of iterated expectations**. An important special case is the known as the Simple Law.

**Theorem 2.7.1 Simple Law of Iterated Expectations**

If  $\mathbb{E}|y| < \infty$  then for any random vector  $\mathbf{x}$ ,

$$\mathbb{E}(\mathbb{E}(y | \mathbf{x})) = \mathbb{E}(y)$$

The simple law states that the expectation of the conditional expectation is the unconditional expectation. In other words, the average of the conditional averages is the unconditional average. When  $\mathbf{x}$  is discrete

$$\mathbb{E}(\mathbb{E}(y | \mathbf{x})) = \sum_{j=1}^{\infty} \mathbb{E}(y | \mathbf{x}_j) \Pr(\mathbf{x} = \mathbf{x}_j)$$

and when  $\mathbf{x}$  is continuous

$$\mathbb{E}(\mathbb{E}(y | \mathbf{x})) = \int_{\mathbb{R}^k} \mathbb{E}(y | \mathbf{x}) f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x}.$$

Going back to our investigation of average log wages for men and women, the simple law states that

$$\begin{aligned} & \mathbb{E}(\log(wage) | sex = man) \Pr(sex = man) \\ & + \mathbb{E}(\log(wage) | sex = woman) \Pr(sex = woman) \\ & = \mathbb{E}(\log(wage)). \end{aligned}$$

Or numerically,

$$3.05 \times 0.57 + 2.79 \times 0.43 = 2.92.$$

The general law of iterated expectations allows two sets of conditioning variables.

**Theorem 2.7.2 Law of Iterated Expectations**

If  $\mathbb{E}|y| < \infty$  then for any random vectors  $\mathbf{x}_1$  and  $\mathbf{x}_2$ ,

$$\mathbb{E}(\mathbb{E}(y | \mathbf{x}_1, \mathbf{x}_2) | \mathbf{x}_1) = \mathbb{E}(y | \mathbf{x}_1)$$

Notice the way the law is applied. The inner expectation conditions on  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , while the outer expectation conditions only on  $\mathbf{x}_1$ . The iterated expectation yields the simple answer  $\mathbb{E}(y | \mathbf{x}_1)$ , the expectation conditional on  $\mathbf{x}_1$  alone. Sometimes we phrase this as: “The smaller information set wins.”

As an example

$$\begin{aligned} & \mathbb{E}(\log(wage) | sex = man, race = white) \Pr(race = white | sex = man) \\ & + \mathbb{E}(\log(wage) | sex = man, race = black) \Pr(race = black | sex = man) \\ & + \mathbb{E}(\log(wage) | sex = man, race = other) \Pr(race = other | sex = man) \\ & = \mathbb{E}(\log(wage) | sex = man) \end{aligned}$$

or numerically

$$3.07 \times 0.84 + 2.86 \times 0.08 + 3.03 \times 0.08 = 3.05.$$

A property of conditional expectations is that when you condition on a random vector  $\mathbf{x}$  you can effectively treat it as if it is constant. For example,  $\mathbb{E}(\mathbf{x} | \mathbf{x}) = \mathbf{x}$  and  $\mathbb{E}(g(\mathbf{x}) | \mathbf{x}) = g(\mathbf{x})$  for any function  $g(\cdot)$ . The general property is known as the Conditioning Theorem.

**Theorem 2.7.3 *Conditioning Theorem****If  $\mathbb{E}|y| < \infty$  then*

$$\mathbb{E}(g(\mathbf{x})y \mid \mathbf{x}) = g(\mathbf{x})\mathbb{E}(y \mid \mathbf{x}). \quad (2.8)$$

*In in addition*

$$\mathbb{E}|g(\mathbf{x})y| < \infty \quad (2.9)$$

*then*

$$\mathbb{E}(g(\mathbf{x})y) = \mathbb{E}(g(\mathbf{x})\mathbb{E}(y \mid \mathbf{x})). \quad (2.10)$$

The proofs of Theorems 2.7.1, 2.7.2 and 2.7.3 are given in Section 2.34.

## 2.8 CEF Error

The CEF error  $e$  is defined as the difference between  $y$  and the CEF evaluated at the random vector  $\mathbf{x}$ :

$$e = y - m(\mathbf{x}).$$

By construction, this yields the formula

$$y = m(\mathbf{x}) + e. \quad (2.11)$$

In (2.11) it is useful to understand that the error  $e$  is derived from the joint distribution of  $(y, \mathbf{x})$ , and so its properties are derived from this construction.

**A key property of the CEF error is that it has a conditional mean of zero.** To see this, by the linearity of expectations, the definition  $m(\mathbf{x}) = \mathbb{E}(y \mid \mathbf{x})$  and the Conditioning Theorem

$$\begin{aligned} \mathbb{E}(e \mid \mathbf{x}) &= \mathbb{E}((y - m(\mathbf{x})) \mid \mathbf{x}) \\ &= \mathbb{E}(y \mid \mathbf{x}) - \mathbb{E}(m(\mathbf{x}) \mid \mathbf{x}) \\ &= m(\mathbf{x}) - m(\mathbf{x}) \\ &= 0. \end{aligned}$$

This fact can be combined with the law of iterated expectations to show that the unconditional mean is also zero.

$$\mathbb{E}(e) = \mathbb{E}(\mathbb{E}(e \mid \mathbf{x})) = \mathbb{E}(0) = 0.$$

We state this and some other results formally.

**Theorem 2.8.1 *Properties of the CEF error****If  $\mathbb{E}|y| < \infty$  then*

1.  $\mathbb{E}(e \mid \mathbf{x}) = 0$ .
2.  $\mathbb{E}(e) = 0$ .
3. If  $\mathbb{E}|y|^r < \infty$  for  $r \geq 1$  then  $\mathbb{E}|e|^r < \infty$ .
4. For any function  $h(\mathbf{x})$  such that  $\mathbb{E}|h(\mathbf{x})e| < \infty$  then  $\mathbb{E}(h(\mathbf{x})e) = 0$ .

The proof of the third result is deferred to Section 2.34.

The fourth result, whose proof is left to Exercise 2.3, implies that  $e$  is uncorrelated with any function of the regressors.

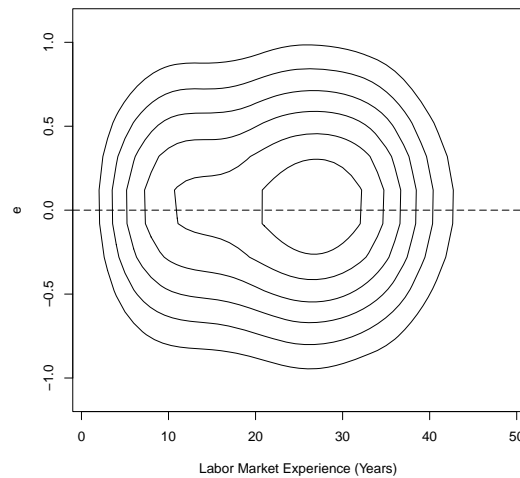


Figure 2.7: Joint density of CEF error  $e$  and *experience* for white men with *education*=12.

The equations

$$\begin{aligned} y &= m(\mathbf{x}) + e \\ \mathbb{E}(e \mid \mathbf{x}) &= 0 \end{aligned}$$

together imply that  $m(\mathbf{x})$  is the CEF of  $y$  given  $\mathbf{x}$ . It is important to understand that this is not a restriction. These equations hold true by definition.

The condition  $\mathbb{E}(e \mid \mathbf{x}) = 0$  is implied by the definition of  $e$  as the difference between  $y$  and the CEF  $m(\mathbf{x})$ . The equation  $\mathbb{E}(e \mid \mathbf{x}) = 0$  is sometimes called a conditional mean restriction, since the conditional mean of the error  $e$  is restricted to equal zero. The property is also sometimes called **mean independence**, for the conditional mean of  $e$  is 0 and thus independent of  $\mathbf{x}$ . However, it does not imply that the distribution of  $e$  is independent of  $\mathbf{x}$ . Sometimes the assumption “ $e$  is independent of  $\mathbf{x}$ ” is added as a convenient simplification, but it is not generic feature of the conditional mean. Typically and generally,  $e$  and  $\mathbf{x}$  are jointly dependent, even though the conditional mean of  $e$  is zero.

As an example, the contours of the joint density of  $e$  and *experience* are plotted in Figure 2.7 for the same population as Figure 2.6. The error  $e$  has a conditional mean of zero for all values of *experience*, but the shape of the conditional distribution varies with the level of *experience*.

As a simple example of a case where  $x$  and  $e$  are mean independent yet dependent, let  $e = x\varepsilon$  where  $x$  and  $\varepsilon$  are independent  $N(0, 1)$ . Then conditional on  $x$ , the error  $e$  has the distribution  $N(0, x^2)$ . Thus  $\mathbb{E}(e \mid x) = 0$  and  $e$  is mean independent of  $x$ , yet  $e$  is not fully independent of  $x$ . Mean independence does not imply full independence.

## 2.9 Intercept-Only Model

A special case of the regression model is when there are no regressors  $\mathbf{x}$ . In this case  $m(\mathbf{x}) = \mathbb{E}(y) = \mu$ , the unconditional mean of  $y$ . We can still write an equation for  $y$  in the regression format:

$$\begin{aligned} y &= \mu + e \\ \mathbb{E}(e) &= 0. \end{aligned}$$

This is useful for it unifies the notation.



## 2.10 Regression Variance

An important measure of the dispersion about the CEF function is the unconditional variance of the CEF error  $e$ . We write this as

$$\sigma^2 = \text{var}(e) = \mathbb{E}\left((e - \mathbb{E}e)^2\right) = \mathbb{E}(e^2).$$

Theorem 2.8.1.3 implies the following simple but useful result.

**Theorem 2.10.1** *If  $\mathbb{E}(y^2) < \infty$  then  $\sigma^2 < \infty$ .*

We can call  $\sigma^2$  the regression variance or the variance of the regression error. The magnitude of  $\sigma^2$  measures the amount of variation in  $y$  which is not “explained” or accounted for in the conditional mean  $\mathbb{E}(y | \mathbf{x})$ .

The regression variance depends on the regressors  $\mathbf{x}$ . Consider two regressions

$$\begin{aligned} y &= \mathbb{E}(y | \mathbf{x}_1) + e_1 \\ y &= \mathbb{E}(y | \mathbf{x}_1, \mathbf{x}_2) + e_2. \end{aligned}$$

We write the two errors distinctly as  $e_1$  and  $e_2$  as they are different – changing the conditioning information changes the conditional mean and therefore the regression error as well.

In our discussion of iterated expectations, we have seen that by increasing the conditioning set, the conditional expectation reveals greater detail about the distribution of  $y$ . What is the implication for the regression error?

It turns out that there is a simple relationship. We can think of the conditional mean  $\mathbb{E}(y | \mathbf{x})$  as the “explained portion” of  $y$ . The remainder  $e = y - \mathbb{E}(y | \mathbf{x})$  is the “unexplained portion”. The simple relationship we now derive shows that the variance of this unexplained portion decreases when we condition on more variables. This relationship is monotonic in the sense that increasing the amount of information always decreases the variance of the unexplained portion.

**Theorem 2.10.2** *If  $\mathbb{E}(y^2) < \infty$  then*

$$\text{var}(y) \geq \text{var}(y - \mathbb{E}(y | \mathbf{x}_1)) \geq \text{var}(y - \mathbb{E}(y | \mathbf{x}_1, \mathbf{x}_2)).$$

Theorem 2.10.2 says that the variance of the difference between  $y$  and its conditional mean (weakly) decreases whenever an additional variable is added to the conditioning information.

The proof of Theorem 2.10.2 is given in Section 2.34.

## 2.11 Best Predictor

Suppose that given a realized value of  $\mathbf{x}$ , we want to create a prediction or forecast of  $y$ . We can write any predictor as a function  $g(\mathbf{x})$  of  $\mathbf{x}$ . The prediction error is the realized difference  $y - g(\mathbf{x})$ . A non-stochastic measure of the magnitude of the prediction error is the expectation of its square

$$\mathbb{E}\left((y - g(\mathbf{x}))^2\right). \tag{2.12}$$

We can define the best predictor as the function  $g(\mathbf{x})$  which minimizes (2.12). What function is the best predictor? It turns out that the answer is the CEF  $m(\mathbf{x})$ . This holds regardless of the joint distribution of  $(y, \mathbf{x})$ .

To see this, note that the mean squared error of a predictor  $g(\mathbf{x})$  is

$$\begin{aligned}
 \mathbb{E} \left( (y - g(\mathbf{x}))^2 \right) &= \mathbb{E} \left( (e + m(\mathbf{x}) - g(\mathbf{x}))^2 \right) \\
 &= \mathbb{E}(e^2) + 2\mathbb{E}(e(m(\mathbf{x}) - g(\mathbf{x}))) + \mathbb{E} \left( (m(\mathbf{x}) - g(\mathbf{x}))^2 \right) \\
 &= \mathbb{E}(e^2) + \mathbb{E} \left( (m(\mathbf{x}) - g(\mathbf{x}))^2 \right) \\
 &\geq \mathbb{E}(e^2) \\
 &= \mathbb{E} \left( (y - m(\mathbf{x}))^2 \right)
 \end{aligned}$$

where the first equality makes the substitution  $y = m(\mathbf{x}) + e$  and the third equality uses Theorem 2.8.1.4. The right-hand-side after the third equality is minimized by setting  $g(\mathbf{x}) = m(\mathbf{x})$ , yielding the inequality in the fourth line. The minimum is finite under the assumption  $\mathbb{E}(y^2) < \infty$  as shown by Theorem 2.10.1.

We state this formally in the following result.

**Theorem 2.11.1** *Conditional Mean as Best Predictor*

If  $\mathbb{E}(y^2) < \infty$ , then for any predictor  $g(\mathbf{x})$ ,

$$\mathbb{E} \left( (y - g(\mathbf{x}))^2 \right) \geq \mathbb{E} \left( (y - m(\mathbf{x}))^2 \right)$$

where  $m(\mathbf{x}) = \mathbb{E}(y \mid \mathbf{x})$ .

It may be helpful to consider this result in the context of the intercept-only model

$$\begin{aligned}
 y &= \mu + e \\
 \mathbb{E}(e) &= 0.
 \end{aligned}$$

Theorem 2.11.1 shows that the best predictor for  $y$  (in the class of constants) is the unconditional mean  $\mu = \mathbb{E}(y)$ , in the sense that the mean minimizes the mean squared prediction error.

## 2.12 Conditional Variance

While the conditional mean is a good measure of the location of a conditional distribution, it does not provide information about the spread of the distribution. A common measure of the dispersion is the **conditional variance**. We first give the general definition of the conditional variance of a random variable  $w$ .

**Definition 2.12.1** If  $\mathbb{E}(w^2) < \infty$ , the **conditional variance** of  $w$  given  $\mathbf{x}$  is

$$\text{var}(w \mid \mathbf{x}) = \mathbb{E} \left( (w - \mathbb{E}(w \mid \mathbf{x}))^2 \mid \mathbf{x} \right)$$

Notice that the conditional variance is the conditional second moment, centered around the conditional first moment. Given this definition, we define the conditional variance of the regression error.

**Definition 2.12.2** If  $\mathbb{E}(e^2) < \infty$ , the **conditional variance** of the regression error  $e$  is

$$\sigma^2(\mathbf{x}) = \text{var}(e \mid \mathbf{x}) = \mathbb{E}(e^2 \mid \mathbf{x}).$$

Generally,  $\sigma^2(\mathbf{x})$  is a non-trivial function of  $\mathbf{x}$  and can take any form subject to the restriction that it is non-negative. One way to think about  $\sigma^2(\mathbf{x})$  is that it is the conditional mean of  $e^2$  given  $\mathbf{x}$ . Notice as well that  $\sigma^2(\mathbf{x}) = \text{var}(y \mid \mathbf{x})$  so it is equivalently the conditional variance of the dependent variable.

The variance is in a different unit of measurement than the original variable. To convert the variance back to the same unit of measure we define the **conditional standard deviation** as its square root  $\sigma(\mathbf{x}) = \sqrt{\sigma^2(\mathbf{x})}$ .

As an example of how the conditional variance depends on observables, compare the conditional log wage densities for men and women displayed in Figure 2.3. The difference between the densities is not purely a location shift, but is also a difference in spread. Specifically, we can see that the density for men's log wages is somewhat more spread out than that for women, while the density for women's wages is somewhat more peaked. Indeed, the conditional standard deviation for men's wages is 3.05 and that for women is 2.81. So while men have higher average wages, they are also somewhat more dispersed.

The unconditional error variance and the conditional variance are related by the law of iterated expectations

$$\sigma^2 = \mathbb{E}(e^2) = \mathbb{E}(\mathbb{E}(e^2 \mid \mathbf{x})) = \mathbb{E}(\sigma^2(\mathbf{x})).$$

That is, the unconditional error variance is the average conditional variance.

Given the conditional variance, we can define a rescaled error

$$\varepsilon = \frac{e}{\sigma(\mathbf{x})}. \quad (2.13)$$

We can calculate that since  $\sigma(\mathbf{x})$  is a function of  $\mathbf{x}$

$$\mathbb{E}(\varepsilon \mid \mathbf{x}) = \mathbb{E}\left(\frac{e}{\sigma(\mathbf{x})} \mid \mathbf{x}\right) = \frac{1}{\sigma(\mathbf{x})} \mathbb{E}(e \mid \mathbf{x}) = 0$$

and

$$\text{var}(\varepsilon \mid \mathbf{x}) = \mathbb{E}(\varepsilon^2 \mid \mathbf{x}) = \mathbb{E}\left(\frac{e^2}{\sigma^2(\mathbf{x})} \mid \mathbf{x}\right) = \frac{1}{\sigma^2(\mathbf{x})} \mathbb{E}(e^2 \mid \mathbf{x}) = \frac{\sigma^2(\mathbf{x})}{\sigma^2(\mathbf{x})} = 1.$$

Thus  $\varepsilon$  has a conditional mean of zero, and a conditional variance of 1.

Notice that (2.13) can be rewritten as

$$e = \sigma(\mathbf{x})\varepsilon.$$

and substituting this for  $e$  in the CEF equation (2.11), we find that

$$y = m(\mathbf{x}) + \sigma(\mathbf{x})\varepsilon. \quad (2.14)$$

This is an alternative (mean-variance) representation of the CEF equation.

Many econometric studies focus on the conditional mean  $m(\mathbf{x})$  and either ignore the conditional variance  $\sigma^2(\mathbf{x})$ , treat it as a constant  $\sigma^2(\mathbf{x}) = \sigma^2$ , or treat it as a nuisance parameter (a parameter not of primary interest). This is appropriate when the primary variation in the conditional distribution is in the mean, but can be short-sighted in other cases. Dispersion is relevant to many economic topics, including income and wealth distribution, economic inequality, and price dispersion. Conditional dispersion (variance) can be a fruitful subject for investigation.

The perverse consequences of a narrow-minded focus on the mean has been parodied in a classic joke:

An economist was standing with one foot in a bucket of boiling water and the other foot in a bucket of ice. When asked how he felt, he replied, “On average I feel just fine.”

Clearly, the economist in question ignored variance!

## 2.13 Homoskedasticity and Heteroskedasticity

An important special case obtains when the conditional variance  $\sigma^2(\mathbf{x})$  is a constant and independent of  $\mathbf{x}$ . This is called **homoskedasticity**.

**Definition 2.13.1** *The error is **homoskedastic** if  $\mathbb{E}(e^2 | \mathbf{x}) = \sigma^2$  does not depend on  $\mathbf{x}$ .*

In the general case where  $\sigma^2(\mathbf{x})$  depends on  $\mathbf{x}$  we say that the error  $e$  is **heteroskedastic**.

**Definition 2.13.2** *The error is **heteroskedastic** if  $\mathbb{E}(e^2 | \mathbf{x}) = \sigma^2(\mathbf{x})$  depends on  $\mathbf{x}$ .*

It is helpful to understand that the concepts homoskedasticity and heteroskedasticity concern the conditional variance, not the unconditional variance. By definition, the unconditional variance  $\sigma^2$  is a constant and independent of the regressors  $\mathbf{x}$ . So when we talk about the variance as a function of the regressors, we are talking about the conditional variance  $\sigma^2(\mathbf{x})$ .

Some older or introductory textbooks describe heteroskedasticity as the case where “the variance of  $e$  varies across observations”. This is a poor and confusing definition. It is more constructive to understand that heteroskedasticity means that the conditional variance  $\sigma^2(\mathbf{x})$  depends on observables.

Older textbooks also tend to describe homoskedasticity as a component of a correct regression specification, and describe heteroskedasticity as an exception or deviance. This description has influenced many generations of economists, but it is unfortunately backwards. The correct view is that heteroskedasticity is generic and “standard”, while homoskedasticity is unusual and exceptional. The default in empirical work should be to assume that the errors are heteroskedastic, not the converse.

In apparent contradiction to the above statement, we will still frequently impose the homoskedasticity assumption when making theoretical investigations into the properties of estimation and inference methods. The reason is that in many cases homoskedasticity greatly simplifies the theoretical calculations, and it is therefore quite advantageous for teaching and learning. It should always be remembered, however, that homoskedasticity is never imposed because it is believed to be a correct feature of an empirical model, but rather because of its simplicity.

## 2.14 Regression Derivative

One way to interpret the CEF  $m(\mathbf{x}) = \mathbb{E}(y \mid \mathbf{x})$  is in terms of how marginal changes in the regressors  $\mathbf{x}$  imply changes in the conditional mean of the response variable  $y$ . It is typical to consider marginal changes in a single regressor, say  $x_1$ , holding the remainder fixed. When a regressor  $x_1$  is continuously distributed, we define the marginal effect of a change in  $x_1$ , holding the variables  $x_2, \dots, x_k$  fixed, as the partial derivative of the CEF

$$\frac{\partial}{\partial x_1} m(x_1, \dots, x_k).$$

When  $x_1$  is discrete we define the marginal effect as a discrete difference. For example, if  $x_1$  is binary, then the marginal effect of  $x_1$  on the CEF is

$$m(1, x_2, \dots, x_k) - m(0, x_2, \dots, x_k).$$

We can unify the continuous and discrete cases with the notation

$$\nabla_1 m(\mathbf{x}) = \begin{cases} \frac{\partial}{\partial x_1} m(x_1, \dots, x_k), & \text{if } x_1 \text{ is continuous} \\ m(1, x_2, \dots, x_k) - m(0, x_2, \dots, x_k), & \text{if } x_1 \text{ is binary.} \end{cases}$$

Collecting the  $k$  effects into one  $k \times 1$  vector, we define the **regression derivative** with respect to  $\mathbf{x}$ :

$$\nabla m(\mathbf{x}) = \begin{bmatrix} \nabla_1 m(\mathbf{x}) \\ \nabla_2 m(\mathbf{x}) \\ \vdots \\ \nabla_k m(\mathbf{x}) \end{bmatrix}.$$

When all elements of  $\mathbf{x}$  are continuous, then we have the simplification  $\nabla m(\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} m(\mathbf{x})$ , the vector of partial derivatives.

There are two important points to remember concerning our definition of the regression derivative.

First, the effect of each variable is calculated holding the other variables constant. This is the **ceteris paribus** concept commonly used in economics. But in the case of a regression derivative, the conditional mean does not literally hold *all else* constant. It only holds constant the variables included in the conditional mean. This means that the regression derivative depends on which regressors are included. For example, in a regression of wages on education, experience, race and sex, the regression derivative with respect to education shows the marginal effect of education on mean wages, holding constant experience, race and sex. But it does not hold constant an individual's unobservable characteristics (such as ability), nor variables not included in the regression (such as the quality of education).

Second, the regression derivative is the change in the conditional expectation of  $y$ , not the change in the actual value of  $y$  for an individual. It is tempting to think of the regression derivative as the change in the actual value of  $y$ , but this is not a correct interpretation. The regression derivative  $\nabla m(\mathbf{x})$  is the change in the actual value of  $y$  only if the error  $e$  is unaffected by the change in the regressor  $\mathbf{x}$ . We return to a discussion of causal effects in Section 2.29.

## 2.15 Linear CEF

An important special case is when the CEF  $m(\mathbf{x}) = \mathbb{E}(y \mid \mathbf{x})$  is linear in  $\mathbf{x}$ . In this case we can write the mean equation as

$$m(\mathbf{x}) = x_1\beta_1 + x_2\beta_2 + \dots + x_k\beta_k + \beta_{k+1}.$$

Notationally it is convenient to write this as a simple function of the vector  $\mathbf{x}$ . An easy way to do so is to augment the regressor vector  $\mathbf{x}$  by listing the number “1” as an element. We call this the “constant” and the corresponding coefficient is called the “intercept”. Equivalently, specify that the final element<sup>10</sup> of the vector  $\mathbf{x}$  is  $x_k = 1$ . Thus (2.5) has been redefined as the  $k \times 1$  vector

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_{k-1} \\ 1 \end{pmatrix}. \quad (2.15)$$

With this redefinition, the CEF is

$$\begin{aligned} m(\mathbf{x}) &= x_1\beta_1 + x_2\beta_2 + \cdots + \beta_k \\ &= \mathbf{x}'\boldsymbol{\beta} \end{aligned} \quad (2.16)$$

where

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} \quad (2.17)$$

is a  $k \times 1$  coefficient vector. This is the **linear CEF model**. It is also often called the **linear regression model**, or the regression of  $y$  on  $\mathbf{x}$ .

In the linear CEF model, the regression derivative is simply the coefficient vector. That is

$$\nabla m(\mathbf{x}) = \boldsymbol{\beta}.$$

This is one of the appealing features of the linear CEF model. The coefficients have simple and natural interpretations as the marginal effects of changing one variable, holding the others constant.

#### Linear CEF Model

$$\begin{aligned} y &= \mathbf{x}'\boldsymbol{\beta} + e \\ \mathbb{E}(e \mid \mathbf{x}) &= 0 \end{aligned}$$

If in addition the error is homoskedastic, we call this the homoskedastic linear CEF model.

#### Homoskedastic Linear CEF Model

$$\begin{aligned} y &= \mathbf{x}'\boldsymbol{\beta} + e \\ \mathbb{E}(e \mid \mathbf{x}) &= 0 \\ \mathbb{E}(e^2 \mid \mathbf{x}) &= \sigma^2 \end{aligned}$$

---

<sup>10</sup>The order doesn't matter. It could be any element.

## 2.16 Linear CEF with Nonlinear Effects

The linear CEF model of the previous section is less restrictive than it might appear, as we can include as regressors nonlinear transformations of the original variables. In this sense, the linear CEF framework is flexible and can capture many nonlinear effects.

For example, suppose we have two scalar variables  $x_1$  and  $x_2$ . The CEF could take the quadratic form

$$m(x_1, x_2) = x_1\beta_1 + x_2\beta_2 + x_1^2\beta_3 + x_2^2\beta_4 + x_1x_2\beta_5 + \beta_6. \quad (2.18)$$

This equation is quadratic in the regressors  $(x_1, x_2)$  yet linear in the coefficients  $\beta = (\beta_1, \dots, \beta_6)'$ . We will descriptively call (2.18) a quadratic CEF, and yet (2.18) is also a linear CEF in the sense of being linear in the coefficients. The key is to understand that (2.18) is quadratic in the variables  $(x_1, x_2)$  yet linear in the coefficients  $\beta$ .

To simplify the expression, we define the transformations  $x_3 = x_1^2$ ,  $x_4 = x_2^2$ ,  $x_5 = x_1x_2$ , and  $x_6 = 1$ , and redefine the regressor vector as  $\mathbf{x} = (x_1, \dots, x_6)'$ . With this redefinition,

$$m(x_1, x_2) = \mathbf{x}'\beta$$

which is linear in  $\beta$ . For most econometric purposes (estimation and inference on  $\beta$ ) the linearity in  $\beta$  is all that is important.

An exception is in the analysis of regression derivatives. In nonlinear equations such as (2.18), the regression derivative should be defined with respect to the original variables, not with respect to the transformed variables. Thus

$$\begin{aligned} \frac{\partial}{\partial x_1} m(x_1, x_2) &= \beta_1 + 2x_1\beta_3 + x_2\beta_5 \\ \frac{\partial}{\partial x_2} m(x_1, x_2) &= \beta_2 + 2x_2\beta_4 + x_1\beta_5. \end{aligned}$$

We see that in the model (2.18), the regression derivatives are not a simple coefficient, but are functions of several coefficients plus the levels of  $(x_1, x_2)$ . Consequently it is difficult to interpret the coefficients individually. It is more useful to interpret them as a group.

We typically call  $\beta_5$  the **interaction effect**. Notice that it appears in both regression derivative equations, and has a symmetric interpretation in each. If  $\beta_5 > 0$  then the regression derivative with respect to  $x_1$  is increasing in the level of  $x_2$  (and the regression derivative with respect to  $x_2$  is increasing in the level of  $x_1$ ), while if  $\beta_5 < 0$  the reverse is true.

## 2.17 Linear CEF with Dummy Variables

When all regressors take a finite set of values, it turns out the CEF can be written as a linear function of regressors.

This simplest example is a **binary** variable, which takes only two distinct values. For example, in most data sets the variable *sex* takes only the values *man* and *woman* (or male and female). Binary variables are extremely common in econometric applications, and are alternatively called **dummy variables** or **indicator variables**.

Consider the simple case of a single binary regressor. In this case, the conditional mean can only take two distinct values. For example,

$$\mathbb{E}(y \mid \text{sex}) = \begin{cases} \mu_0 & \text{if } \text{sex}=\text{man} \\ \mu_1 & \text{if } \text{sex}=\text{woman} \end{cases}.$$

To facilitate a mathematical treatment, we typically record dummy variables with the values  $\{0, 1\}$ . For example

$$x_1 = \begin{cases} 0 & \text{if } \text{sex}=\text{man} \\ 1 & \text{if } \text{sex}=\text{woman} \end{cases}. \quad (2.19)$$

Given this notation we can write the conditional mean as a linear function of the dummy variable  $x_1$ , that is

$$\mathbb{E}(y \mid x_1) = \beta_1 x_1 + \beta_2$$

where  $\beta_1 = \mu_1 - \mu_0$  and  $\beta_2 = \mu_0$ . In this simple regression equation the intercept  $\beta_2$  is equal to the conditional mean of  $y$  for the  $x_1 = 0$  subpopulation (men) and the slope  $\beta_1$  is equal to the difference in the conditional means between the two subpopulations.

Equivalently, we could have defined  $x_1$  as

$$x_1 = \begin{cases} 1 & \text{if } sex=man \\ 0 & \text{if } sex=woman \end{cases} . \quad (2.20)$$

In this case, the regression intercept is the mean for women (rather than for men) and the regression slope has switched signs. The two regressions are equivalent but the interpretation of the coefficients has changed. Therefore it is always important to understand the precise definitions of the variables, and illuminating labels are helpful. For example, labelling  $x_1$  as “sex” does not help distinguish between definitions (2.19) and (2.20). Instead, it is better to label  $x_1$  as “women” or “female” if definition (2.19) is used, or as “men” or “male” if (2.20) is used.

Now suppose we have two dummy variables  $x_1$  and  $x_2$ . For example,  $x_2 = 1$  if the person is married, else  $x_2 = 0$ . The conditional mean given  $x_1$  and  $x_2$  takes at most four possible values:

$$\mathbb{E}(y \mid x_1, x_2) = \begin{cases} \mu_{00} & \text{if } x_1 = 0 \text{ and } x_2 = 0 & (\text{unmarried men}) \\ \mu_{01} & \text{if } x_1 = 0 \text{ and } x_2 = 1 & (\text{married men}) \\ \mu_{10} & \text{if } x_1 = 1 \text{ and } x_2 = 0 & (\text{unmarried women}) \\ \mu_{11} & \text{if } x_1 = 1 \text{ and } x_2 = 1 & (\text{married women}) \end{cases} .$$

In this case we can write the conditional mean as a linear function of  $x_1$ ,  $x_2$  and their product  $x_1 x_2$ :

$$\mathbb{E}(y \mid x_1, x_2) = \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \beta_4$$

where  $\beta_1 = \mu_{10} - \mu_{00}$ ,  $\beta_2 = \mu_{01} - \mu_{00}$ ,  $\beta_3 = \mu_{11} - \mu_{10} - \mu_{01} + \mu_{00}$ , and  $\beta_4 = \mu_{00}$ .

We can view the coefficient  $\beta_1$  as the effect of sex on expected log wages for unmarried wage earners, the coefficient  $\beta_2$  as the effect of marriage on expected log wages for men wage earners, and the coefficient  $\beta_3$  as the difference between the effects of marriage on expected log wages among women and among men. Alternatively, it can also be interpreted as the difference between the effects of sex on expected log wages among married and non-married wage earners. Both interpretations are equally valid. We often describe  $\beta_3$  as measuring the **interaction** between the two dummy variables, or the **interaction effect**, and describe  $\beta_3 = 0$  as the case when the interaction effect is zero.

In this setting we can see that the CEF is linear in the three variables  $(x_1, x_2, x_1 x_2)$ . Thus to put the model in the framework of Section 2.15, we would define the regressor  $x_3 = x_1 x_2$  and the regressor vector as

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ 1 \end{pmatrix} .$$

So even though we started with only 2 dummy variables, the number of regressors (including the intercept) is 4.

If there are 3 dummy variables  $x_1, x_2, x_3$ , then  $\mathbb{E}(y \mid x_1, x_2, x_3)$  takes at most  $2^3 = 8$  distinct values and can be written as the linear function

$$\mathbb{E}(y \mid x_1, x_2, x_3) = \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_1 x_2 + \beta_5 x_1 x_3 + \beta_6 x_2 x_3 + \beta_7 x_1 x_2 x_3 + \beta_8$$

which has eight regressors including the intercept.



In general, if there are  $p$  dummy variables  $x_1, \dots, x_p$  then the CEF  $\mathbb{E}(y \mid x_1, x_2, \dots, x_p)$  takes at most  $2^p$  distinct values, and can be written as a linear function of the  $2^p$  regressors including  $x_1, x_2, \dots, x_p$  and all cross-products. This might be excessive in practice if  $p$  is modestly large. In the next section we will discuss projection approximations which yield more parsimonious parameterizations.

We started this section by saying that the conditional mean is linear whenever all regressors take only a finite number of possible values. How can we see this? Take a **categorical** variable, such as *race*. For example, we earlier divided race into three categories. We can record categorical variables using numbers to indicate each category, for example

$$x_3 = \begin{cases} 1 & \text{if } white \\ 2 & \text{if } black \\ 3 & \text{if } other \end{cases}.$$

When doing so, the values of  $x_3$  have no meaning in terms of magnitude, they simply indicate the relevant category.

When the regressor is categorical the conditional mean of  $y$  given  $x_3$  takes a distinct value for each possibility:

$$\mathbb{E}(y \mid x_3) = \begin{cases} \mu_1 & \text{if } x_3 = 1 \\ \mu_2 & \text{if } x_3 = 2 \\ \mu_3 & \text{if } x_3 = 3 \end{cases}.$$

This is not a linear function of  $x_3$  itself, but it can be made a linear function by constructing dummy variables for two of the three categories. For example

$$x_4 = \begin{cases} 1 & \text{if } black \\ 0 & \text{if } not\ black \end{cases}$$

$$x_5 = \begin{cases} 1 & \text{if } other \\ 0 & \text{if } not\ other \end{cases}.$$

In this case, the categorical variable  $x_3$  is equivalent to the pair of dummy variables  $(x_4, x_5)$ . The explicit relationship is

$$x_3 = \begin{cases} 1 & \text{if } x_4 = 0 \text{ and } x_5 = 0 \\ 2 & \text{if } x_4 = 1 \text{ and } x_5 = 0 \\ 3 & \text{if } x_4 = 0 \text{ and } x_5 = 1 \end{cases}.$$

Given these transformations, we can write the conditional mean of  $y$  as a linear function of  $x_4$  and  $x_5$

$$\mathbb{E}(y \mid x_3) = \mathbb{E}(y \mid x_4, x_5) = \beta_1 x_4 + \beta_2 x_5 + \beta_3.$$

We can write the CEF as either  $\mathbb{E}(y \mid x_3)$  or  $\mathbb{E}(y \mid x_4, x_5)$  (they are equivalent), but it is only linear as a function of  $x_4$  and  $x_5$ .

This setting is similar to the case of two dummy variables, with the difference that we have not included the interaction term  $x_4 x_5$ . This is because the event  $\{x_4 = 1 \text{ and } x_5 = 1\}$  is empty by construction, so  $x_4 x_5 = 0$  by definition.

## 2.18 Best Linear Predictor

While the conditional mean  $m(\mathbf{x}) = \mathbb{E}(y \mid \mathbf{x})$  is the best predictor of  $y$  among all functions of  $\mathbf{x}$ , its functional form is typically unknown. In particular, the linear CEF model is empirically unlikely to be accurate unless  $\mathbf{x}$  is discrete and low-dimensional so all interactions are included. Consequently in most cases it is more realistic to view the linear specification (2.16) as an approximation. In this section we derive a specific approximation with a simple interpretation.

Theorem 2.11.1 showed that the conditional mean  $m(\mathbf{x})$  is the best predictor in the sense that it has the lowest mean squared error among all predictors. By extension, we can define an approximation to the CEF by the linear function with the lowest mean squared error among all linear predictors.

For this derivation we require the following regularity condition.

**Assumption 2.18.1**

1.  $\mathbb{E}(y^2) < \infty$ .
2.  $\mathbb{E}\|\mathbf{x}\|^2 < \infty$ .
3.  $\mathbf{Q}_{\mathbf{x}\mathbf{x}} = \mathbb{E}(\mathbf{x}\mathbf{x}')$  is positive definite.

In Assumption 2.18.1.2 we use the notation  $\|\mathbf{x}\| = (\mathbf{x}'\mathbf{x})^{1/2}$  to denote the Euclidean length of the vector  $\mathbf{x}$ .

The first two parts of Assumption 2.18.1 imply that the variables  $y$  and  $\mathbf{x}$  have finite means, variances, and covariances. The third part of the assumption is more technical, and its role will become apparent shortly. It is equivalent to imposing that the columns of the matrix  $\mathbf{Q}_{\mathbf{x}\mathbf{x}} = \mathbb{E}(\mathbf{x}\mathbf{x}')$  are linearly independent, or that the matrix is invertible.

A linear predictor for  $y$  is a function of the form  $\mathbf{x}'\boldsymbol{\beta}$  for some  $\boldsymbol{\beta} \in \mathbb{R}^k$ . The mean squared prediction error is

$$S(\boldsymbol{\beta}) = \mathbb{E}\left((y - \mathbf{x}'\boldsymbol{\beta})^2\right).$$

The **best linear predictor** of  $y$  given  $\mathbf{x}$ , written  $\mathcal{P}(y \mid \mathbf{x})$ , is found by selecting the vector  $\boldsymbol{\beta}$  to minimize  $S(\boldsymbol{\beta})$ .

**Definition 2.18.1** *The **Best Linear Predictor** of  $y$  given  $\mathbf{x}$  is*

$$\mathcal{P}(y \mid \mathbf{x}) = \mathbf{x}'\boldsymbol{\beta}$$

*where  $\boldsymbol{\beta}$  minimizes the mean squared prediction error*

$$S(\boldsymbol{\beta}) = \mathbb{E}\left((y - \mathbf{x}'\boldsymbol{\beta})^2\right).$$

*The minimizer*

$$\boldsymbol{\beta} = \underset{\mathbf{b} \in \mathbb{R}^k}{\operatorname{argmin}} S(\mathbf{b}) \tag{2.21}$$

*is called the **Linear Projection Coefficient**.*

We now calculate an explicit expression for its value. The mean squared prediction error can be written out as a quadratic function of  $\boldsymbol{\beta}$ :

$$S(\boldsymbol{\beta}) = \mathbb{E}(y^2) - 2\boldsymbol{\beta}'\mathbb{E}(\mathbf{x}y) + \boldsymbol{\beta}'\mathbb{E}(\mathbf{x}\mathbf{x}')\boldsymbol{\beta}.$$

The quadratic structure of  $S(\boldsymbol{\beta})$  means that we can solve explicitly for the minimizer. The first-order condition for minimization (from Appendix A.15) is

$$\mathbf{0} = \frac{\partial}{\partial \boldsymbol{\beta}} S(\boldsymbol{\beta}) = -2\mathbb{E}(\mathbf{x}y) + 2\mathbb{E}(\mathbf{x}\mathbf{x}')\boldsymbol{\beta}. \tag{2.22}$$

Rewriting (2.22) as

$$2\mathbb{E}(xy) = 2\mathbb{E}(xx')\beta$$

and dividing by 2, this equation takes the form

$$\mathbf{Q}_{xy} = \mathbf{Q}_{xx}\beta \quad (2.23)$$

where  $\mathbf{Q}_{xy} = \mathbb{E}(xy)$  is  $k \times 1$  and  $\mathbf{Q}_{xx} = \mathbb{E}(xx')$  is  $k \times k$ . The solution is found by inverting the matrix  $\mathbf{Q}_{xx}$ , and is written

$$\beta = \mathbf{Q}_{xx}^{-1}\mathbf{Q}_{xy}$$

or

$$\beta = (\mathbb{E}(xx'))^{-1}\mathbb{E}(xy). \quad (2.24)$$

It is worth taking the time to understand the notation involved in the expression (2.24).  $\mathbf{Q}_{xx}$  is a  $k \times k$  matrix and  $\mathbf{Q}_{xy}$  is a  $k \times 1$  column vector. Therefore, alternative expressions such as  $\frac{\mathbb{E}(xy)}{\mathbb{E}(xx')}$  or  $\mathbb{E}(xy)(\mathbb{E}(xx'))^{-1}$  are incoherent and incorrect. We also can now see the role of Assumption 2.18.1.3. It is equivalent to assuming that  $\mathbf{Q}_{xx}$  has an inverse  $\mathbf{Q}_{xx}^{-1}$  which is necessary for the normal equations (2.23) to have a solution or equivalently for (2.24) to be uniquely defined. In the absence of Assumption 2.18.1.3 there could be multiple solutions to the equation (2.23).

We now have an explicit expression for the best linear predictor:

$$\mathcal{P}(y | x) = x'(\mathbb{E}(xx'))^{-1}\mathbb{E}(xy).$$

This expression is also referred to as the **linear projection** of  $y$  on  $x$ .

The **projection error** is

$$e = y - x'\beta. \quad (2.25)$$

This equals the error (2.11) from the regression equation when (and only when) the conditional mean is linear in  $x$ , otherwise they are distinct.

Rewriting, we obtain a decomposition of  $y$  into linear predictor and error

$$y = x'\beta + e. \quad (2.26)$$

In general we call equation (2.26) or  $x'\beta$  the best linear predictor of  $y$  given  $x$ , or the linear projection of  $y$  on  $x$ . Equation (2.26) is also often called the **regression** of  $y$  on  $x$  but this can sometimes be confusing as economists use the term *regression* in many contexts. (Recall that we said in Section 2.15 that the linear CEF model is also called the linear regression model.)

An important property of the projection error  $e$  is

$$\mathbb{E}(xe) = \mathbf{0}. \quad (2.27)$$

To see this, using the definitions (2.25) and (2.24) and the matrix properties  $\mathbf{A}\mathbf{A}^{-1} = \mathbf{I}$  and  $\mathbf{I}\mathbf{a} = \mathbf{a}$ ,

$$\begin{aligned} \mathbb{E}(xe) &= \mathbb{E}(x(y - x'\beta)) \\ &= \mathbb{E}(xy) - \mathbb{E}(xx')(\mathbb{E}(xx'))^{-1}\mathbb{E}(xy) \\ &= \mathbf{0} \end{aligned} \quad (2.28)$$

as claimed.

Equation (2.27) is a set of  $k$  equations, one for each regressor. In other words, (2.27) is equivalent to

$$\mathbb{E}(x_j e) = 0 \quad (2.29)$$

for  $j = 1, \dots, k$ . As in (2.15), the regressor vector  $\mathbf{x}$  typically contains a constant, e.g.  $x_k = 1$ . In this case (2.29) for  $j = k$  is the same as

$$\mathbb{E}(e) = 0. \quad (2.30)$$

Thus the projection error has a mean of zero when the regressor vector contains a constant. (When  $\mathbf{x}$  does not have a constant, (2.30) is not guaranteed. As it is desirable for  $e$  to have a zero mean, this is a good reason to always include a constant in any regression model.)

It is also useful to observe that since  $\text{cov}(x_j, e) = \mathbb{E}(x_j e) - \mathbb{E}(x_j)\mathbb{E}(e)$ , then (2.29)-(2.30) together imply that the variables  $x_j$  and  $e$  are uncorrelated.

This completes the derivation of the model. We summarize some of the most important properties.

**Theorem 2.18.1 Properties of Linear Projection Model**

*Under Assumption 2.18.1,*

1. The moments  $\mathbb{E}(\mathbf{x}\mathbf{x}')$  and  $\mathbb{E}(\mathbf{x}y)$  exist with finite elements.
2. The Linear Projection Coefficient (2.21) exists, is unique, and equals

$$\boldsymbol{\beta} = (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}y).$$

3. The best linear predictor of  $y$  given  $\mathbf{x}$  is

$$\mathcal{P}(y | \mathbf{x}) = \mathbf{x}' (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}y).$$

4. The projection error  $e = y - \mathbf{x}'\boldsymbol{\beta}$  exists and satisfies

$$\mathbb{E}(e^2) < \infty$$

and

$$\mathbb{E}(\mathbf{x}e) = \mathbf{0}.$$

5. If  $\mathbf{x}$  contains an constant, then

$$\mathbb{E}(e) = 0.$$

6. If  $\mathbb{E}|y|^r < \infty$  and  $\mathbb{E}\|\mathbf{x}\|^r < \infty$  for  $r \geq 2$  then  $\mathbb{E}|e|^r < \infty$ .

A complete proof of Theorem 2.18.1 is given in Section 2.34.

It is useful to reflect on the generality of Theorem 2.18.1. The only restriction is Assumption 2.18.1. Thus for any random variables  $(y, \mathbf{x})$  with finite variances we can define a linear equation (2.26) with the properties listed in Theorem 2.18.1. Stronger assumptions (such as the linear CEF model) are not necessary. In this sense the linear model (2.26) exists quite generally. However, it is important not to misinterpret the generality of this statement. The linear equation (2.26) is defined as the best linear predictor. It is not necessarily a conditional mean, nor a parameter of a structural or causal economic model.

**Linear Projection Model**

$$y = \mathbf{x}'\boldsymbol{\beta} + e.$$

$$\mathbb{E}(\mathbf{x}e) = \mathbf{0}$$

$$\boldsymbol{\beta} = (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}y)$$

We illustrate projection using three log wage equations introduced in earlier sections.

For our first example, we consider a model with the two dummy variables for sex and race similar to Table 2.1. As we learned in Section 2.17, the entries in this table can be equivalently expressed by a linear CEF. For simplicity, let's consider the CEF of  $\log(\text{wage})$  as a function of *Black* and *Female*.

$$\mathbb{E}(\log(\text{wage}) \mid \text{Black}, \text{Female}) = -0.20\text{Black} - 0.24\text{Female} + 0.10\text{Black} \times \text{Female} + 3.06. \quad (2.31)$$

This is a CEF as the variables are binary and all interactions are included.

Now consider a simpler model omitting the interaction effect. This is the linear projection on the variables *Black* and *Female*

$$\mathcal{P}(\log(\text{wage}) \mid \text{Black}, \text{Female}) = -0.15\text{Black} - 0.23\text{Female} + 3.06. \quad (2.32)$$

What is the difference? The full CEF (2.31) shows that the race gap is differentiated by sex: it is 20% for black men (relative to non-black men) and 10% for black women (relative to non-black women). The projection model (2.32) simplifies this analysis, calculating an average 15% wage gap for blacks, ignoring the role of sex. Notice that this is despite the fact that the sex variable is included in (2.32).

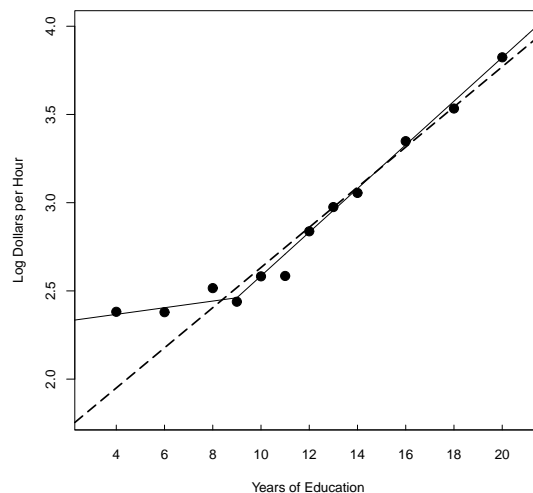


Figure 2.8: Projections of  $\log(\text{wage})$  onto Education

For our second example we consider the CEF of log wages as a function of years of education for white men which was illustrated in Figure 2.5 and is repeated in Figure 2.8. Superimposed on the figure are two projections. The first (given by the dashed line) is the linear projection of log wages on years of education

$$\mathcal{P}(\log(\text{wage}) \mid \text{Education}) = 0.11\text{Education} + 1.5.$$

This simple equation indicates an average 11% increase in wages for every year of education. An inspection of the Figure shows that this approximation works well for  $education \geq 9$ , but under-predicts for individuals with lower levels of education. To correct this imbalance we use a linear spline equation which allows different rates of return above and below 9 years of education:

$$\begin{aligned} \mathcal{P}(\log(wage) \mid Education, (Education - 9) \times 1(Education > 9)) \\ = 0.02Education + 0.10 \times (Education - 9) \times 1(Education > 9) + 2.3. \end{aligned}$$

This equation is displayed in Figure 2.8 using the solid line, and appears to fit much better. It indicates a 2% increase in mean wages for every year of education below 9, and a 12% increase in mean wages for every year of education above 9. It is still an approximation to the conditional mean but it appears to be fairly reasonable.

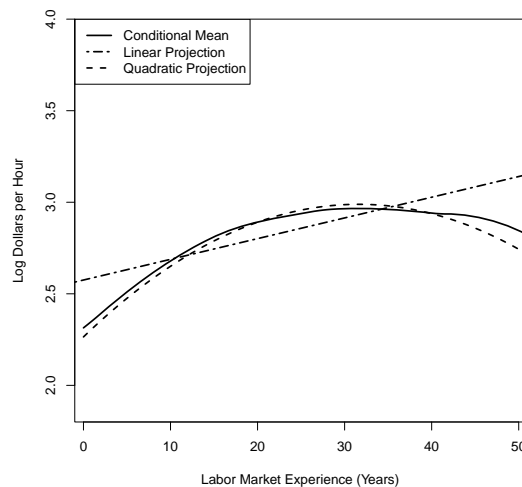


Figure 2.9: Linear and Quadratic Projections of  $\log(wage)$  onto Experience

For our third example we take the CEF of log wages as a function of years of experience for white men with 12 years of education, which was illustrated in Figure 2.6 and is repeated as the solid line in Figure 2.9. Superimposed on the figure are two projections. The first (given by the dot-dashed line) is the linear projection on experience

$$\mathcal{P}(\log(wage) \mid Experience) = 0.011Experience + 2.5$$

and the second (given by the dashed line) is the linear projection on experience and its square

$$\mathcal{P}(\log(wage) \mid Experience) = 0.046Experience - 0.0007Experience^2 + 2.3.$$

It is fairly clear from an examination of Figure 2.9 that the first linear projection is a poor approximation. It over-predicts wages for young and old workers, and under-predicts for the rest. Most importantly, it misses the strong downturn in expected wages for older wage-earners. The second projection fits much better. We can call this equation a **quadratic projection** since the function is quadratic in *experience*.

### Invertibility and Identification

The linear projection coefficient  $\beta = (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}y)$  exists and is unique as long as the  $k \times k$  matrix  $\mathbf{Q}_{xx} = \mathbb{E}(\mathbf{x}\mathbf{x}')$  is invertible. The matrix  $\mathbf{Q}_{xx}$  is sometimes called the **design matrix**, as in experimental settings the researcher is able to control  $\mathbf{Q}_{xx}$  by manipulating the distribution of the regressors  $\mathbf{x}$ .

Observe that for any non-zero  $\alpha \in \mathbb{R}^k$ ,

$$\alpha' \mathbf{Q}_{xx} \alpha = \mathbb{E}(\alpha' \mathbf{x} \mathbf{x}' \alpha) = \mathbb{E}(\alpha' \mathbf{x})^2 \geq 0$$

so  $\mathbf{Q}_{xx}$  by construction is positive semi-definite. The assumption that it is positive definite means that this is a strict inequality,  $\mathbb{E}(\alpha' \mathbf{x})^2 > 0$ . Equivalently, there cannot exist a non-zero vector  $\alpha$  such that  $\alpha' \mathbf{x} = 0$  identically. This occurs when redundant variables are included in  $\mathbf{x}$ . Positive semi-definite matrices are invertible if and only if they are positive definite. When  $\mathbf{Q}_{xx}$  is invertible then  $\beta = (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}y)$  exists and is uniquely defined. In other words, in order for  $\beta$  to be uniquely defined, we must exclude the degenerate situation of redundant variables.

Theorem 2.18.1 shows that the linear projection coefficient  $\beta$  is **identified** (uniquely determined) under Assumption 2.18.1. The key is invertibility of  $\mathbf{Q}_{xx}$ . Otherwise, there is no unique solution to the equation

$$\mathbf{Q}_{xx} \beta = \mathbf{Q}_{xy}. \quad (2.33)$$

When  $\mathbf{Q}_{xx}$  is not invertible there are multiple solutions to (2.33), all of which yield an equivalent best linear predictor  $\mathbf{x}'\beta$ . In this case the coefficient  $\beta$  is **not identified** as it does not have a unique value. Even so, the best linear predictor  $\mathbf{x}'\beta$  still identified. One solution is to set

$$\beta = (\mathbb{E}(\mathbf{x}\mathbf{x}'))^- \mathbb{E}(\mathbf{x}y)$$

where  $\mathbf{A}^-$  denotes the generalized inverse of  $\mathbf{A}$  (see Appendix A.6).

## 2.19 Linear Predictor Error Variance

As in the CEF model, we define the error variance as

$$\sigma^2 = \mathbb{E}(e^2).$$

Setting  $Q_{yy} = \mathbb{E}(y^2)$  and  $\mathbf{Q}_{yx} = \mathbb{E}(y\mathbf{x}')$  we can write  $\sigma^2$  as

$$\begin{aligned} \sigma^2 &= \mathbb{E}\left((y - \mathbf{x}'\beta)^2\right) \\ &= \mathbb{E}(y^2) - 2\mathbb{E}(y\mathbf{x}')\beta + \beta'\mathbb{E}(\mathbf{x}\mathbf{x}')\beta \\ &= Q_{yy} - 2\mathbf{Q}_{yx}\mathbf{Q}_{xx}^{-1}\mathbf{Q}_{xy} + \mathbf{Q}_{yx}\mathbf{Q}_{xx}^{-1}\mathbf{Q}_{xx}\mathbf{Q}_{xx}^{-1}\mathbf{Q}_{xy} \\ &= Q_{yy} - \mathbf{Q}_{yx}\mathbf{Q}_{xx}^{-1}\mathbf{Q}_{xy} \\ &\stackrel{def}{=} Q_{yy \cdot x}. \end{aligned} \quad (2.34)$$

One useful feature of this formula is that it shows that  $Q_{yy \cdot x} = Q_{yy} - \mathbf{Q}_{yx}\mathbf{Q}_{xx}^{-1}\mathbf{Q}_{xy}$  equals the variance of the error from the linear projection of  $y$  on  $\mathbf{x}$ .

## 2.20 Regression Coefficients

Sometimes it is useful to separate the constant from the other regressors, and write the linear projection equation in the format

$$y = \mathbf{x}'\boldsymbol{\beta} + \alpha + e \quad (2.35)$$

where  $\alpha$  is the intercept and  $\mathbf{x}$  does not contain a constant.

Taking expectations of this equation, we find

$$\mathbb{E}(y) = \mathbb{E}(\mathbf{x}'\boldsymbol{\beta}) + \mathbb{E}(\alpha) + \mathbb{E}(e)$$

or

$$\mu_y = \boldsymbol{\mu}_x'\boldsymbol{\beta} + \alpha$$

where  $\mu_y = \mathbb{E}(y)$  and  $\boldsymbol{\mu}_x = \mathbb{E}(\mathbf{x})$ , since  $\mathbb{E}(e) = 0$  from (2.30). (While  $\mathbf{x}$  does not contain a constant, the equation does so (2.30) still applies.) Rearranging, we find

$$\alpha = \mu_y - \boldsymbol{\mu}_x'\boldsymbol{\beta}.$$

Subtracting this equation from (2.35) we find

$$y - \mu_y = (\mathbf{x} - \boldsymbol{\mu}_x)'\boldsymbol{\beta} + e, \quad (2.36)$$

a linear equation between the centered variables  $y - \mu_y$  and  $\mathbf{x} - \boldsymbol{\mu}_x$ . (They are centered at their means, so are mean-zero random variables.) Because  $\mathbf{x} - \boldsymbol{\mu}_x$  is uncorrelated with  $e$ , (2.36) is also a linear projection, thus by the formula for the linear projection model,

$$\begin{aligned} \boldsymbol{\beta} &= \left( \mathbb{E}((\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)') \right)^{-1} \mathbb{E}((\mathbf{x} - \boldsymbol{\mu}_x)(y - \mu_y)) \\ &= \text{var}(\mathbf{x})^{-1} \text{cov}(\mathbf{x}, y) \end{aligned}$$

a function only of the covariances<sup>11</sup> of  $\mathbf{x}$  and  $y$ .

**Theorem 2.20.1** *In the linear projection model*

$$y = \mathbf{x}'\boldsymbol{\beta} + \alpha + e,$$

*then*

$$\alpha = \mu_y - \boldsymbol{\mu}_x'\boldsymbol{\beta} \quad (2.37)$$

*and*

$$\boldsymbol{\beta} = \text{var}(\mathbf{x})^{-1} \text{cov}(\mathbf{x}, y). \quad (2.38)$$

## 2.21 Regression Sub-Vectors

Let the regressors be partitioned as

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}. \quad (2.39)$$

<sup>11</sup>The **covariance matrix** between vectors  $\mathbf{x}$  and  $\mathbf{z}$  is  $\text{cov}(\mathbf{x}, \mathbf{z}) = \mathbb{E}((\mathbf{x} - \mathbb{E}\mathbf{x})(\mathbf{z} - \mathbb{E}\mathbf{z})')$ . The (co)variance matrix of the vector  $\mathbf{x}$  is  $\text{var}(\mathbf{x}) = \text{cov}(\mathbf{x}, \mathbf{x}) = \mathbb{E}((\mathbf{x} - \mathbb{E}\mathbf{x})(\mathbf{x} - \mathbb{E}\mathbf{x})')$ .



We can write the projection of  $y$  on  $\mathbf{x}$  as

$$\begin{aligned} y &= \mathbf{x}'\boldsymbol{\beta} + e \\ &= \mathbf{x}'_1\beta_1 + \mathbf{x}'_2\beta_2 + e \\ \mathbb{E}(\mathbf{x}e) &= \mathbf{0}. \end{aligned} \tag{2.40}$$

In this section we derive formula for the sub-vectors  $\beta_1$  and  $\beta_2$ .

Partition  $\mathbf{Q}_{\mathbf{x}\mathbf{x}}$  conformably with  $\mathbf{x}$

$$\mathbf{Q}_{\mathbf{x}\mathbf{x}} = \begin{bmatrix} \mathbf{Q}_{11} & \mathbf{Q}_{12} \\ \mathbf{Q}_{21} & \mathbf{Q}_{22} \end{bmatrix} = \begin{bmatrix} \mathbb{E}(\mathbf{x}_1\mathbf{x}'_1) & \mathbb{E}(\mathbf{x}_1\mathbf{x}'_2) \\ \mathbb{E}(\mathbf{x}_2\mathbf{x}'_1) & \mathbb{E}(\mathbf{x}_2\mathbf{x}'_2) \end{bmatrix}$$

and similarly  $\mathbf{Q}_{\mathbf{x}y}$

$$\mathbf{Q}_{\mathbf{x}y} = \begin{bmatrix} \mathbf{Q}_{1y} \\ \mathbf{Q}_{2y} \end{bmatrix} = \begin{bmatrix} \mathbb{E}(\mathbf{x}_1y) \\ \mathbb{E}(\mathbf{x}_2y) \end{bmatrix}.$$

By the partitioned matrix inversion formula (A.4)

$$\mathbf{Q}_{\mathbf{x}\mathbf{x}}^{-1} = \begin{bmatrix} \mathbf{Q}_{11} & \mathbf{Q}_{12} \\ \mathbf{Q}_{21} & \mathbf{Q}_{22} \end{bmatrix}^{-1} \stackrel{def}{=} \begin{bmatrix} \mathbf{Q}_{11}^{-1} & \mathbf{Q}_{12}^{-1} \\ \mathbf{Q}_{21}^{-1} & \mathbf{Q}_{22}^{-1} \end{bmatrix} = \begin{bmatrix} \mathbf{Q}_{11}^{-1} & -\mathbf{Q}_{11}^{-1}\mathbf{Q}_{12}\mathbf{Q}_{22}^{-1} \\ -\mathbf{Q}_{22}^{-1}\mathbf{Q}_{21}\mathbf{Q}_{11}^{-1} & \mathbf{Q}_{22}^{-1} \end{bmatrix}. \tag{2.41}$$

where  $\mathbf{Q}_{11.2} \stackrel{def}{=} \mathbf{Q}_{11} - \mathbf{Q}_{12}\mathbf{Q}_{22}^{-1}\mathbf{Q}_{21}$  and  $\mathbf{Q}_{22.1} \stackrel{def}{=} \mathbf{Q}_{22} - \mathbf{Q}_{21}\mathbf{Q}_{11}^{-1}\mathbf{Q}_{12}$ . Thus

$$\begin{aligned} \boldsymbol{\beta} &= \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} \\ &= \begin{bmatrix} \mathbf{Q}_{11.2}^{-1} & -\mathbf{Q}_{11.2}^{-1}\mathbf{Q}_{12}\mathbf{Q}_{22}^{-1} \\ -\mathbf{Q}_{22.1}^{-1}\mathbf{Q}_{21}\mathbf{Q}_{11}^{-1} & \mathbf{Q}_{22.1}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{Q}_{1y} \\ \mathbf{Q}_{2y} \end{bmatrix} \\ &= \begin{pmatrix} \mathbf{Q}_{11.2}^{-1}(\mathbf{Q}_{1y} - \mathbf{Q}_{12}\mathbf{Q}_{22}^{-1}\mathbf{Q}_{2y}) \\ \mathbf{Q}_{22.1}^{-1}(\mathbf{Q}_{2y} - \mathbf{Q}_{21}\mathbf{Q}_{11}^{-1}\mathbf{Q}_{1y}) \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{Q}_{11.2}^{-1}\mathbf{Q}_{1y.2} \\ \mathbf{Q}_{22.1}^{-1}\mathbf{Q}_{2y.1} \end{pmatrix}. \end{aligned}$$

We have shown that

$$\begin{aligned} \beta_1 &= \mathbf{Q}_{11.2}^{-1}\mathbf{Q}_{1y.2} \\ \beta_2 &= \mathbf{Q}_{22.1}^{-1}\mathbf{Q}_{2y.1}. \end{aligned}$$

## 2.22 Coefficient Decomposition

In the previous section we derived formulae for the coefficient sub-vectors  $\beta_1$  and  $\beta_2$ . We now use these formulae to give a useful interpretation of the coefficients in terms of an iterated projection.

Take equation (2.40) for the case  $\dim(\mathbf{x}_1) = 1$  so that  $\beta_1 \in \mathbb{R}$ .

$$y = x_1\beta_1 + \mathbf{x}'_2\beta_2 + e. \tag{2.42}$$

Now consider the projection of  $x_1$  on  $\mathbf{x}_2$ :

$$\begin{aligned} x_1 &= \mathbf{x}'_2\gamma_2 + u_1 \\ \mathbb{E}(\mathbf{x}_2u_1) &= \mathbf{0}. \end{aligned}$$

From (2.24) and (2.34),  $\gamma_2 = \mathbf{Q}_{22}^{-1}\mathbf{Q}_{21}$  and  $\mathbb{E}u_1^2 = \mathbf{Q}_{11.2} = \mathbf{Q}_{11} - \mathbf{Q}_{12}\mathbf{Q}_{22}^{-1}\mathbf{Q}_{21}$ . We can also calculate that

$$\mathbb{E}(u_1y) = \mathbb{E}((x_1 - \gamma'_2\mathbf{x}_2)y) = \mathbb{E}(x_1y) - \gamma'_2\mathbb{E}(\mathbf{x}_2y) = \mathbf{Q}_{1y} - \mathbf{Q}_{12}\mathbf{Q}_{22}^{-1}\mathbf{Q}_{2y} = \mathbf{Q}_{1y.2}.$$

We have found that

$$\beta_1 = \mathbf{Q}_{11.2}^{-1} \mathbf{Q}_{1y.2} = \frac{\mathbb{E}(u_1 y)}{\mathbb{E}(u_1^2)}$$

the coefficient from the simple regression of  $y$  on  $u_1$ .

What this means is that in the multivariate projection equation (2.42), the coefficient  $\beta_1$  equals the projection coefficient from a regression of  $y$  on  $u_1$ , the error from a projection of  $x_1$  on the other regressors  $\mathbf{x}_2$ . The error  $u_1$  can be thought of as the component of  $x_1$  which is not linearly explained by the other regressors. Thus the coefficient  $\beta_1$  equals the linear effect of  $x_1$  on  $y$ , after stripping out the effects of the other variables.

There was nothing special in the choice of the variable  $x_1$ . This derivation applies symmetrically to all coefficients in a linear projection. Each coefficient equals the simple regression of  $y$  on the error from a projection of that regressor on all the other regressors. Each coefficient equals the linear effect of that variable on  $y$ , after linearly controlling for all the other regressors.

## 2.23 Omitted Variable Bias

Again, let the regressors be partitioned as in (2.39). Consider the projection of  $y$  on  $\mathbf{x}_1$  only. Perhaps this is done because the variables  $\mathbf{x}_2$  are not observed. This is the equation

$$\begin{aligned} y &= \mathbf{x}_1' \boldsymbol{\gamma}_1 + u \\ \mathbb{E}(\mathbf{x}_1 u) &= \mathbf{0}. \end{aligned} \tag{2.43}$$

Notice that we have written the coefficient on  $\mathbf{x}_1$  as  $\boldsymbol{\gamma}_1$  rather than  $\boldsymbol{\beta}_1$  and the error as  $u$  rather than  $e$ . This is because (2.43) is different than (2.40). Goldberger (1991) introduced the catchy labels **long regression** for (2.40) and **short regression** for (2.43) to emphasize the distinction.

Typically,  $\boldsymbol{\beta}_1 \neq \boldsymbol{\gamma}_1$ , except in special cases. To see this, we calculate

$$\begin{aligned} \boldsymbol{\gamma}_1 &= (\mathbb{E}(\mathbf{x}_1 \mathbf{x}_1'))^{-1} \mathbb{E}(\mathbf{x}_1 y) \\ &= (\mathbb{E}(\mathbf{x}_1 \mathbf{x}_1'))^{-1} \mathbb{E}(\mathbf{x}_1 (\mathbf{x}_1' \boldsymbol{\beta}_1 + \mathbf{x}_2' \boldsymbol{\beta}_2 + e)) \\ &= \boldsymbol{\beta}_1 + (\mathbb{E}(\mathbf{x}_1 \mathbf{x}_1'))^{-1} \mathbb{E}(\mathbf{x}_1 \mathbf{x}_2') \boldsymbol{\beta}_2 \\ &= \boldsymbol{\beta}_1 + \boldsymbol{\Gamma}_{12} \boldsymbol{\beta}_2 \end{aligned}$$

where  $\boldsymbol{\Gamma}_{12} = \mathbf{Q}_{11}^{-1} \mathbf{Q}_{12}$  is the coefficient matrix from a projection of  $\mathbf{x}_2$  on  $\mathbf{x}_1$ , where we use the notation from Section 2.21.

Observe that  $\boldsymbol{\gamma}_1 = \boldsymbol{\beta}_1 + \boldsymbol{\Gamma}_{12} \boldsymbol{\beta}_2 \neq \boldsymbol{\beta}_1$  unless  $\boldsymbol{\Gamma}_{12} = \mathbf{0}$  or  $\boldsymbol{\beta}_2 = \mathbf{0}$ . Thus the short and long regressions have different coefficients on  $\mathbf{x}_1$ . They are the same only under one of two conditions. First, if the projection of  $\mathbf{x}_2$  on  $\mathbf{x}_1$  yields a set of zero coefficients (they are uncorrelated), or second, if the coefficient on  $\mathbf{x}_2$  in (2.40) is zero. In general, the coefficient in (2.43) is  $\boldsymbol{\gamma}_1$  rather than  $\boldsymbol{\beta}_1$ . The difference  $\boldsymbol{\Gamma}_{12} \boldsymbol{\beta}_2$  between  $\boldsymbol{\gamma}_1$  and  $\boldsymbol{\beta}_1$  is known as **omitted variable bias**. It is the consequence of omission of a relevant correlated variable.

To avoid omitted variables bias the standard advice is to include all potentially relevant variables in estimated models. By construction, the general model will be free of such bias. Unfortunately in many cases it is not feasible to completely follow this advice as many desired variables are not observed. In this case, the possibility of omitted variables bias should be acknowledged and discussed in the course of an empirical investigation.

For example, suppose  $y$  is log wages,  $x_1$  is education, and  $x_2$  is intellectual ability. It seems reasonable to suppose that education and intellectual ability are positively correlated (highly able individuals attain higher levels of education) which means  $\boldsymbol{\Gamma}_{12} > 0$ . It also seems reasonable to suppose that conditional on education, individuals with higher intelligence will earn higher wages on average, so that  $\beta_2 > 0$ . This implies that  $\boldsymbol{\Gamma}_{12} \boldsymbol{\beta}_2 > 0$  and  $\gamma_1 = \beta_1 + \boldsymbol{\Gamma}_{12} \beta_2 > \beta_1$ . Therefore,

it seems reasonable to expect that in a regression of wages on education with ability omitted, the coefficient on education is higher than in a regression where ability is included. In other words, in this context the omitted variable biases the regression coefficient upwards. It is possible, for example, that  $\beta_1 = 0$  so that education has no direct effect on wages yet  $\gamma_1 = \mathbf{\Gamma}_{12}\beta_2 > 0$  meaning that the regression coefficient on education alone is positive, but is a consequence of the unmodeled correlation between education and intellectual ability.

Unfortunately the above simple characterization of omitted variable bias does not immediately carry over to more complicated settings, as discovered by Luca, Magnus, and Peracchi (2017). For example, suppose we compare three nested projections

$$\begin{aligned} y &= \mathbf{x}'_1 \gamma_1 + u_1 \\ y &= \mathbf{x}'_1 \delta_1 + \mathbf{x}'_2 \delta_2 + u_2 \\ y &= \mathbf{x}'_1 \beta_1 + \mathbf{x}'_2 \beta_2 + \mathbf{x}'_3 \beta_3 + e. \end{aligned}$$

We can call them the short, medium, and long regressions. Suppose that the parameter of interest is  $\beta_1$  in the long regression. We are interested in the consequences of omitting  $\mathbf{x}_3$  when estimating the medium regression, and of omitting both  $\mathbf{x}_2$  and  $\mathbf{x}_3$  when estimating the short regression. In particular we are interested in the question: Is it better to estimate the short or medium regression, given that both omit  $\mathbf{x}_3$ ? Intuition suggests that the medium regression should be “less biased” but it is worth investigating in greater detail. By similar calculations to those above, we find that

$$\begin{aligned} \gamma_1 &= \beta_1 + \mathbf{\Gamma}_{12}\beta_2 + \mathbf{\Gamma}_{13}\beta_3 \\ \delta_1 &= \beta_1 + \mathbf{\Gamma}_{13.2}\beta_3 \end{aligned}$$

where  $\mathbf{\Gamma}_{13.2} = \mathbf{Q}_{11.2}^{-1} \mathbf{Q}_{13.2}$  using the notation from Section 2.21.

We see that the bias in the short regression coefficient is  $\mathbf{\Gamma}_{12}\beta_2 + \mathbf{\Gamma}_{13}\beta_3$  which depends on both  $\beta_2$  and  $\beta_3$ , while that for the medium regression coefficient is  $\mathbf{\Gamma}_{13.2}\beta_3$  which only depends on  $\beta_3$ . So the bias for the medium regression is less complicated, and intuitively seems more likely to be smaller than that of the short regression. However it is impossible to strictly rank the two. It is quite possible that  $\gamma_1$  is less biased than  $\delta_1$ . Thus as a general rule it is strictly impossible to state that estimation of the medium regression will be less biased than estimation of the short regression.

## 2.24 Best Linear Approximation

There are alternative ways we could construct a linear approximation  $\mathbf{x}'\beta$  to the conditional mean  $m(\mathbf{x})$ . In this section we show that one alternative approach turns out to yield the same answer as the best linear predictor.

We start by defining the mean-square approximation error of  $\mathbf{x}'\beta$  to  $m(\mathbf{x})$  as the expected squared difference between  $\mathbf{x}'\beta$  and the conditional mean  $m(\mathbf{x})$

$$d(\beta) = \mathbb{E} \left( (m(\mathbf{x}) - \mathbf{x}'\beta)^2 \right). \quad (2.44)$$

The function  $d(\beta)$  is a measure of the deviation of  $\mathbf{x}'\beta$  from  $m(\mathbf{x})$ . If the two functions are identical then  $d(\beta) = 0$ , otherwise  $d(\beta) > 0$ . We can also view the mean-square difference  $d(\beta)$  as a density-weighted average of the function  $(m(\mathbf{x}) - \mathbf{x}'\beta)^2$ , since

$$d(\beta) = \int_{\mathbb{R}^k} (m(\mathbf{x}) - \mathbf{x}'\beta)^2 f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x}$$

where  $f_{\mathbf{x}}(\mathbf{x})$  is the marginal density of  $\mathbf{x}$ .

We can then define the best linear approximation to the conditional  $m(\mathbf{x})$  as the function  $\mathbf{x}'\beta$  obtained by selecting  $\beta$  to minimize  $d(\beta)$  :

$$\beta = \underset{\mathbf{b} \in \mathbb{R}^k}{\operatorname{argmin}} d(\mathbf{b}). \quad (2.45)$$

Similar to the best linear predictor we are measuring accuracy by expected squared error. The difference is that the best linear predictor (2.21) selects  $\beta$  to minimize the expected squared prediction error, while the best linear approximation (2.45) selects  $\beta$  to minimize the expected squared approximation error.

Despite the different definitions, it turns out that the best linear predictor and the best linear approximation are identical. By the same steps as in (2.18) plus an application of conditional expectations we can find that

$$\beta = (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}m(\mathbf{x})) \quad (2.46)$$

$$= (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}y) \quad (2.47)$$

(see Exercise 2.19). Thus (2.45) equals (2.21). We conclude that the definition (2.45) can be viewed as an alternative motivation for the linear projection coefficient.

## 2.25 Regression to the Mean

The term **regression** originated in an influential paper by Francis Galton (1886), where he examined the joint distribution of the stature (height) of parents and children. Effectively, he was estimating the conditional mean of children's height given their parent's height. Galton discovered that this conditional mean was approximately linear with a slope of  $2/3$ . This implies that *on average* a child's height is more mediocre (average) than his or her parent's height. Galton called this phenomenon **regression to the mean**, and the label **regression** has stuck to this day to describe most conditional relationships.

One of Galton's fundamental insights was to recognize that if the marginal distributions of  $y$  and  $x$  are the same (e.g. the heights of children and parents in a stable environment) then the regression slope in a linear projection is always less than one.

To be more precise, take the simple linear projection

$$y = x\beta + \alpha + e \quad (2.48)$$

where  $y$  equals the height of the child and  $x$  equals the height of the parent. Assume that  $y$  and  $x$  have the same mean, so that  $\mu_y = \mu_x = \mu$ . Then from (2.37)

$$\alpha = (1 - \beta)\mu$$

so we can write the linear projection (2.48) as

$$\mathcal{P}(y | x) = (1 - \beta)\mu + x\beta.$$

This shows that the projected height of the child is a weighted average of the population average height  $\mu$  and the parent's height  $x$ , with the weight equal to the regression slope  $\beta$ . When the height distribution is stable across generations, so that  $\text{var}(y) = \text{var}(x)$ , then this slope is the simple correlation of  $y$  and  $x$ . Using (2.38)

$$\beta = \frac{\text{cov}(x, y)}{\text{var}(x)} = \text{corr}(x, y).$$

By the properties of correlation (e.g. equation (??) in the Appendix),  $-1 \leq \text{corr}(x, y) \leq 1$ , with  $\text{corr}(x, y) = 1$  only in the degenerate case  $y = x$ . Thus if we exclude degeneracy,  $\beta$  is strictly less than 1.

This means that on average a child's height is more mediocre (closer to the population average) than the parent's.

### Sir Francis Galton

Sir Francis Galton (1822-1911) of England was one of the leading figures in late 19th century statistics. In addition to inventing the concept of regression, he is credited with introducing the concepts of correlation, the standard deviation, and the bivariate normal distribution. His work on heredity made a significant intellectual advance by examining the joint distributions of observables, allowing the application of the tools of mathematical statistics to the social sciences.

A common error – known as the **regression fallacy** – is to infer from  $\beta < 1$  that the population is **converging**, meaning that its variance is declining towards zero. This is a fallacy because we derived the implication  $\beta < 1$  under the assumption of constant means and variances. So certainly  $\beta < 1$  does not imply that the variance  $y$  is less than the variance of  $x$ .

Another way of seeing this is to examine the conditions for convergence in the context of equation (2.48). Since  $x$  and  $e$  are uncorrelated, it follows that

$$\text{var}(y) = \beta^2 \text{var}(x) + \text{var}(e).$$

Then  $\text{var}(y) < \text{var}(x)$  if and only if

$$\beta^2 < 1 - \frac{\text{var}(e)}{\text{var}(x)}$$

which is not implied by the simple condition  $|\beta| < 1$ .

The regression fallacy arises in related empirical situations. Suppose you sort families into groups by the heights of the parents, and then plot the average heights of each subsequent generation over time. If the population is stable, the regression property implies that the plots lines will converge – children’s height will be more average than their parents. The regression fallacy is to incorrectly conclude that the population is converging. A message to be learned from this example is that such plots are misleading for inferences about convergence.

The regression fallacy is subtle. It is easy for intelligent economists to succumb to its temptation. A famous example is *The Triumph of Mediocrity in Business* by Horace Secrist, published in 1933. In this book, Secrist carefully and with great detail documented that in a sample of department stores over 1920-1930, when he divided the stores into groups based on 1920-1921 profits, and plotted the average profits of these groups for the subsequent 10 years, he found clear and persuasive evidence for convergence “toward mediocrity”. Of course, there was no discovery – regression to the mean is a necessary feature of stable distributions.

## 2.26 Reverse Regression

Galton noticed another interesting feature of the bivariate distribution. There is nothing special about a regression of  $y$  on  $x$ . We can also regress  $x$  on  $y$ . (In his heredity example this is the best linear predictor of the height of parents given the height of their children.) This regression takes the form

$$x = y\beta^* + \alpha^* + e^*. \tag{2.49}$$

This is sometimes called the **reverse regression**. In this equation, the coefficients  $\alpha^*$ ,  $\beta^*$  and error  $e^*$  are defined by linear projection. In a stable population we find that

$$\beta^* = \text{corr}(x, y) = \beta$$

$$\alpha^* = (1 - \beta)\mu = \alpha$$

which are exactly the same as in the projection of  $y$  on  $x$ ! The intercept and slope have exactly the same values in the forward and reverse projections!

While this algebraic discovery is quite simple, it is counter-intuitive. Instead, a common yet mistaken guess for the form of the reverse regression is to take the equation (2.48), divide through by  $\beta$  and rewrite to find the equation

$$x = y \frac{1}{\beta} - \frac{\alpha}{\beta} - \frac{1}{\beta} e \quad (2.50)$$

suggesting that the projection of  $x$  on  $y$  should have a slope coefficient of  $1/\beta$  instead of  $\beta$ , and intercept of  $-\alpha/\beta$  rather than  $\alpha$ . What went wrong? Equation (2.50) is perfectly valid, because it is a simple manipulation of the valid equation (2.48). The trouble is that (2.50) is neither a CEF nor a linear projection. **Inverting a projection (or CEF) does not yield a projection (or CEF).** Instead, (2.49) is a valid projection, not (2.50).

In any event, Galton's finding was that when the variables are standardized, the slope in both projections ( $y$  on  $x$ , and  $x$  and  $y$ ) equals the correlation, and both equations exhibit regression to the mean. It is not a causal relation, but a natural feature of all joint distributions.

## 2.27 Limitations of the Best Linear Projection

Let's compare the linear projection and linear CEF models.

From Theorem 2.8.1.4 we know that the CEF error has the property  $\mathbb{E}(\mathbf{x}e) = \mathbf{0}$ . Thus a linear CEF is the best linear projection. However, the converse is not true as the projection error does not necessarily satisfy  $\mathbb{E}(e | \mathbf{x}) = 0$ . Furthermore, the linear projection may be a poor approximation to the CEF.

To see these points in a simple example, suppose that the true process is  $y = x + x^2$  with  $x \sim N(0, 1)$ . In this case the true CEF is  $m(x) = x + x^2$  and there is no error. Now consider the linear projection of  $y$  on  $x$  and a constant, namely the model  $y = \beta x + \alpha + u$ . Since  $x \sim N(0, 1)$  then  $x$  and  $x^2$  are uncorrelated and the linear projection takes the form  $\mathcal{P}(y | x) = x + 1$ . This is quite different from the true CEF  $m(x) = x + x^2$ . The projection error equals  $e = x^2 - 1$ , which is a deterministic function of  $x$ , yet is uncorrelated with  $x$ . We see in this example that a projection error need not be a CEF error, and a linear projection can be a poor approximation to the CEF.

Another defect of linear projection is that it is sensitive to the marginal distribution of the regressors when the conditional mean is non-linear. We illustrate the issue in Figure 2.10 for a constructed<sup>12</sup> joint distribution of  $y$  and  $x$ . The solid line is the non-linear CEF of  $y$  given  $x$ . The data are divided in two groups – Group 1 and Group 2 – which have different marginal distributions for the regressor  $x$ , and Group 1 has a lower mean value of  $x$  than Group 2. The separate linear projections of  $y$  on  $x$  for these two groups are displayed in the Figure by the dashed lines. These two projections are distinct approximations to the CEF. A defect with linear projection is that it leads to the incorrect conclusion that the effect of  $x$  on  $y$  is different for individuals in the two groups. This conclusion is incorrect because in fact there is no difference in the conditional mean function. The apparent difference is a by-product of a linear approximation to a nonlinear mean, combined with different marginal distributions for the conditioning variables.

## 2.28 Random Coefficient Model

A model which is notationally similar to but conceptually distinct from the linear CEF model is the linear random coefficient model. It takes the form

$$y = \mathbf{x}'\boldsymbol{\eta}$$

---

<sup>12</sup>The  $x$  in Group 1 are  $N(2, 1)$  and those in Group 2 are  $N(4, 1)$ , and the conditional distribution of  $y$  given  $x$  is  $N(m(x), 1)$  where  $m(x) = 2x - x^2/6$ .

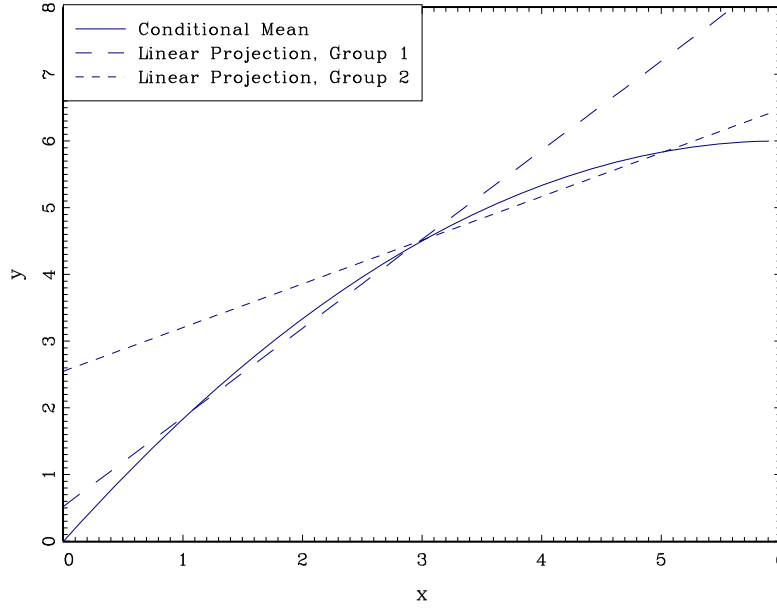


Figure 2.10: Conditional Mean and Two Linear Projections

where the individual-specific coefficient  $\boldsymbol{\eta}$  is random and independent of  $\mathbf{x}$ . For example, if  $\mathbf{x}$  is years of schooling and  $y$  is log wages, then  $\boldsymbol{\eta}$  is the individual-specific returns to schooling. If a person obtains an extra year of schooling,  $\boldsymbol{\eta}$  is the actual change in their wage. The random coefficient model allows the returns to schooling to vary in the population. Some individuals might have a high return to education (a high  $\boldsymbol{\eta}$ ) and others a low return, possibly 0, or even negative.

In the linear CEF model the regressor coefficient equals the regression derivative – the change in the conditional mean due to a change in the regressors,  $\boldsymbol{\beta} = \nabla m(\mathbf{x})$ . This is not the effect on a given individual, it is the effect on the population average. In contrast, in the random coefficient model, the random vector  $\boldsymbol{\eta} = \nabla(\mathbf{x}'\boldsymbol{\eta})$  is the true causal effect – the change in the response variable  $y$  itself due to a change in the regressors.

It is interesting, however, to discover that the linear random coefficient model implies a linear CEF. To see this, let  $\boldsymbol{\beta}$  and  $\boldsymbol{\Sigma}$  denote the mean and covariance matrix of  $\boldsymbol{\eta}$  :

$$\begin{aligned}\boldsymbol{\beta} &= \mathbb{E}(\boldsymbol{\eta}) \\ \boldsymbol{\Sigma} &= \text{var}(\boldsymbol{\eta})\end{aligned}$$

and then decompose the random coefficient as

$$\boldsymbol{\eta} = \boldsymbol{\beta} + \mathbf{u}$$

where  $\mathbf{u}$  is distributed independently of  $\mathbf{x}$  with mean zero and covariance matrix  $\boldsymbol{\Sigma}$ . Then we can write

$$\mathbb{E}(y \mid \mathbf{x}) = \mathbf{x}'\mathbb{E}(\boldsymbol{\eta} \mid \mathbf{x}) = \mathbf{x}'\mathbb{E}(\boldsymbol{\eta}) = \mathbf{x}'\boldsymbol{\beta}$$

so the CEF is linear in  $\mathbf{x}$ , and the coefficients  $\boldsymbol{\beta}$  equal the mean of the random coefficient  $\boldsymbol{\eta}$ .

We can thus write the equation as a linear CEF

$$y = \mathbf{x}'\boldsymbol{\beta} + e \tag{2.51}$$

where  $e = \mathbf{x}'\mathbf{u}$  and  $\mathbf{u} = \boldsymbol{\eta} - \boldsymbol{\beta}$ . The error is conditionally mean zero:

$$\mathbb{E}(e \mid \mathbf{x}) = 0.$$

Furthermore

$$\begin{aligned}\text{var}(e \mid \mathbf{x}) &= \mathbf{x}' \text{var}(\boldsymbol{\eta}) \mathbf{x} \\ &= \mathbf{x}' \boldsymbol{\Sigma} \mathbf{x}\end{aligned}$$

so the error is conditionally heteroskedastic with its variance a quadratic function of  $\mathbf{x}$ .

**Theorem 2.28.1** *In the linear random coefficient model  $y = \mathbf{x}'\boldsymbol{\eta}$  with  $\boldsymbol{\eta}$  independent of  $\mathbf{x}$ ,  $\mathbb{E}\|\mathbf{x}\|^2 < \infty$ , and  $\mathbb{E}\|\boldsymbol{\eta}\|^2 < \infty$ , then*

$$\begin{aligned}\mathbb{E}(y \mid \mathbf{x}) &= \mathbf{x}'\boldsymbol{\beta} \\ \text{var}(y \mid \mathbf{x}) &= \mathbf{x}'\boldsymbol{\Sigma} \mathbf{x}\end{aligned}$$

where  $\boldsymbol{\beta} = \mathbb{E}(\boldsymbol{\eta})$  and  $\boldsymbol{\Sigma} = \text{var}(\boldsymbol{\eta})$ .

## 2.29 Causal Effects

So far we have avoided the concept of causality, yet often the underlying goal of an econometric analysis is to uncover a causal relationship between variables. It is often of great interest to understand the causes and effects of decisions, actions, and policies. For example, we may be interested in the effect of class sizes on test scores, police expenditures on crime rates, climate change on economic activity, years of schooling on wages, institutional structure on growth, the effectiveness of rewards on behavior, the consequences of medical procedures for health outcomes, or any variety of possible causal relationships. In each case, the goal is to understand what is the actual effect on the outcome  $y$  due to a change in the input  $x$ . We are not just interested in the conditional mean or linear projection, we would like to know the actual change.

Two inherent barriers are that the causal effect is typically specific to an individual and that it is unobserved.

Consider the effect of schooling on wages. The causal effect is the actual difference a person would receive in wages if we could change their level of education *holding all else constant*. This is specific to each individual as their employment outcomes in these two distinct situations is individual. The causal effect is unobserved because the most we can observe is their actual level of education and their actual wage, but not the counterfactual wage if their education had been different.

To be even more specific, suppose that there are two individuals, Jennifer and George, and both have the possibility of being high-school graduates or college graduates, but both would have received different wages given their choices. For example, suppose that Jennifer would have earned \$10 an hour as a high-school graduate and \$20 an hour as a college graduate while George would have earned \$8 as a high-school graduate and \$12 as a college graduate. In this example the causal effect of schooling is \$10 a hour for Jennifer and \$4 an hour for George. The causal effects are specific to the individual and neither causal effect is observed.

A variable  $x_1$  can be said to have a causal effect on the response variable  $y$  if the latter changes when all other inputs are held constant. To make this precise we need a mathematical formulation. We can write a full model for the response variable  $y$  as

$$y = h(x_1, \mathbf{x}_2, \mathbf{u}) \tag{2.52}$$

where  $x_1$  and  $\mathbf{x}_2$  are the observed variables,  $\mathbf{u}$  is an  $\ell \times 1$  unobserved random factor, and  $h$  is a functional relationship. This framework, called the **potential outcomes** framework, includes as



a special case the random coefficient model (2.28) studied earlier. We define the causal effect of  $x_1$  within this model as the change in  $y$  due to a change in  $x_1$  holding the other variables  $\mathbf{x}_2$  and  $\mathbf{u}$  constant.

**Definition 2.29.1** *In the model (2.52) the **causal effect** of  $x_1$  on  $y$  is*

$$C(x_1, \mathbf{x}_2, \mathbf{u}) = \nabla_1 h(x_1, \mathbf{x}_2, \mathbf{u}), \quad (2.53)$$

*the change in  $y$  due to a change in  $x_1$ , holding  $\mathbf{x}_2$  and  $\mathbf{u}$  constant.*

To understand this concept, imagine taking a single individual. As far as our structural model is concerned, this person is described by their observables  $x_1$  and  $\mathbf{x}_2$  and their unobservables  $\mathbf{u}$ . In a wage regression the unobservables would include characteristics such as the person's abilities, skills, work ethic, interpersonal connections, and preferences. The causal effect of  $x_1$  (say, education) is the change in the wage as  $x_1$  changes, holding constant all other observables **and** unobservables.

It may be helpful to understand that (2.53) is a definition, and does not necessarily describe causality in a fundamental or experimental sense. Perhaps it would be more appropriate to label (2.53) as a **structural effect** (the effect within the structural model).

Sometimes it is useful to write this relationship as a potential outcome function

$$y(x_1) = h(x_1, \mathbf{x}_2, \mathbf{u})$$

where the notation implies that  $y(x_1)$  is holding  $\mathbf{x}_2$  and  $\mathbf{u}$  constant.

A popular example arises in the analysis of treatment effects with a binary regressor  $x_1$ . Let  $x_1 = 1$  indicate treatment (e.g. a medical procedure) and  $x_1 = 0$  indicate non-treatment. In this case  $y(x_1)$  can be written

$$\begin{aligned} y(0) &= h(0, \mathbf{x}_2, \mathbf{u}) \\ y(1) &= h(1, \mathbf{x}_2, \mathbf{u}). \end{aligned}$$

In the literature on treatment effects, it is common to refer to  $y(0)$  and  $y(1)$  as the latent outcomes associated with non-treatment and treatment, respectively. That is, for a given individual,  $y(0)$  is the health outcome if there is no treatment, and  $y(1)$  is the health outcome if there is treatment. The causal effect of treatment for the individual is the change in their health outcome due to treatment – the change in  $y$  as we hold both  $\mathbf{x}_2$  and  $\mathbf{u}$  constant:

$$C(\mathbf{x}_2, \mathbf{u}) = y(1) - y(0).$$

This is random (a function of  $\mathbf{x}_2$  and  $\mathbf{u}$ ) as both potential outcomes  $y(0)$  and  $y(1)$  are different across individuals.

In a sample, we cannot observe both outcomes from the same individual, we only observe the realized value

$$y = \begin{cases} y(0) & \text{if } x_1 = 0 \\ y(1) & \text{if } x_1 = 1. \end{cases}$$

As the causal effect varies across individuals and is not observable, it cannot be measured on the individual level. We therefore focus on aggregate causal effects, in particular what is known as the average causal effect.

**Definition 2.29.2** In the model (2.52) the *average causal effect* of  $x_1$  on  $y$  conditional on  $\mathbf{x}_2$  is

$$\begin{aligned} ACE(x_1, \mathbf{x}_2) &= \mathbb{E}(C(x_1, \mathbf{x}_2, \mathbf{u}) \mid x_1, \mathbf{x}_2) \\ &= \int_{\mathbb{R}^\ell} \nabla_1 h(x_1, \mathbf{x}_2, \mathbf{u}) f(\mathbf{u} \mid x_1, \mathbf{x}_2) d\mathbf{u} \end{aligned} \quad (2.54)$$

where  $f(\mathbf{u} \mid x_1, \mathbf{x}_2)$  is the conditional density of  $\mathbf{u}$  given  $x_1, \mathbf{x}_2$ .

We can think of the average causal effect  $ACE(x_1, \mathbf{x}_2)$  as the average effect in the general population. In our Jennifer & George schooling example given earlier, supposing that half of the population are Jennifer's and the other half George's, then the average causal effect of college is  $(10+4)/2 = \$7$  an hour. This is not the individual causal effect, it is the average of the causal effect across all individuals in the population. Given data on only educational attainment and wages, the ACE of \$7 is the best we can hope to learn.

When we conduct a regression analysis (that is, consider the regression of observed wages on educational attainment) we might hope that the regression reveals the average causal effect. Technically, that the regression derivative (the coefficient on education) equals the ACE. Is this the case? In other words, what is the relationship between the average causal effect  $ACE(x_1, \mathbf{x}_2)$  and the regression derivative  $\nabla_1 m(x_1, \mathbf{x}_2)$ ? Equation (2.52) implies that the CEF is

$$\begin{aligned} m(x_1, \mathbf{x}_2) &= \mathbb{E}(h(x_1, \mathbf{x}_2, \mathbf{u}) \mid x_1, \mathbf{x}_2) \\ &= \int_{\mathbb{R}^\ell} h(x_1, \mathbf{x}_2, \mathbf{u}) f(\mathbf{u} \mid x_1, \mathbf{x}_2) d\mathbf{u}, \end{aligned}$$

the average causal equation, averaged over the conditional distribution of the unobserved component  $\mathbf{u}$ .

Applying the marginal effect operator, the regression derivative is

$$\begin{aligned} \nabla_1 m(x_1, \mathbf{x}_2) &= \int_{\mathbb{R}^\ell} \nabla_1 h(x_1, \mathbf{x}_2, \mathbf{u}) f(\mathbf{u} \mid x_1, \mathbf{x}_2) d\mathbf{u} \\ &\quad + \int_{\mathbb{R}^\ell} h(x_1, \mathbf{x}_2, \mathbf{u}) \nabla_1 f(\mathbf{u} \mid x_1, \mathbf{x}_2) d\mathbf{u} \\ &= ACE(x_1, \mathbf{x}_2) + \int_{\mathbb{R}^\ell} h(x_1, \mathbf{x}_2, \mathbf{u}) \nabla_1 f(\mathbf{u} \mid x_1, \mathbf{x}_2) d\mathbf{u}. \end{aligned} \quad (2.55)$$

Equation (2.55) shows that in general, the regression derivative does not equal the average causal effect. The difference is the second term on the right-hand-side of (2.55). The regression derivative and ACE equal in the special case when this term equals zero, which occurs when  $\nabla_1 f(\mathbf{u} \mid x_1, \mathbf{x}_2) = 0$ , that is, when the conditional density of  $\mathbf{u}$  given  $(x_1, \mathbf{x}_2)$  does not depend on  $x_1$ . When this condition holds then the regression derivative equals the ACE, which means that regression analysis can be interpreted causally, in the sense that it uncovers average causal effects.

The condition is sufficiently important that it has a special name in the treatment effects literature.

**Definition 2.29.3** *Conditional Independence Assumption (CIA).* Conditional on  $\mathbf{x}_2$ , the random variables  $x_1$  and  $\mathbf{u}$  are statistically independent.

The CIA implies  $f(\mathbf{u} \mid x_1, \mathbf{x}_2) = f(\mathbf{u} \mid \mathbf{x}_2)$  does not depend on  $x_1$ , and thus  $\nabla_1 f(\mathbf{u} \mid x_1, \mathbf{x}_2) = 0$ . Thus the CIA implies that  $\nabla_1 m(x_1, \mathbf{x}_2) = ACE(x_1, \mathbf{x}_2)$ , the regression derivative equals the average causal effect.

**Theorem 2.29.1** *In the structural model (2.52), the Conditional Independence Assumption implies*

$$\nabla_1 m(x_1, \mathbf{x}_2) = ACE(x_1, \mathbf{x}_2)$$

*the regression derivative equals the average causal effect for  $x_1$  on  $y$  conditional on  $\mathbf{x}_2$ .*

This is a fascinating result. It shows that whenever the unobservable is independent of the treatment variable (after conditioning on appropriate regressors) the regression derivative equals the average causal effect. In this case, the CEF has causal economic meaning, giving strong justification to estimation of the CEF. Our derivation also shows the critical role of the CIA. If CIA fails, then the equality of the regression derivative and ACE fails.

This theorem is quite general. It applies equally to the treatment-effects model where  $x_1$  is binary or to more general settings where  $x_1$  is continuous.

It is also helpful to understand that the CIA is weaker than full independence of  $\mathbf{u}$  from the regressors  $(x_1, \mathbf{x}_2)$ . The CIA was introduced precisely as a minimal sufficient condition to obtain the desired result. Full independence implies the CIA and implies that each regression derivative equals that variable's average causal effect, but full independence is not necessary in order to causally interpret a subset of the regressors.

To illustrate, let's return to our education example involving a population with equal numbers of Jennifer's and George's. Recall that Jennifer earns \$10 as a high-school graduate and \$20 as a college graduate (and so has a causal effect of \$10) while George earns \$8 as a high-school graduate and \$12 as a college graduate (so has a causal effect of \$4). Given this information, the average causal effect of college is \$7, which is what we hope to learn from a regression analysis.

Now suppose that while in high school all students take an aptitude test, and if a student gets a high (H) score he or she goes to college with probability 3/4, and if a student gets a low (L) score he or she goes to college with probability 1/4. Suppose further that Jennifer's get an aptitude score of H with probability 3/4, while George's get a score of H with probability 1/4. Given this situation, 62.5% of Jennifer's will go to college<sup>13</sup>, while 37.5% of George's will go to college<sup>14</sup>.

An econometrician who randomly samples 32 individuals and collects data on educational attainment and wages will find the following wage distribution:

	\$8	\$10	\$12	\$20	Mean
High-School Graduate	10	6	0	0	\$8.75
College Graduate	0	0	6	10	\$17.00

Let *college* denote a dummy variable taking the value of 1 for a college graduate, otherwise 0. Thus the regression of wages on college attendance takes the form

$$\mathbb{E}(\text{wage} \mid \text{college}) = 8.25\text{college} + 8.75.$$

The coefficient on the college dummy, \$8.25, is the regression derivative, and the implied wage effect of college attendance. But \$8.25 overstates the average causal effect of \$7. The reason is because

<sup>13</sup>  $\Pr(\text{College} \mid \text{Jennifer}) = \Pr(\text{College} \mid H) \Pr(H \mid \text{Jennifer}) + \Pr(\text{College} \mid L) \Pr(L \mid \text{Jennifer}) = (3/4)^2 + (1/4)^2$

<sup>14</sup>  $\Pr(\text{College} \mid \text{George}) = \Pr(\text{College} \mid H) \Pr(H \mid \text{George}) + \Pr(\text{College} \mid L) \Pr(L \mid \text{George}) = (3/4)(1/4) + (1/4)(3/4)$

the CIA fails. In this model the unobservable  $\mathbf{u}$  is the individual's type (Jennifer or George) which is not independent of the regressor  $x_1$  (education), since Jennifer is more likely to go to college than George. Since Jennifer's causal effect is higher than George's, the regression derivative overstates the ACE. The coefficient \$8.25 is not the average benefit of college attendance, rather it is the observed difference in realized wages in a population whose decision to attend college is correlated with their individual causal effect. At the risk of repeating myself, in this example, \$8.25 is the true regression derivative, it is the difference in average wages between those with a college education and those without. It is not, however, the average causal effect of college education in the population.

This does not mean that it is impossible to estimate the ACE. The key is conditioning on the appropriate variables. The CIA says that we need to find a variable  $x_2$  such that conditional on  $x_2$ ,  $\mathbf{u}$  and  $x_1$  (type and education) are independent. In this example a variable which will achieve this is the aptitude test score. The decision to attend college was based on the test score, not on an individual's type. Thus educational attainment and type are independent once we condition on the test score.

This also alters the ACE. Notice that Definition 2.29.2 is a function of  $x_2$  (the test score). Among the students who receive a high test score,  $3/4$  are Jennifer's and  $1/4$  are George's. Thus the ACE for students with a score of H is  $(3/4) \times 10 + (1/4) \times 4 = \$8.50$ . Among the students who receive a low test score,  $1/4$  are Jennifer's and  $3/4$  are George's. Thus the ACE for students with a score of L is  $(1/4) \times 10 + (3/4) \times 4 = \$5.50$ . The ACE varies between these two observable groups (those with high test scores and those with low test scores). Again, we would hope to be able to learn the ACE from a regression analysis, this time from a regression of wages on education and test scores.

To see this in the wage distribution, suppose that the econometrician collects data on the aptitude test score as well as education and wages. Given a random sample of 32 individuals we would expect to find the following wage distribution:

	\$8	\$10	\$12	\$20	Mean
High-School Graduate + High Test Score	1	3	0	0	\$9.50
College Graduate + High Test Score	0	0	3	9	\$18.00
High-School Graduate + Low Test Score	9	3	0	0	\$8.50
College Graduate + Low Test Score	0	0	3	1	\$14.00

Define the dummy variable *highscore* which takes the value 1 for students who received a high test score, else zero. The regression of wages on college attendance and test scores (with interactions) takes the form

$$\mathbb{E}(\text{wage} \mid \text{college}, \text{highscore}) = 1.00\text{highscore} + 5.50\text{college} + 3.00\text{highscore} \times \text{college} + 8.50.$$

The coefficient on *college*, \$5.50, is the regression derivative of college attendance for those with low test scores, and the sum of this coefficient with the interaction coefficient, \$8.50, is the regression derivative for college attendance for those with high test scores. These equal the average causal effect as calculated above. Furthermore, since  $1/2$  of the population achieves a high test score and  $1/2$  achieve a low test score, the measured average causal effect in the entire population is \$7, which precisely equals the true value.

In this example, by conditioning on the aptitude test score, the average causal effect of education on wages can be learned from a regression analysis. What this shows is that by conditioning on the proper variables, it may be possible to achieve the CIA, in which case regression analysis measures average causal effects.

## 2.30 Expectation: Mathematical Details\*

We define the **mean** or **expectation**  $\mathbb{E}(y)$  of a random variable  $y$  as follows. If  $y$  is discrete on the set  $\{\tau_1, \tau_2, \dots\}$  then

$$\mathbb{E}(y) = \sum_{j=1}^{\infty} \tau_j \Pr(y = \tau_j),$$

and if  $y$  is continuous with density  $f$  then

$$\mathbb{E}(y) = \int_{-\infty}^{\infty} y f(y) dy.$$

We can unify these definitions by writing the expectation as the Lebesgue integral with respect to the distribution function  $F$

$$\mathbb{E}(y) = \int_{-\infty}^{\infty} y dF(y). \quad (2.56)$$

In the event that the integral (2.56) is not finite, separately evaluate the two integrals

$$I_1 = \int_0^{\infty} y dF(y) \quad (2.57)$$

$$I_2 = - \int_{-\infty}^0 y dF(y). \quad (2.58)$$

If  $I_1 = \infty$  and  $I_2 < \infty$  then it is typical to define  $\mathbb{E}(y) = \infty$ . If  $I_1 < \infty$  and  $I_2 = \infty$  then we define  $\mathbb{E}(y) = -\infty$ . However, if both  $I_1 = \infty$  and  $I_2 = \infty$  then  $\mathbb{E}(y)$  is undefined. If

$$\mathbb{E}|y| = \int_{-\infty}^{\infty} |y| dF(y) = I_1 + I_2 < \infty$$

then  $\mathbb{E}(y)$  exists and is finite. In this case it is common to say that the mean  $\mathbb{E}(y)$  is “well-defined”.

More generally,  $y$  has a finite  $r^{th}$  moment if

$$\mathbb{E}|y|^r < \infty. \quad (2.59)$$

By Liapunov’s Inequality (B.13), (2.59) implies  $\mathbb{E}|y|^s < \infty$  for all  $1 \leq s \leq r$ . Thus, for example, if the fourth moment is finite then the first, second and third moments are also finite, and so is the 3.9<sup>th</sup> moment.

It is common in econometric theory to assume that the variables, or certain transformations of the variables, have finite moments of a certain order. How should we interpret this assumption? How restrictive is it?

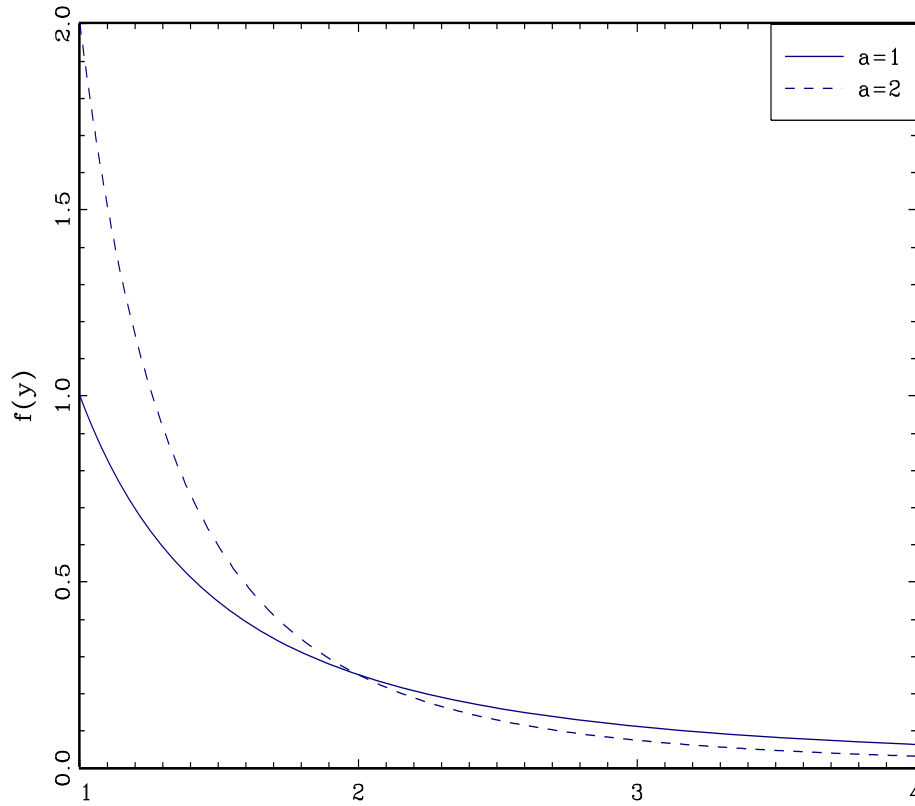
One way to visualize the importance is to consider the class of Pareto densities given by

$$f(y) = ay^{-a-1}, \quad y > 1.$$

The parameter  $a$  of the Pareto distribution indexes the rate of decay of the tail of the density. Larger  $a$  means that the tail declines to zero more quickly. See Figure 2.11 below where we plot the Pareto density for  $a = 1$  and  $a = 2$ . The parameter  $a$  also determines which moments are finite. We can calculate that

$$\mathbb{E}|y|^r = \begin{cases} a \int_1^{\infty} y^{r-a-1} dy = \frac{a}{a-r} & \text{if } r < a \\ \infty & \text{if } r \geq a. \end{cases}$$

This shows that if  $y$  is Pareto distributed with parameter  $a$ , then the  $r^{th}$  moment of  $y$  is finite if and only if  $r < a$ . Higher  $a$  means higher finite moments. Equivalently, the faster the tail of the density declines to zero, the more moments are finite.

Figure 2.11: Pareto Densities,  $a = 1$  and  $a = 2$ 

This connection between tail decay and finite moments is not limited to the Pareto distribution. We can make a similar analysis using a tail bound. Suppose that  $y$  has density  $f(y)$  which satisfies the bound  $f(y) \leq A|y|^{-a-1}$  for some  $A < \infty$  and  $a > 0$ . Since  $f(y)$  is bounded below a scale of a Pareto density, its tail behavior is similarly bounded. This means that for  $r < a$

$$\mathbb{E}|y|^r = \int_{-\infty}^{\infty} |y|^r f(y) dy \leq \int_{-1}^1 f(y) dy + 2A \int_1^{\infty} y^{r-a-1} dy \leq 1 + \frac{2A}{a-r} < \infty.$$

Thus if the tail of the density declines at the rate  $|y|^{-a-1}$  or faster, then  $y$  has finite moments up to (but not including)  $a$ . Broadly speaking, the restriction that  $y$  has a finite  $r^{\text{th}}$  moment means that the tail of  $y$ 's density declines to zero faster than  $y^{-r-1}$ . The faster decline of the tail means that the probability of observing an extreme value of  $y$  is a more rare event.

We complete this section by adding an alternative representation of expectation in terms of the distribution function.

**Theorem 2.30.1** *For any non-negative random variable  $y$*

$$\mathbb{E}(y) = \int_0^{\infty} \Pr(y > u) du$$

**Proof of Theorem 2.30.1:** Let  $F^*(x) = \Pr(y > x) = 1 - F(x)$ , where  $F(x)$  is the distribution function. By integration by parts

$$\mathbb{E}(y) = \int_0^{\infty} y dF(y) = - \int_0^{\infty} y dF^*(y) = -[yF^*(y)]_0^{\infty} + \int_0^{\infty} F^*(y) dy = \int_0^{\infty} \Pr(y > u) du$$

as stated. ■

### 2.31 Moment Generating and Characteristic Functions\*

For a random variable  $z$  with distribution  $F$  its **moment generating function** (MGF) is

$$M(t) = \mathbb{E}(\exp(tz)) = \int \exp(tz) dF(z). \quad (2.60)$$

This is also known as the Laplace transformation of the density of  $z$ . The MGF is a function of the argument  $t$ , and is an alternative representation of the distribution  $F$ . It is called the moment generating function since the  $r^{th}$  derivative evaluated at zero is the  $r^{th}$  uncentered moment. Indeed,

$$M^{(r)}(t) = \mathbb{E}\left(\frac{d^r}{dt^r} \exp(tz)\right) = \mathbb{E}(z^r \exp(tz))$$

and thus the  $r^{th}$  derivative at  $t = 0$  is

$$M^{(r)}(0) = \mathbb{E}(z^r).$$

A major limitation with the MGF is that it does not exist for many random variables. Essentially, existence of the integral (2.60) requires the tail of the density of  $z$  to decline exponentially. This excludes thick-tailed distributions such as the Pareto.

This limitation is removed if we consider the **characteristic function** (CF) of  $z$ , which is defined as

$$C(t) = \mathbb{E}(\exp(itz)) = \int \exp(itz) dF(z)$$

where  $i = \sqrt{-1}$ . Like the MGF, the CF is a function of its argument  $t$  and is a representation of the distribution function  $F$ . The CF is also known as the Fourier transformation of the density of  $z$ . Unlike the MGF, the CF exists for all random variables and all values of  $t$  since  $\exp(itz) = \cos(tz) + i \sin(tz)$  is bounded.

Similarly to the MGF, the  $r^{th}$  derivative of the characteristic function evaluated at zero takes the simple form

$$C^{(r)}(0) = i^r \mathbb{E}(z^r) \quad (2.61)$$

when such expectations exist. A further connection is that the  $r^{th}$  moment is finite if and only if  $C^{(r)}(t)$  is continuous at zero.

For random vectors  $\mathbf{z}$  with distribution  $F$  we define the multivariate MGF as

$$M(\mathbf{t}) = \mathbb{E}(\exp(\mathbf{t}'\mathbf{z})) = \int \exp(\mathbf{t}'\mathbf{z}) dF(\mathbf{z}) \quad (2.62)$$

when it exists. Similarly, we define the multivariate CF as

$$C(\mathbf{t}) = \mathbb{E}(\exp(i\mathbf{t}'\mathbf{z})) = \int \exp(i\mathbf{t}'\mathbf{z}) dF(\mathbf{z}).$$

### 2.32 Existence and Uniqueness of the Conditional Expectation\*

In Sections 2.3 and 2.6 we defined the conditional mean when the conditioning variables  $\mathbf{x}$  are discrete and when the variables  $(y, \mathbf{x})$  have a joint density. We have explored these cases because these are the situations where the conditional mean is easiest to describe and understand. However, the conditional mean exists quite generally without appealing to the properties of either discrete or continuous random variables.

To justify this claim we now present a deep result from probability theory. What it says is that the conditional mean exists for all joint distributions  $(y, \mathbf{x})$  for which  $y$  has a finite mean.

**Theorem 2.32.1 Existence of the Conditional Mean**

If  $\mathbb{E}|y| < \infty$  then there exists a function  $m(\mathbf{x})$  such that for all sets  $\mathcal{X}$  for which  $\Pr(\mathbf{x} \in \mathcal{X})$  is defined,

$$\mathbb{E}(1(\mathbf{x} \in \mathcal{X})y) = \mathbb{E}(1(\mathbf{x} \in \mathcal{X})m(\mathbf{x})). \quad (2.63)$$

The function  $m(\mathbf{x})$  is almost everywhere unique, in the sense that if  $h(\mathbf{x})$  satisfies (2.63), then there is a set  $S$  such that  $\Pr(S) = 1$  and  $m(\mathbf{x}) = h(\mathbf{x})$  for  $\mathbf{x} \in S$ . The function  $m(\mathbf{x})$  is called the **conditional mean** and is written  $m(\mathbf{x}) = \mathbb{E}(y | \mathbf{x})$ .

See, for example, Ash (1972), Theorem 6.3.3.

The conditional mean  $m(\mathbf{x})$  defined by (2.63) specializes to (2.7) when  $(y, \mathbf{x})$  have a joint density. The usefulness of definition (2.63) is that Theorem 2.32.1 shows that the conditional mean  $m(\mathbf{x})$  exists for all finite-mean distributions. This definition allows  $y$  to be discrete or continuous, for  $\mathbf{x}$  to be scalar or vector-valued, and for the components of  $\mathbf{x}$  to be discrete or continuously distributed.

You may have noticed that Theorem 2.32.1 applies only to sets  $\mathcal{X}$  for which  $\Pr(\mathbf{x} \in \mathcal{X})$  is defined. This is a technical issue – measurability – which we largely side-step in this textbook. Formal probability theory only applies to sets which are measurable – for which probabilities are defined, as it turns out that not all sets satisfy measurability. This is not a practical concern for econometrics, so we defer such distinctions for formal theoretical treatments.

## 2.33 Identification\*

A critical and important issue in structural econometric modeling is identification, meaning that a parameter is uniquely determined by the distribution of the observed variables. It is relatively straightforward in the context of the unconditional and conditional mean, but it is worthwhile to introduce and explore the concept at this point for clarity.

Let  $F$  denote the distribution of the observed data, for example the distribution of the pair  $(y, x)$ . Let  $\mathcal{F}$  be a collection of distributions  $F$ . Let  $\theta$  be a parameter of interest (for example, the mean  $\mathbb{E}(y)$ ).

**Definition 2.33.1** A parameter  $\theta \in \mathbb{R}$  is identified on  $\mathcal{F}$  if for all  $F \in \mathcal{F}$ , there is a uniquely determined value of  $\theta$ .

Equivalently,  $\theta$  is identified if we can write it as a mapping  $\theta = g(F)$  on the set  $\mathcal{F}$ . The restriction to the set  $\mathcal{F}$  is important. Most parameters are identified only on a strict subset of the space of all distributions.

Take, for example, the mean  $\mu = \mathbb{E}(y)$ . It is uniquely determined if  $\mathbb{E}|y| < \infty$ , so it is clear that  $\mu$  is identified for the set  $\mathcal{F} = \{F : \int_{-\infty}^{\infty} |y| dF(y) < \infty\}$ . However,  $\mu$  is also well defined when it is either positive or negative infinity. Hence, defining  $I_1$  and  $I_2$  as in (2.57) and (2.58), we can deduce that  $\mu$  is identified on the set  $\mathcal{F} = \{F : \{I_1 < \infty\} \cup \{I_2 < \infty\}\}$ .

Next, consider the conditional mean. Theorem 2.32.1 demonstrates that  $\mathbb{E}|y| < \infty$  is a sufficient condition for identification.



**Theorem 2.33.1 Identification of the Conditional Mean**

If  $\mathbb{E}|y| < \infty$ , the conditional mean  $m(\mathbf{x}) = \mathbb{E}(y | \mathbf{x})$  is identified almost everywhere.

It might seem as if identification is a general property for parameters, so long as we exclude degenerate cases. This is true for moments of observed data, but not necessarily for more complicated models. As a case in point, consider the context of censoring. Let  $y$  be a random variable with distribution  $F$ . Instead of observing  $y$ , we observe  $y^*$  defined by the censoring rule

$$y^* = \begin{cases} y & \text{if } y \leq \tau \\ \tau & \text{if } y > \tau \end{cases}.$$

That is,  $y^*$  is capped at the value  $\tau$ . A common example is income surveys, where income responses are “top-coded”, meaning that incomes above the top code  $\tau$  are recorded as the top code. The observed variable  $y^*$  has distribution

$$F^*(u) = \begin{cases} F(u) & \text{for } u \leq \tau \\ 1 & \text{for } u \geq \tau. \end{cases}$$

We are interested in features of the distribution  $F$  not the censored distribution  $F^*$ . For example, we are interested in the mean wage  $\mu = \mathbb{E}(y)$ . The difficulty is that we cannot calculate  $\mu$  from  $F^*$  except in the trivial case where there is no censoring  $\Pr(y \geq \tau) = 0$ . Thus the mean  $\mu$  is not generically identified from the censored distribution.

A typical solution to the identification problem is to assume a parametric distribution. For example, let  $\mathcal{F}$  be the set of normal distributions  $y \sim N(\mu, \sigma^2)$ . It is possible to show that the parameters  $(\mu, \sigma^2)$  are identified for all  $F \in \mathcal{F}$ . That is, if we know that the uncensored distribution is normal, we can uniquely determine the parameters from the censored distribution. This is often called **parametric identification** as identification is restricted to a parametric class of distributions. In modern econometrics this is generally viewed as a second-best solution, as identification has been achieved only through the use of an arbitrary and unverifiable parametric assumption.

A pessimistic conclusion might be that it is impossible to identify parameters of interest from censored data without parametric assumptions. Interestingly, this pessimism is unwarranted. It turns out that we can identify the quantiles  $q_\alpha$  of  $F$  for  $\alpha \leq \Pr(y \leq \tau)$ . For example, if 20% of the distribution is censored, we can identify all quantiles for  $\alpha \in (0, 0.8)$ . This is often called **nonparametric identification** as the parameters are identified without restriction to a parametric class.

What we have learned from this little exercise is that in the context of censored data, moments can only be parametrically identified, while non-censored quantiles are nonparametrically identified. Part of the message is that a study of identification can help focus attention on what can be learned from the data distributions available.

## 2.34 Technical Proofs\*

**Proof of Theorem 2.7.1:** For convenience, assume that the variables have a joint density  $f(y, \mathbf{x})$ . Since  $\mathbb{E}(y | \mathbf{x})$  is a function of the random vector  $\mathbf{x}$  only, to calculate its expectation we integrate with respect to the density  $f_{\mathbf{x}}(\mathbf{x})$  of  $\mathbf{x}$ , that is

$$\mathbb{E}(\mathbb{E}(y | \mathbf{x})) = \int_{\mathbb{R}^k} \mathbb{E}(y | \mathbf{x}) f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x}.$$

Substituting in (2.7) and noting that  $f_{y|\mathbf{x}}(y|\mathbf{x})f_{\mathbf{x}}(\mathbf{x}) = f(y, \mathbf{x})$ , we find that the above expression equals

$$\int_{\mathbb{R}^k} \left( \int_{\mathbb{R}} y f_{y|\mathbf{x}}(y|\mathbf{x}) dy \right) f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} = \int_{\mathbb{R}^k} \int_{\mathbb{R}} y f(y, \mathbf{x}) dy d\mathbf{x} = \mathbb{E}(y)$$

the unconditional mean of  $y$ . ■

**Proof of Theorem 2.7.2:** Again assume that the variables have a joint density. It is useful to observe that

$$f(y|\mathbf{x}_1, \mathbf{x}_2) f(\mathbf{x}_2|\mathbf{x}_1) = \frac{f(y, \mathbf{x}_1, \mathbf{x}_2)}{f(\mathbf{x}_1, \mathbf{x}_2)} \frac{f(\mathbf{x}_1, \mathbf{x}_2)}{f(\mathbf{x}_1)} = f(y, \mathbf{x}_2|\mathbf{x}_1), \quad (2.64)$$

the density of  $(y, \mathbf{x}_2)$  given  $\mathbf{x}_1$ . Here, we have abused notation and used a single symbol  $f$  to denote the various unconditional and conditional densities to reduce notational clutter.

Note that

$$\mathbb{E}(y | \mathbf{x}_1, \mathbf{x}_2) = \int_{\mathbb{R}} y f(y|\mathbf{x}_1, \mathbf{x}_2) dy. \quad (2.65)$$

Integrating (2.65) with respect to the conditional density of  $\mathbf{x}_2$  given  $\mathbf{x}_1$ , and applying (2.64) we find that

$$\begin{aligned} \mathbb{E}(\mathbb{E}(y | \mathbf{x}_1, \mathbf{x}_2) | \mathbf{x}_1) &= \int_{\mathbb{R}^{k_2}} \mathbb{E}(y | \mathbf{x}_1, \mathbf{x}_2) f(\mathbf{x}_2|\mathbf{x}_1) d\mathbf{x}_2 \\ &= \int_{\mathbb{R}^{k_2}} \left( \int_{\mathbb{R}} y f(y|\mathbf{x}_1, \mathbf{x}_2) dy \right) f(\mathbf{x}_2|\mathbf{x}_1) d\mathbf{x}_2 \\ &= \int_{\mathbb{R}^{k_2}} \int_{\mathbb{R}} y f(y|\mathbf{x}_1, \mathbf{x}_2) f(\mathbf{x}_2|\mathbf{x}_1) dy d\mathbf{x}_2 \\ &= \int_{\mathbb{R}^{k_2}} \int_{\mathbb{R}} y f(y, \mathbf{x}_2|\mathbf{x}_1) dy d\mathbf{x}_2 \\ &= \mathbb{E}(y | \mathbf{x}_1) \end{aligned}$$

as stated. ■

**Proof of Theorem 2.7.3:**

$$\mathbb{E}(g(\mathbf{x}) y | \mathbf{x}) = \int_{\mathbb{R}} g(\mathbf{x}) y f_{y|\mathbf{x}}(y|\mathbf{x}) dy = g(\mathbf{x}) \int_{\mathbb{R}} y f_{y|\mathbf{x}}(y|\mathbf{x}) dy = g(\mathbf{x}) \mathbb{E}(y | \mathbf{x})$$

This is (2.8). Equation (2.10) follows by applying the Simple Law of Iterated Expectations to (2.8). ■

**Proof of Theorem 2.8.1.** Applying Minkowski's Inequality (B.12) to  $e = y - m(\mathbf{x})$ ,

$$(\mathbb{E}|e|^r)^{1/r} = (\mathbb{E}|y - m(\mathbf{x})|^r)^{1/r} \leq (\mathbb{E}|y|^r)^{1/r} + (\mathbb{E}|m(\mathbf{x})|^r)^{1/r} < \infty,$$

where the two parts on the right-hand are finite since  $\mathbb{E}|y|^r < \infty$  by assumption and  $\mathbb{E}|m(\mathbf{x})|^r < \infty$  by the Conditional Expectation Inequality (B.7). The fact that  $(\mathbb{E}|e|^r)^{1/r} < \infty$  implies  $\mathbb{E}|e|^r < \infty$ . ■

**Proof of Theorem 2.10.2:** The assumption that  $\mathbb{E}(y^2) < \infty$  implies that all the conditional expectations below exist.

Using the law of iterated expectations  $\mathbb{E}(y | \mathbf{x}_1) = \mathbb{E}(\mathbb{E}(y | \mathbf{x}_1, \mathbf{x}_2) | \mathbf{x}_1)$  and the conditional Jensen's inequality (B.6),

$$(\mathbb{E}(y | \mathbf{x}_1))^2 = (\mathbb{E}(\mathbb{E}(y | \mathbf{x}_1, \mathbf{x}_2) | \mathbf{x}_1))^2 \leq \mathbb{E}((\mathbb{E}(y | \mathbf{x}_1, \mathbf{x}_2))^2 | \mathbf{x}_1).$$

Taking unconditional expectations, this implies

$$\mathbb{E} \left( (\mathbb{E}(y \mid \mathbf{x}_1))^2 \right) \leq \mathbb{E} \left( (\mathbb{E}(y \mid \mathbf{x}_1, \mathbf{x}_2))^2 \right).$$

Similarly,

$$(\mathbb{E}(y))^2 \leq \mathbb{E} \left( (\mathbb{E}(y \mid \mathbf{x}_1))^2 \right) \leq \mathbb{E} \left( (\mathbb{E}(y \mid \mathbf{x}_1, \mathbf{x}_2))^2 \right). \quad (2.66)$$

The variables  $y$ ,  $\mathbb{E}(y \mid \mathbf{x}_1)$  and  $\mathbb{E}(y \mid \mathbf{x}_1, \mathbf{x}_2)$  all have the same mean  $\mathbb{E}(y)$ , so the inequality (2.66) implies that the variances are ranked monotonically:

$$0 \leq \text{var}(\mathbb{E}(y \mid \mathbf{x}_1)) \leq \text{var}(\mathbb{E}(y \mid \mathbf{x}_1, \mathbf{x}_2)). \quad (2.67)$$

Define  $e = y - \mathbb{E}(y \mid \mathbf{x})$  and  $u = \mathbb{E}(y \mid \mathbf{x}) - \mu$  so that we have the decomposition

$$y - \mu = e + u.$$

Notice  $\mathbb{E}(e \mid \mathbf{x}) = 0$  and  $u$  is a function of  $\mathbf{x}$ . Thus by the Conditioning Theorem,  $\mathbb{E}(eu) = 0$  so  $e$  and  $u$  are uncorrelated. It follows that

$$\text{var}(y) = \text{var}(e) + \text{var}(u) = \text{var}(y - \mathbb{E}(y \mid \mathbf{x})) + \text{var}(\mathbb{E}(y \mid \mathbf{x})). \quad (2.68)$$

The monotonicity of the variances of the conditional mean (2.67) applied to the variance decomposition (2.68) implies the reverse monotonicity of the variances of the differences, completing the proof. ■

**Proof of Theorem 2.18.1.** For part 1, by the Expectation Inequality (B.8), (A.24) and Assumption 2.18.1,

$$\|\mathbb{E}(\mathbf{x}\mathbf{x}')\| \leq \mathbb{E}\|\mathbf{x}\mathbf{x}'\| = \mathbb{E}(\|\mathbf{x}\|^2) < \infty.$$

Similarly, using the Expectation Inequality (B.8), the Cauchy-Schwarz Inequality (B.10) and Assumption 2.18.1,

$$\|\mathbb{E}(\mathbf{x}y)\| \leq \mathbb{E}\|\mathbf{x}y\| \leq \left( \mathbb{E}(\|\mathbf{x}\|^2) \right)^{1/2} (\mathbb{E}(y^2))^{1/2} < \infty.$$

Thus the moments  $\mathbb{E}(\mathbf{x}y)$  and  $\mathbb{E}(\mathbf{x}\mathbf{x}')$  are finite and well defined.

For part 2, the coefficient  $\beta = (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}y)$  is well defined since  $(\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1}$  exists under Assumption 2.18.1.

Part 3 follows from Definition 2.18.1 and part 2.

For part 4, first note that

$$\begin{aligned} \mathbb{E}(e^2) &= \mathbb{E} \left( (y - \mathbf{x}'\beta)^2 \right) \\ &= \mathbb{E}(y^2) - 2\mathbb{E}(y\mathbf{x}')\beta + \beta'\mathbb{E}(\mathbf{x}\mathbf{x}')\beta \\ &= \mathbb{E}(y^2) - 2\mathbb{E}(y\mathbf{x}')(\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1}\mathbb{E}(\mathbf{x}y) \\ &\leq \mathbb{E}(y^2) \\ &< \infty. \end{aligned}$$

The first inequality holds because  $\mathbb{E}(y\mathbf{x}')(\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1}\mathbb{E}(\mathbf{x}y)$  is a quadratic form and therefore necessarily non-negative. Second, by the Expectation Inequality (B.8), the Cauchy-Schwarz Inequality (B.10) and Assumption 2.18.1,

$$\|\mathbb{E}(\mathbf{x}e)\| \leq \mathbb{E}\|\mathbf{x}e\| = \left( \mathbb{E}(\|\mathbf{x}\|^2) \right)^{1/2} (\mathbb{E}(e^2))^{1/2} < \infty.$$

It follows that the expectation  $\mathbb{E}(\mathbf{x}e)$  is finite, and is zero by the calculation (2.28).

For part 6, Applying Minkowski's Inequality (B.12) to  $e = y - \mathbf{x}'\boldsymbol{\beta}$ ,

$$\begin{aligned} (\mathbb{E} |e|^r)^{1/r} &= (\mathbb{E} |y - \mathbf{x}'\boldsymbol{\beta}|^r)^{1/r} \\ &\leq (\mathbb{E} |y|^r)^{1/r} + (\mathbb{E} |\mathbf{x}'\boldsymbol{\beta}|^r)^{1/r} \\ &\leq (\mathbb{E} |y|^r)^{1/r} + (\mathbb{E} \|\mathbf{x}\|^r)^{1/r} \|\boldsymbol{\beta}\| \\ &< \infty, \end{aligned}$$

the final inequality by assumption. ■

## Exercises

**Exercise 2.1** Find  $\mathbb{E}(\mathbb{E}(\mathbb{E}(y \mid \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) \mid \mathbf{x}_1, \mathbf{x}_2) \mid \mathbf{x}_1)$ .

**Exercise 2.2** If  $\mathbb{E}(y \mid x) = a + bx$ , find  $\mathbb{E}(yx)$  as a function of moments of  $x$ .

**Exercise 2.3** Prove Theorem 2.8.1.4 using the law of iterated expectations.

**Exercise 2.4** Suppose that the random variables  $y$  and  $x$  only take the values 0 and 1, and have the following joint probability distribution

	$x = 0$	$x = 1$
$y = 0$	.1	.2
$y = 1$	.4	.3

Find  $\mathbb{E}(y \mid x)$ ,  $\mathbb{E}(y^2 \mid x)$  and  $\text{var}(y \mid x)$  for  $x = 0$  and  $x = 1$ .

**Exercise 2.5** Show that  $\sigma^2(\mathbf{x})$  is the best predictor of  $e^2$  given  $\mathbf{x}$ :

- (a) Write down the mean-squared error of a predictor  $h(\mathbf{x})$  for  $e^2$ .
- (b) What does it mean to be predicting  $e^2$ ?
- (c) Show that  $\sigma^2(\mathbf{x})$  minimizes the mean-squared error and is thus the best predictor.

**Exercise 2.6** Use  $y = m(\mathbf{x}) + e$  to show that

$$\text{var}(y) = \text{var}(m(\mathbf{x})) + \sigma^2$$

**Exercise 2.7** Show that the conditional variance can be written as

$$\sigma^2(\mathbf{x}) = \mathbb{E}(y^2 \mid \mathbf{x}) - (\mathbb{E}(y \mid \mathbf{x}))^2.$$

**Exercise 2.8** Suppose that  $y$  is discrete-valued, taking values only on the non-negative integers, and the conditional distribution of  $y$  given  $\mathbf{x}$  is Poisson:

$$\Pr(y = j \mid \mathbf{x}) = \frac{\exp(-\mathbf{x}'\boldsymbol{\beta})(\mathbf{x}'\boldsymbol{\beta})^j}{j!}, \quad j = 0, 1, 2, \dots$$

Compute  $\mathbb{E}(y \mid \mathbf{x})$  and  $\text{var}(y \mid \mathbf{x})$ . Does this justify a linear regression model of the form  $y = \mathbf{x}'\boldsymbol{\beta} + e$ ?

Hint: If  $\Pr(y = j) = \frac{\exp(-\lambda)\lambda^j}{j!}$ , then  $\mathbb{E}(y) = \lambda$  and  $\text{var}(y) = \lambda$ .

**Exercise 2.9** Suppose you have two regressors:  $x_1$  is binary (takes values 0 and 1) and  $x_2$  is categorical with 3 categories ( $A, B, C$ ). Write  $\mathbb{E}(y \mid x_1, x_2)$  as a linear regression.

**Exercise 2.10** True or False. If  $y = x\beta + e$ ,  $x \in \mathbb{R}$ , and  $\mathbb{E}(e \mid x) = 0$ , then  $\mathbb{E}(x^2e) = 0$ .

**Exercise 2.11** True or False. If  $y = x\beta + e$ ,  $x \in \mathbb{R}$ , and  $\mathbb{E}(xe) = 0$ , then  $\mathbb{E}(x^2e) = 0$ .

**Exercise 2.12** True or False. If  $y = \mathbf{x}'\boldsymbol{\beta} + e$  and  $\mathbb{E}(e \mid \mathbf{x}) = 0$ , then  $e$  is independent of  $\mathbf{x}$ .

**Exercise 2.13** True or False. If  $y = \mathbf{x}'\boldsymbol{\beta} + e$  and  $\mathbb{E}(xe) = \mathbf{0}$ , then  $\mathbb{E}(e \mid \mathbf{x}) = 0$ .

**Exercise 2.14** True or False. If  $y = \mathbf{x}'\boldsymbol{\beta} + e$ ,  $\mathbb{E}(e | \mathbf{x}) = 0$ , and  $\mathbb{E}(e^2 | \mathbf{x}) = \sigma^2$ , a constant, then  $e$  is independent of  $\mathbf{x}$ .

**Exercise 2.15** Consider the intercept-only model  $y = \alpha + e$  defined as the best linear predictor. Show that  $\alpha = \mathbb{E}(y)$ .

**Exercise 2.16** Let  $x$  and  $y$  have the joint density  $f(x, y) = \frac{3}{2}(x^2 + y^2)$  on  $0 \leq x \leq 1$ ,  $0 \leq y \leq 1$ . Compute the coefficients of the best linear predictor  $y = \alpha + \beta x + e$ . Compute the conditional mean  $m(x) = \mathbb{E}(y | x)$ . Are the best linear predictor and conditional mean different?

**Exercise 2.17** Let  $x$  be a random variable with  $\mu = \mathbb{E}(x)$  and  $\sigma^2 = \text{var}(x)$ . Define

$$g(x | \mu, \sigma^2) = \left( \frac{x - \mu}{(x - \mu)^2 - \sigma^2} \right).$$

Show that  $\mathbb{E}g(x | m, s) = 0$  if and only if  $m = \mu$  and  $s = \sigma^2$ .

**Exercise 2.18** Suppose that

$$\mathbf{x} = \begin{pmatrix} 1 \\ x_2 \\ x_3 \end{pmatrix}$$

and  $x_3 = \alpha_1 + \alpha_2 x_2$  is a linear function of  $x_2$ .

- Show that  $\mathbf{Q}_{\mathbf{x}\mathbf{x}} = \mathbb{E}(\mathbf{x}\mathbf{x}')$  is not invertible.
- Use a linear transformation of  $\mathbf{x}$  to find an expression for the best linear predictor of  $y$  given  $\mathbf{x}$ . (Be explicit, do not just use the generalized inverse formula.)

**Exercise 2.19** Show (2.46)-(2.47), namely that for

$$d(\boldsymbol{\beta}) = \mathbb{E}(m(\mathbf{x}) - \mathbf{x}'\boldsymbol{\beta})^2$$

then

$$\begin{aligned} \boldsymbol{\beta} &= \underset{\mathbf{b} \in \mathbb{R}^k}{\text{argmin}} d(\mathbf{b}) \\ &= (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}m(\mathbf{x})) \\ &= (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}y). \end{aligned}$$

Hint: To show  $\mathbb{E}(\mathbf{x}m(\mathbf{x})) = \mathbb{E}(\mathbf{x}y)$  use the law of iterated expectations.

**Exercise 2.20** Verify that (2.63) holds with  $m(\mathbf{x})$  defined in (2.7) when  $(y, \mathbf{x})$  have a joint density  $f(y, \mathbf{x})$ .

**Exercise 2.21** Consider the short and long projections

$$y = x\gamma_1 + e$$

$$y = x\beta_1 + x^2\beta_2 + u$$

- Under what condition does  $\gamma_1 = \beta_1$ ?
- Now suppose the long projection is

$$y = x\theta_1 + x^3\theta_2 + v$$

Is there a similar condition under which  $\gamma_1 = \theta_1$ ?

**Exercise 2.22** Take the homoskedastic model

$$\begin{aligned}y &= \mathbf{x}'_1\boldsymbol{\beta}_1 + \mathbf{x}'_2\boldsymbol{\beta}_2 + e \\ \mathbb{E}(e \mid \mathbf{x}_1, \mathbf{x}_2) &= 0 \\ \mathbb{E}(e^2 \mid \mathbf{x}_1, \mathbf{x}_2) &= \sigma^2 \\ \mathbb{E}(\mathbf{x}_2 \mid \mathbf{x}_1) &= \boldsymbol{\Gamma}\mathbf{x}_1 \\ \boldsymbol{\Gamma} &\neq 0\end{aligned}$$

Suppose the parameter  $\boldsymbol{\beta}_1$  is of interest. We know that the exclusion of  $\mathbf{x}_2$  creates omitted variable bias in the projection coefficient on  $\mathbf{x}_2$ . It also changes the equation error. Our question is: what is the effect on the homoskedasticity property of the induced equation error? Does the exclusion of  $\mathbf{x}_2$  induce heteroskedasticity or not? Be specific.

## Chapter 3

# The Algebra of Least Squares

### 3.1 Introduction

In this chapter we introduce the popular least-squares estimator. Most of the discussion will be algebraic, with questions of distribution and inference deferred to later chapters.

### 3.2 Samples

In Section 2.18 we derived and discussed the best linear predictor of  $y$  given  $\mathbf{x}$  for a pair of random variables  $(y, \mathbf{x}) \in \mathbb{R} \times \mathbb{R}^k$ , and called this the linear projection model. We are now interested in **estimating** the parameters of this model, in particular the projection coefficient

$$\beta = (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}y). \quad (3.1)$$

We can estimate  $\beta$  from observational data which includes joint measurements on the variables  $(y, \mathbf{x})$ . For example, supposing we are interested in estimating a wage equation, we would use a dataset with observations on wages (or weekly earnings), education, experience (or age), and demographic characteristics (gender, race, location). One possible dataset is the Current Population Survey (CPS), a survey of U.S. households which includes questions on employment, income, education, and demographic characteristics.

Notationally we wish to distinguish observations from the underlying random variables. The convention in econometrics is to denote observations by appending a subscript  $i$  which runs from 1 to  $n$ , thus the  $i^{\text{th}}$  observation is  $(y_i, \mathbf{x}_i)$ , and  $n$  denotes the sample size. The dataset is then  $\{(y_i, \mathbf{x}_i); i = 1, \dots, n\}$ . We call this the **sample** or the **observations**.

From the viewpoint of empirical analysis, a dataset is an array of numbers often organized as a table, where the columns of the table correspond to distinct variables and the rows correspond to distinct observations. For empirical analysis, the dataset and observations are fixed in the sense that they are numbers presented to the researcher. For statistical analysis we need to view the dataset as random, or more precisely as a realization of a random process.

In order for the coefficient  $\beta$  defined in (3.1) to make sense as defined, the expectations over the random variables  $(\mathbf{x}, y)$  need to be common across the observations. The most elegant approach to ensure this is to assume that the observations are draws from an identical underlying population  $F$ . This is the standard assumption that the observations are identically distributed:

**Assumption 3.2.1** *The observations  $\{(y_1, \mathbf{x}_1), \dots, (y_i, \mathbf{x}_i), \dots, (y_n, \mathbf{x}_n)\}$  are identically distributed; they are draws from a common distribution  $F$ .*



This assumption does not need to be viewed as literally true, rather it is a useful modeling device so that parameters such as  $\beta$  are well defined. This assumption should be interpreted as how we view an observation *a priori*, before we actually observe it. If I tell you that we have a sample with  $n = 59$  observations set in no particular order, then it makes sense to view two observations, say 17 and 58, as draws from the same distribution. We have no reason to expect anything special about either observation.

In econometric theory, we refer to the underlying common distribution  $F$  as the **population**. Some authors prefer the label the **data-generating-process** (DGP). You can think of it as a theoretical concept or an infinitely-large potential population. In contrast we refer to the observations available to us  $\{(y_i, \mathbf{x}_i); i = 1, \dots, n\}$  as the **sample** or **dataset**. In some contexts the dataset consists of all potential observations, for example administrative tax records may contain every single taxpayer in a political unit. Even in this case we view the observations as if they are random draws from an underlying infinitely-large population, as this will allow us to apply the tools of statistical theory.

The linear projection model applies to the random observations  $(y_i, \mathbf{x}_i)$ . This means that the probability model for the observations is the same as that described in Section 2.18. We can write the model as

$$y_i = \mathbf{x}_i' \beta + e_i \quad (3.2)$$

where the linear projection coefficient  $\beta$  is defined as

$$\beta = \underset{b \in \mathbb{R}^k}{\operatorname{argmin}} S(b), \quad (3.3)$$

the minimizer of the expected squared error

$$S(\beta) = \mathbb{E} \left( (y_i - \mathbf{x}_i' \beta)^2 \right), \quad (3.4)$$

and has the explicit solution

$$\beta = \left( \mathbb{E} (\mathbf{x}_i \mathbf{x}_i') \right)^{-1} \mathbb{E} (\mathbf{x}_i y_i). \quad (3.5)$$

### 3.3 Moment Estimators

We want to estimate the coefficient  $\beta$  defined in (3.5) from the sample of observations. Notice that  $\beta$  is written as a function of certain population expectations. In this context an appropriate estimator is the same function of the sample moments. Let's explain this in detail.

To start, suppose that we are interested in the population mean  $\mu$  of a random variable  $y_i$  with distribution function  $F$

$$\mu = \mathbb{E}(y_i) = \int_{-\infty}^{\infty} y dF(y). \quad (3.6)$$

The mean  $\mu$  is a function of the distribution  $F$  as written in (3.6). To estimate  $\mu$  given a sample  $\{y_1, \dots, y_n\}$  a natural estimator is the sample mean

$$\hat{\mu} = \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i.$$

Notice that we have written this using two pieces of notation. The notation  $\bar{y}$  with the bar on top is conventional for a sample mean. The notation  $\hat{\mu}$  with the hat “ $\wedge$ ” is conventional in econometrics to denote an estimator of the parameter  $\mu$ . In this case, the sample mean of  $y_i$  is the estimator of  $\mu$ , so  $\hat{\mu}$  and  $\bar{y}$  are the same. The sample mean  $\bar{y}$  can be viewed as the natural analog of the population mean (3.6) because  $\bar{y}$  equals the expectation (3.6) with respect to the empirical distribution – the discrete distribution which puts weight  $1/n$  on each observation  $y_i$ . There are many other justifications for  $\bar{y}$  as an estimator for  $\mu$ , we will defer these discussions for now. Suffice it to say

that it is the conventional estimator in the lack of other information about  $\mu$  or the distribution of  $y_i$ .

Now suppose that we are interested in a set of population means of possibly non-linear functions of a random vector  $\mathbf{y}$ , say  $\boldsymbol{\mu} = \mathbb{E}(\mathbf{h}(\mathbf{y}_i))$ . For example, we may be interested in the first two moments of  $y_i$ ,  $\mathbb{E}(y_i)$  and  $\mathbb{E}(y_i^2)$ . In this case the natural estimator is the vector of sample means,

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{h}(\mathbf{y}_i).$$

For example,  $\hat{\mu}_1 = \frac{1}{n} \sum_{i=1}^n y_i$  and  $\hat{\mu}_2 = \frac{1}{n} \sum_{i=1}^n y_i^2$ . This is not really a substantive change. We call  $\hat{\boldsymbol{\mu}}$  the **moment estimator** for  $\boldsymbol{\mu}$ .

Now suppose that we are interested in a nonlinear function of a set of moments. For example, consider the variance of  $y$

$$\sigma^2 = \text{var}(y_i) = \mathbb{E}(y_i^2) - (\mathbb{E}(y_i))^2.$$

In general, many parameters of interest, say  $\boldsymbol{\beta}$ , can be written as a function of moments of  $\mathbf{y}$ . Notationally,

$$\begin{aligned} \boldsymbol{\beta} &= \mathbf{g}(\boldsymbol{\mu}) \\ \boldsymbol{\mu} &= \mathbb{E}(\mathbf{h}(\mathbf{y}_i)). \end{aligned}$$

Here,  $\mathbf{y}_i$  are the random variables,  $\mathbf{h}(\mathbf{y}_i)$  are functions (transformations) of the random variables, and  $\boldsymbol{\mu}$  is the mean (expectation) of these functions.  $\boldsymbol{\beta}$  is the parameter of interest, and is the (nonlinear) function  $\mathbf{g}(\cdot)$  of these means.

In this context a natural estimator of  $\boldsymbol{\beta}$  is obtained by replacing  $\boldsymbol{\mu}$  with  $\hat{\boldsymbol{\mu}}$ .

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \mathbf{g}(\hat{\boldsymbol{\mu}}) \\ \hat{\boldsymbol{\mu}} &= \frac{1}{n} \sum_{i=1}^n \mathbf{h}(\mathbf{y}_i) \end{aligned}$$

The estimator  $\hat{\boldsymbol{\beta}}$  is often called a “plug-in” estimator, and sometimes a “substitution” estimator. We typically call  $\hat{\boldsymbol{\beta}}$  a moment, or moment-based, estimator of  $\boldsymbol{\beta}$ , since it is a natural extension of the moment estimator  $\hat{\boldsymbol{\mu}}$ .

Take the example of the variance  $\sigma^2 = \text{var}(y_i)$ . Its moment estimator is

$$\hat{\sigma}^2 = \hat{\mu}_2 - \hat{\mu}_1^2 = \frac{1}{n} \sum_{i=1}^n y_i^2 - \left( \frac{1}{n} \sum_{i=1}^n y_i \right)^2.$$

This is not the only possible estimator for  $\sigma^2$  (there is the well-known bias-corrected version appropriate for independent observations) but it is a straightforward and simple choice.

### 3.4 Least Squares Estimator

The linear projection coefficient  $\boldsymbol{\beta}$  is defined in (3.3) as the minimizer of the expected squared error  $S(\boldsymbol{\beta})$  defined in (3.4). For given  $\boldsymbol{\beta}$ , the expected squared error is the expectation of the squared error  $(y_i - \mathbf{x}_i' \boldsymbol{\beta})^2$ . The moment estimator of  $S(\boldsymbol{\beta})$  is the sample average:

$$\begin{aligned} \hat{S}(\boldsymbol{\beta}) &= \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i' \boldsymbol{\beta})^2 \\ &= \frac{1}{n} SSE(\boldsymbol{\beta}) \end{aligned} \tag{3.7}$$

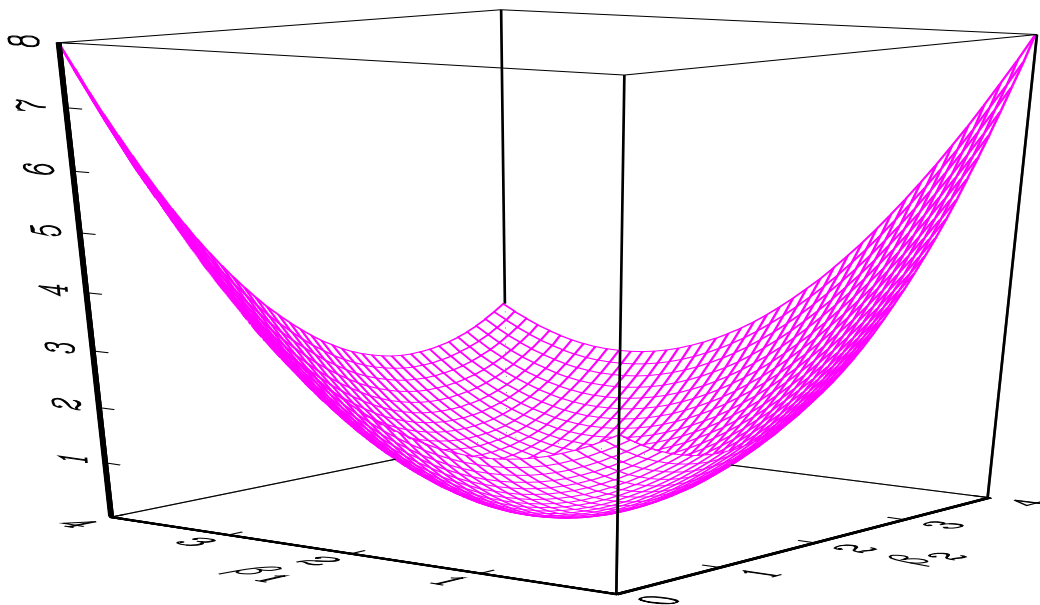


Figure 3.1: Sum-of-Squared Errors Function

where

$$SSE(\beta) = \sum_{i=1}^n (y_i - \mathbf{x}_i' \beta)^2 \quad (3.8)$$

is called the **sum-of-squared-errors** function.

Since  $\hat{S}(\beta)$  is a sample average, we can interpret it as an estimator of the expected squared error  $S(\beta)$ . Examining  $\hat{S}(\beta)$  as a function of  $\beta$  is informative about how  $S(\beta)$  varies with  $\beta$ . Since the projection coefficient minimizes  $S(\beta)$ , an analog estimator minimizes (3.7):

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^k}{\operatorname{argmin}} \hat{S}(\beta).$$

Alternatively, as  $\hat{S}(\beta)$  is a scale multiple of  $SSE(\beta)$ , we may equivalently define  $\hat{\beta}$  as the minimizer of  $SSE_n(\beta)$ . Hence  $\hat{\beta}$  is commonly called the **least-squares (LS)** estimator of  $\beta$ . (The estimator is also commonly referred to as the **ordinary least-squares OLS** estimator. For the origin of this label see the historical discussion on Adrien-Marie Legendre below.) Here, as is common in econometrics, we put a hat “^” over the parameter  $\beta$  to indicate that  $\hat{\beta}$  is a sample estimate of  $\beta$ . This is a helpful convention. Just by seeing the symbol  $\hat{\beta}$  we can immediately interpret it as an estimator (because of the hat) of the parameter  $\beta$ . Sometimes when we want to be explicit about the estimation method, we will write  $\hat{\beta}_{\text{ols}}$  to signify that it is the OLS estimator. It is also common to see the notation  $\hat{\beta}_n$ , where the subscript “ $n$ ” indicates that the estimator depends on the sample size  $n$ .

It is important to understand the distinction between population parameters such as  $\beta$  and sample estimates such as  $\hat{\beta}$ . The population parameter  $\beta$  is a non-random feature of the population while the sample estimate  $\hat{\beta}$  is a random feature of a random sample.  $\beta$  is fixed, while  $\hat{\beta}$  varies across samples.

To visualize the quadratic function  $\hat{S}(\beta)$ , Figure 3.1 displays an example sum-of-squared errors function  $SSE(\beta)$  for the case  $k = 2$ . The least-squares estimator  $\hat{\beta}$  is the pair  $(\hat{\beta}_1, \hat{\beta}_2)$  which minimize this function.

### 3.5 Solving for Least Squares with One Regressor

For simplicity, we start by considering the case  $k = 1$  so that the coefficient  $\beta$  is a scalar. Then the sum of squared errors is a simple quadratic

$$\begin{aligned} SSE(\beta) &= \sum_{i=1}^n (y_i - x_i\beta)^2 \\ &= \left( \sum_{i=1}^n y_i^2 \right) - 2\beta \left( \sum_{i=1}^n x_i y_i \right) + \beta^2 \left( \sum_{i=1}^n x_i^2 \right). \end{aligned}$$

The OLS estimator  $\hat{\beta}$  minimizes this function. From elementary algebra we know that the minimizer of the quadratic function  $a - 2bx + cx^2$  is  $x = b/c$ . Thus the minimizer of  $SSE(\beta)$  is

$$\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}. \quad (3.9)$$

The intercept-only model is the special case  $x_i = 1$ . In this case we find

$$\hat{\beta} = \frac{\sum_{i=1}^n 1 y_i}{\sum_{i=1}^n 1^2} = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}, \quad (3.10)$$

the sample mean of  $y_i$ . Here, as is common, we put a bar “ $-$ ” over  $y$  to indicate that the quantity is a sample mean. This calculation shows that the OLS estimator in the intercept-only model is the sample mean.

### 3.6 Solving for Least Squares with Multiple Regressors

We now consider the case with  $k \geq 1$  so that the coefficient  $\beta$  is a vector.

To solve for  $\hat{\beta}$ , expand the SSE function to find

$$SSE(\beta) = \sum_{i=1}^n y_i^2 - 2\beta' \sum_{i=1}^n \mathbf{x}_i y_i + \beta' \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \beta.$$

This is a quadratic expression in the vector argument  $\beta$ . The first-order-condition for minimization of  $SSE(\beta)$  is

$$0 = \frac{\partial}{\partial \beta} SSE(\hat{\beta}) = -2 \sum_{i=1}^n \mathbf{x}_i y_i + 2 \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \hat{\beta}. \quad (3.11)$$

We have written this using a single expression, but it is actually a system of  $k$  equations with  $k$  unknowns (the elements of  $\hat{\beta}$ ).

The solution for  $\hat{\beta}$  may be found by solving the system of  $k$  equations in (3.11). We can write this solution compactly using matrix algebra. Inverting the  $k \times k$  matrix  $\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'$  we find an explicit formula for the least-squares estimator

$$\hat{\beta} = \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left( \sum_{i=1}^n \mathbf{x}_i y_i \right). \quad (3.12)$$

This is the natural estimator of the best linear projection coefficient  $\beta$  defined in (3.3), and can also be called the linear projection estimator.

We see that (3.12) simplifies to the expression (3.9) when  $k = 1$ . The expression (3.12) is a notationally simple generalization but requires a careful attention to vector and matrix manipulations.

Alternatively, equation (3.5) writes the projection coefficient  $\beta$  as an explicit function of the population moments  $Q_{xy}$  and  $Q_{xx}$ . Their moment estimators are the sample moments

$$\begin{aligned}\widehat{Q}_{xy} &= \frac{1}{n} \sum_{i=1}^n x_i y_i \\ \widehat{Q}_{xx} &= \frac{1}{n} \sum_{i=1}^n x_i x'_i.\end{aligned}$$

The moment estimator of  $\beta$  replaces the population moments in (3.5) with the sample moments:

$$\begin{aligned}\widehat{\beta} &= \widehat{Q}_{xx}^{-1} \widehat{Q}_{xy} \\ &= \left( \frac{1}{n} \sum_{i=1}^n x_i x'_i \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n x_i y_i \right) \\ &= \left( \sum_{i=1}^n x_i x'_i \right)^{-1} \left( \sum_{i=1}^n x_i y_i \right)\end{aligned}$$

which is identical with (3.12).

### Least Squares Estimation

**Definition 3.6.1** *The least-squares estimator  $\widehat{\beta}$  is*

$$\widehat{\beta} = \underset{\beta \in \mathbb{R}^k}{\operatorname{argmin}} \widehat{S}(\beta)$$

where

$$\widehat{S}(\beta) = \frac{1}{n} \sum_{i=1}^n (y_i - x'_i \beta)^2$$

and has the solution

$$\widehat{\beta} = \left( \sum_{i=1}^n x_i x'_i \right)^{-1} \left( \sum_{i=1}^n x_i y_i \right).$$

### Adrien-Marie Legendre

The method of least-squares was first published in 1805 by the French mathematician Adrien-Marie Legendre (1752-1833). Legendre proposed least-squares as a solution to the algebraic problem of solving a system of equations when the number of equations exceeded the number of unknowns. This was a vexing and common problem in astronomical measurement. As viewed by Legendre, (3.2) is a set of  $n$  equations with  $k$  unknowns. As the equations cannot be solved exactly, Legendre's goal was to select  $\beta$  to make the set of errors as small as possible. He proposed the sum of squared error criterion, and derived the algebraic solution presented above. As he noted, the first-order conditions (3.11) is a system of  $k$  equations with  $k$  unknowns, which can be solved by "ordinary" methods. Hence the method became known as **Ordinary Least Squares** and to this day we still use the abbreviation OLS to refer to Legendre's estimation method.

## 3.7 Illustration

We illustrate the least-squares estimator in practice with the data set used to calculate the estimates reported in Chapter 2. This is the March 2009 Current Population Survey, which has extensive information on the U.S. population. This data set is described in more detail in Section 3.19. For this illustration, we use the sub-sample of married (spouse present) black female wage earners with 12 years potential work experience. This sub-sample has 20 observations. Let  $y_i$  be log wages and  $\mathbf{x}_i$  be years of education and an intercept. Then

$$\sum_{i=1}^n \mathbf{x}_i y_i = \begin{pmatrix} 995.86 \\ 62.64 \end{pmatrix},$$

$$\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' = \begin{pmatrix} 5010 & 314 \\ 314 & 20 \end{pmatrix},$$

and

$$\left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} = \begin{pmatrix} 0.0125 & -0.196 \\ -0.196 & 3.124 \end{pmatrix}.$$

Thus

$$\begin{aligned} \hat{\beta} &= \begin{pmatrix} 0.0125 & -0.196 \\ -0.196 & 3.124 \end{pmatrix} \begin{pmatrix} 995.86 \\ 62.64 \end{pmatrix} \\ &= \begin{pmatrix} 0.155 \\ 0.698 \end{pmatrix}. \end{aligned} \tag{3.13}$$

We often write the estimated equation using the format

$$\log(\widehat{Wage}) = 0.155 \text{ education} + 0.698. \tag{3.14}$$

An interpretation of the estimated equation is that each year of education is associated with a 16% increase in mean wages.

Equation (3.14) is called a **bivariate regression** as there are only two variables. A **multivariate regression** has two or more regressors, and allows a more detailed investigation. Let's take

an example similar to (3.14) but include all levels of experience. This time, we use the sub-sample of single (never married) Asian men, which has 268 observations. Including as regressors years of potential work experience (*experience*) and its square ( $experience^2/100$ ) (we divide by 100 to simplify reporting), we obtain the estimates

$$\widehat{\log(Wage)} = 0.143 \text{ education} + 0.036 \text{ experience} - 0.071 \text{ experience}^2/100 + 0.575. \quad (3.15)$$

These estimates suggest a 14% increase in mean wages per year of education, holding experience constant.

### 3.8 Least Squares Residuals

As a by-product of estimation, we define the **fitted value**

$$\hat{y}_i = \mathbf{x}_i' \hat{\boldsymbol{\beta}}$$

and the **residual**

$$\hat{e}_i = y_i - \hat{y}_i = y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}. \quad (3.16)$$

Sometimes  $\hat{y}_i$  is called the predicted value, but this is a misleading label. The fitted value  $\hat{y}_i$  is a function of the entire sample, including  $y_i$ , and thus cannot be interpreted as a valid prediction of  $y_i$ . It is thus more accurate to describe  $\hat{y}_i$  as a *fitted* rather than a *predicted* value.

Note that  $y_i = \hat{y}_i + \hat{e}_i$  and

$$y_i = \mathbf{x}_i' \hat{\boldsymbol{\beta}} + \hat{e}_i. \quad (3.17)$$

We make a distinction between the **error**  $e_i$  and the **residual**  $\hat{e}_i$ . The error  $e_i$  is unobservable while the residual  $\hat{e}_i$  is a by-product of estimation. These two variables are frequently mislabeled, which can cause confusion.

Equation (3.11) implies that

$$\sum_{i=1}^n \mathbf{x}_i \hat{e}_i = \mathbf{0}. \quad (3.18)$$

To see this by a direct calculation, using (3.16) and (3.12),

$$\begin{aligned} \sum_{i=1}^n \mathbf{x}_i \hat{e}_i &= \sum_{i=1}^n \mathbf{x}_i (y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}) \\ &= \sum_{i=1}^n \mathbf{x}_i y_i - \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \hat{\boldsymbol{\beta}} \\ &= \sum_{i=1}^n \mathbf{x}_i y_i - \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left( \sum_{i=1}^n \mathbf{x}_i y_i \right) \\ &= \sum_{i=1}^n \mathbf{x}_i y_i - \sum_{i=1}^n \mathbf{x}_i y_i \\ &= \mathbf{0}. \end{aligned}$$

When  $\mathbf{x}_i$  contains a constant, an implication of (3.18) is

$$\frac{1}{n} \sum_{i=1}^n \hat{e}_i = 0. \quad (3.19)$$

Thus the residuals have a sample mean of zero and the sample correlation between the regressors and the residual is zero. These are algebraic results, and hold true for all linear regression estimates.

### 3.9 Demeaned Regressors

Sometimes it is useful to separate the constant from the other regressors, and write the linear projection equation in the format

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + \alpha + e_i$$

where  $\alpha$  is the intercept and  $\mathbf{x}_i$  does not contain a constant. The least-squares estimates and residuals can be written as

$$y_i = \mathbf{x}_i' \hat{\boldsymbol{\beta}} + \hat{\alpha} + \hat{e}_i$$

In this case (3.18) can be written as the equation system

$$\begin{aligned} \sum_{i=1}^n (y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}} - \hat{\alpha}) &= 0 \\ \sum_{i=1}^n \mathbf{x}_i (y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}} - \hat{\alpha}) &= \mathbf{0} \end{aligned}$$

The first equation implies

$$\hat{\alpha} = \bar{y} - \bar{\mathbf{x}}' \hat{\boldsymbol{\beta}}.$$

Subtracting from the second we obtain

$$\sum_{i=1}^n \mathbf{x}_i \left( (y_i - \bar{y}) - (\mathbf{x}_i - \bar{\mathbf{x}})' \hat{\boldsymbol{\beta}} \right) = \mathbf{0}.$$

Solving for  $\hat{\boldsymbol{\beta}}$  we find

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \left( \sum_{i=1}^n \mathbf{x}_i (\mathbf{x}_i - \bar{\mathbf{x}})' \right)^{-1} \left( \sum_{i=1}^n \mathbf{x}_i (y_i - \bar{y}) \right) \\ &= \left( \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})' \right)^{-1} \left( \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (y_i - \bar{y}) \right). \end{aligned} \quad (3.20)$$

Thus the OLS estimator for the slope coefficients is a regression with demeaned data.

The representation (3.20) is known as the demeaned formula for the least-squares estimator.

### 3.10 Model in Matrix Notation

For many purposes, including computation, it is convenient to write the model and statistics in matrix notation. The linear equation (2.26) is a system of  $n$  equations, one for each observation. We can stack these  $n$  equations together as

$$\begin{aligned} y_1 &= \mathbf{x}_1' \boldsymbol{\beta} + e_1 \\ y_2 &= \mathbf{x}_2' \boldsymbol{\beta} + e_2 \\ &\vdots \\ y_n &= \mathbf{x}_n' \boldsymbol{\beta} + e_n. \end{aligned}$$

Now define

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} \mathbf{x}_1' \\ \mathbf{x}_2' \\ \vdots \\ \mathbf{x}_n' \end{pmatrix}, \quad \mathbf{e} = \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}.$$



Observe that  $\mathbf{y}$  and  $\mathbf{e}$  are  $n \times 1$  vectors, and  $\mathbf{X}$  is an  $n \times k$  matrix. Then the system of  $n$  equations can be compactly written in the single equation

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}. \quad (3.21)$$

Sample sums can be written in matrix notation. For example

$$\begin{aligned} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' &= \mathbf{X}'\mathbf{X} \\ \sum_{i=1}^n \mathbf{x}_i y_i &= \mathbf{X}'\mathbf{y}. \end{aligned}$$

Therefore the least-squares estimator can be written as

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{y}). \quad (3.22)$$

The matrix version of (3.17) and estimated version of (3.21) is

$$\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\mathbf{e}},$$

or equivalently the residual vector is

$$\hat{\mathbf{e}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}.$$

Using the residual vector, we can write (3.18) as

$$\mathbf{X}'\hat{\mathbf{e}} = \mathbf{0}. \quad (3.24)$$

Using matrix notation we have simple expressions for most estimators. This is particularly convenient for computer programming, as most languages allow matrix notation and manipulation.

#### Important Matrix Expressions

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{e} \\ \hat{\boldsymbol{\beta}} &= (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{y}) \\ \hat{\mathbf{e}} &= \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} \\ \mathbf{X}'\hat{\mathbf{e}} &= \mathbf{0}. \end{aligned}$$

#### Early Use of Matrices

The earliest known treatment of the use of matrix methods to solve simultaneous systems is found in Chapter 8 of the Chinese text *The Nine Chapters on the Mathematical Art*, written by several generations of scholars from the 10th to 2nd century BCE.

### 3.11 Projection Matrix

Define the matrix

$$\mathbf{P} = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'.$$

Observe that

$$\mathbf{P}\mathbf{X} = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X} = \mathbf{X}.$$

This is a property of a **projection matrix**. More generally, for any matrix  $\mathbf{Z}$  which can be written as  $\mathbf{Z} = \mathbf{X}\mathbf{\Gamma}$  for some matrix  $\mathbf{\Gamma}$  (we say that  $\mathbf{Z}$  lies in the **range space** of  $\mathbf{X}$ ), then

$$\mathbf{P}\mathbf{Z} = \mathbf{P}\mathbf{X}\mathbf{\Gamma} = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X}\mathbf{\Gamma} = \mathbf{X}\mathbf{\Gamma} = \mathbf{Z}.$$

As an important example, if we partition the matrix  $\mathbf{X}$  into two matrices  $\mathbf{X}_1$  and  $\mathbf{X}_2$  so that

$$\mathbf{X} = [\mathbf{X}_1 \quad \mathbf{X}_2],$$

then  $\mathbf{P}\mathbf{X}_1 = \mathbf{X}_1$ . (See Exercise 3.7.)

The matrix  $\mathbf{P}$  is **symmetric** ( $\mathbf{P}' = \mathbf{P}$ ) and **idempotent** ( $\mathbf{P}\mathbf{P} = \mathbf{P}$ ). (See Section ??.) To see that it is symmetric,

$$\begin{aligned} \mathbf{P}' &= \left( \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \right)' \\ &= (\mathbf{X}')' \left( (\mathbf{X}'\mathbf{X})^{-1} \right)' (\mathbf{X})' \\ &= \mathbf{X} \left( (\mathbf{X}'\mathbf{X})' \right)^{-1} \mathbf{X}' \\ &= \mathbf{X} \left( (\mathbf{X})' (\mathbf{X}')' \right)^{-1} \mathbf{X}' \\ &= \mathbf{P}. \end{aligned}$$

To establish that it is idempotent, the fact that  $\mathbf{P}\mathbf{X} = \mathbf{X}$  implies that

$$\begin{aligned} \mathbf{P}\mathbf{P} &= \mathbf{P}\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \\ &= \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \\ &= \mathbf{P}. \end{aligned}$$

The matrix  $\mathbf{P}$  has the property that it creates the fitted values in a least-squares regression:

$$\mathbf{P}\mathbf{y} = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} = \hat{\mathbf{y}}.$$

Because of this property,  $\mathbf{P}$  is also known as the “hat matrix”.

A special example of a projection matrix occurs when  $\mathbf{X} = \mathbf{1}$  is an  $n$ -vector of ones. Then

$$\begin{aligned} \mathbf{P}_1 &= \mathbf{1} (\mathbf{1}'\mathbf{1})^{-1} \mathbf{1}' \\ &= \frac{1}{n} \mathbf{1}\mathbf{1}'. \end{aligned}$$

Note that

$$\begin{aligned} \mathbf{P}_1\mathbf{y} &= \mathbf{1} (\mathbf{1}'\mathbf{1})^{-1} \mathbf{1}'\mathbf{y} \\ &= \mathbf{1}\bar{y} \end{aligned}$$

creates an  $n$ -vector whose elements are the sample mean  $\bar{y}$  of  $y_i$ .

The  $i^{th}$  diagonal element of  $\mathbf{P} = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$  is

$$h_{ii} = \mathbf{x}'_i (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i \tag{3.25}$$

which is called the **leverage** of the  $i^{th}$  observation.

Two useful properties of the the matrix  $\mathbf{P}$  and the leverage values  $h_{ii}$  are now summarized.

**Theorem 3.11.1**

$$\sum_{i=1}^n h_{ii} = \text{tr } \mathbf{P} = k \quad (3.26)$$

and

$$0 \leq h_{ii} \leq 1. \quad (3.27)$$

To show (3.26),

$$\begin{aligned} \text{tr } \mathbf{P} &= \text{tr} \left( \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \right) \\ &= \text{tr} \left( (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{X} \right) \\ &= \text{tr} (\mathbf{I}_k) \\ &= k. \end{aligned}$$

See Appendix A.5 for definition and properties of the trace operator. The proof of (3.27) is deferred to Section 3.21. One implication is that the rank of  $\mathbf{P}$  is  $k$ .

### 3.12 Orthogonal Projection

Define

$$\begin{aligned} \mathbf{M} &= \mathbf{I}_n - \mathbf{P} \\ &= \mathbf{I}_n - \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \end{aligned}$$

where  $\mathbf{I}_n$  is the  $n \times n$  identity matrix. Note that

$$\mathbf{M} \mathbf{X} = (\mathbf{I}_n - \mathbf{P}) \mathbf{X} = \mathbf{X} - \mathbf{P} \mathbf{X} = \mathbf{X} - \mathbf{X} = \mathbf{0}. \quad (3.28)$$

Thus  $\mathbf{M}$  and  $\mathbf{X}$  are orthogonal. We call  $\mathbf{M}$  an **orthogonal projection matrix**, or more colorfully an **annihilator matrix**, due to the property that for any matrix  $\mathbf{Z}$  in the range space of  $\mathbf{X}$  then

$$\mathbf{M} \mathbf{Z} = \mathbf{Z} - \mathbf{P} \mathbf{Z} = \mathbf{0}.$$

For example,  $\mathbf{M} \mathbf{X}_1 = \mathbf{0}$  for any subcomponent  $\mathbf{X}_1$  of  $\mathbf{X}$ , and  $\mathbf{M} \mathbf{P} = \mathbf{0}$  (see Exercise 3.7).

The orthogonal projection matrix  $\mathbf{M}$  has similar properties with  $\mathbf{P}$ , including that  $\mathbf{M}$  is symmetric ( $\mathbf{M}' = \mathbf{M}$ ) and idempotent ( $\mathbf{M} \mathbf{M} = \mathbf{M}$ ). Similarly to (3.26) we can calculate

$$\text{tr } \mathbf{M} = n - k. \quad (3.29)$$

(See Exercise 3.9.) One implication is that the rank of  $\mathbf{M}$  is  $n - k$ .

While  $\mathbf{P}$  creates fitted values,  $\mathbf{M}$  creates least-squares residuals:

$$\mathbf{M} \mathbf{y} = \mathbf{y} - \mathbf{P} \mathbf{y} = \mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}} = \hat{\mathbf{e}}. \quad (3.30)$$

As discussed in the previous section, a special example of a projection matrix occurs when  $\mathbf{X} = \mathbf{1}$  is an  $n$ -vector of ones, so that  $\mathbf{P}_1 = \mathbf{1} (\mathbf{1}' \mathbf{1})^{-1} \mathbf{1}'$ . Similarly, set

$$\begin{aligned} \mathbf{M}_1 &= \mathbf{I}_n - \mathbf{P}_1 \\ &= \mathbf{I}_n - \mathbf{1} (\mathbf{1}' \mathbf{1})^{-1} \mathbf{1}'. \end{aligned}$$

While  $\mathbf{P}_1$  creates a vector of sample means,  $\mathbf{M}_1$  creates demeaned values:

$$\mathbf{M}_1 \mathbf{y} = \mathbf{y} - \mathbf{1}\bar{y}.$$

For simplicity we will often write the right-hand-side as  $\mathbf{y} - \bar{y}$ . The  $i^{th}$  element is  $y_i - \bar{y}$ , the **demeaned** value of  $y_i$ .

We can also use (3.30) to write an alternative expression for the residual vector. Substituting  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$  into  $\hat{\mathbf{e}} = \mathbf{M}\mathbf{y}$  and using  $\mathbf{M}\mathbf{X} = \mathbf{0}$  we find

$$\hat{\mathbf{e}} = \mathbf{M}\mathbf{y} = \mathbf{M}(\mathbf{X}\boldsymbol{\beta} + \mathbf{e}) = \mathbf{M}\mathbf{e} \quad (3.31)$$

which is free of dependence on the regression coefficient  $\boldsymbol{\beta}$ .

### 3.13 Estimation of Error Variance

The error variance  $\sigma^2 = \mathbb{E}(e_i^2)$  is a moment, so a natural estimator is a moment estimator. If  $e_i$  were observed we would estimate  $\sigma^2$  by

$$\tilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n e_i^2. \quad (3.32)$$

However, this is infeasible as  $e_i$  is not observed. In this case it is common to take a two-step approach to estimation. The residuals  $\hat{e}_i$  are calculated in the first step, and then we substitute  $\hat{e}_i$  for  $e_i$  in expression (3.32) to obtain the feasible estimator

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{e}_i^2. \quad (3.33)$$

In matrix notation, we can write (3.32) and (3.33) as

$$\tilde{\sigma}^2 = n^{-1} \mathbf{e}'\mathbf{e}$$

and

$$\hat{\sigma}^2 = n^{-1} \hat{\mathbf{e}}'\hat{\mathbf{e}}. \quad (3.34)$$

Recall the expressions  $\hat{\mathbf{e}} = \mathbf{M}\mathbf{y} = \mathbf{M}\mathbf{e}$  from (3.30) and (3.31). Applied to (3.34) we find

$$\begin{aligned} \hat{\sigma}^2 &= n^{-1} \hat{\mathbf{e}}'\hat{\mathbf{e}} \\ &= n^{-1} \mathbf{y}'\mathbf{M}\mathbf{M}\mathbf{y} \\ &= n^{-1} \mathbf{y}'\mathbf{M}\mathbf{y} \\ &= n^{-1} \mathbf{e}'\mathbf{M}\mathbf{e} \end{aligned} \quad (3.35)$$

the third equality since  $\mathbf{M}\mathbf{M} = \mathbf{M}$ .

An interesting implication is that

$$\begin{aligned} \tilde{\sigma}^2 - \hat{\sigma}^2 &= n^{-1} \mathbf{e}'\mathbf{e} - n^{-1} \mathbf{e}'\mathbf{M}\mathbf{e} \\ &= n^{-1} \mathbf{e}'\mathbf{P}\mathbf{e} \\ &\geq 0. \end{aligned}$$

The final inequality holds because  $\mathbf{P}$  is positive semi-definite and  $\mathbf{e}'\mathbf{P}\mathbf{e}$  is a quadratic form. This shows that the feasible estimator  $\hat{\sigma}^2$  is numerically smaller than the idealized estimator (3.32).

### 3.14 Analysis of Variance

Another way of writing (3.30) is

$$\mathbf{y} = \mathbf{P}\mathbf{y} + \mathbf{M}\mathbf{y} = \hat{\mathbf{y}} + \hat{\mathbf{e}}. \quad (3.36)$$

This decomposition is **orthogonal**, that is

$$\hat{\mathbf{y}}'\hat{\mathbf{e}} = (\mathbf{P}\mathbf{y})'(\mathbf{M}\mathbf{y}) = \mathbf{y}'\mathbf{P}\mathbf{M}\mathbf{y} = 0.$$

It follows that

$$\mathbf{y}'\mathbf{y} = \hat{\mathbf{y}}'\hat{\mathbf{y}} + 2\hat{\mathbf{y}}'\hat{\mathbf{e}} + \hat{\mathbf{e}}'\hat{\mathbf{e}} = \hat{\mathbf{y}}'\hat{\mathbf{y}} + \hat{\mathbf{e}}'\hat{\mathbf{e}}$$

or

$$\sum_{i=1}^n y_i^2 = \sum_{i=1}^n \hat{y}_i^2 + \sum_{i=1}^n \hat{e}_i^2.$$

Subtracting  $\bar{y}$  from both sides of (3.36) we obtain

$$\mathbf{y} - \mathbf{1}\bar{y} = \hat{\mathbf{y}} - \mathbf{1}\bar{y} + \hat{\mathbf{e}}.$$

This decomposition is also orthogonal when  $\mathbf{X}$  contains a constant, as

$$(\hat{\mathbf{y}} - \mathbf{1}\bar{y})'\hat{\mathbf{e}} = \hat{\mathbf{y}}'\hat{\mathbf{e}} - \bar{y}\mathbf{1}'\hat{\mathbf{e}} = 0$$

under (3.19). It follows that

$$(\mathbf{y} - \mathbf{1}\bar{y})'(\mathbf{y} - \mathbf{1}\bar{y}) = (\hat{\mathbf{y}} - \mathbf{1}\bar{y})'(\hat{\mathbf{y}} - \mathbf{1}\bar{y}) + \hat{\mathbf{e}}'\hat{\mathbf{e}}$$

or

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n \hat{e}_i^2.$$

This is commonly called the **analysis-of-variance** formula for least squares regression.

A commonly reported statistic is the **coefficient of determination** or **R-squared**:

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\sum_{i=1}^n \hat{e}_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

It is often described as the fraction of the sample variance of  $y_i$  which is explained by the least-squares fit.  $R^2$  is a crude measure of regression fit. We have better measures of fit, but these require a statistical (not just algebraic) analysis and we will return to these issues later. One deficiency with  $R^2$  is that it increases when regressors are added to a regression (see Exercise 3.16) so the “fit” can be always increased by increasing the number of regressors.

### 3.15 Regression Components

Partition

$$\mathbf{X} = [\mathbf{X}_1 \quad \mathbf{X}_2]$$

and

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}.$$

Then the regression model can be rewritten as

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{e}. \quad (3.37)$$

The OLS estimator of  $\beta = (\beta'_1, \beta'_2)'$  is obtained by regression of  $\mathbf{y}$  on  $\mathbf{X} = [\mathbf{X}_1 \ \mathbf{X}_2]$  and can be written as

$$\mathbf{y} = \mathbf{X}\hat{\beta} + \hat{\mathbf{e}} = \mathbf{X}_1\hat{\beta}_1 + \mathbf{X}_2\hat{\beta}_2 + \hat{\mathbf{e}}. \quad (3.38)$$

We are interested in algebraic expressions for  $\hat{\beta}_1$  and  $\hat{\beta}_2$ .

The algebra for the estimator is identical as that for the population coefficients as presented in Section 2.21.

Partition  $\hat{\mathbf{Q}}_{xx}$  as

$$\hat{\mathbf{Q}}_{xx} = \begin{bmatrix} \hat{\mathbf{Q}}_{11} & \hat{\mathbf{Q}}_{12} \\ \hat{\mathbf{Q}}_{21} & \hat{\mathbf{Q}}_{22} \end{bmatrix} = \begin{bmatrix} \frac{1}{n}\mathbf{X}'_1\mathbf{X}_1 & \frac{1}{n}\mathbf{X}'_1\mathbf{X}_2 \\ \frac{1}{n}\mathbf{X}'_2\mathbf{X}_1 & \frac{1}{n}\mathbf{X}'_2\mathbf{X}_2 \end{bmatrix}$$

and similarly  $\hat{\mathbf{Q}}_{xy}$

$$\hat{\mathbf{Q}}_{xy} = \begin{bmatrix} \hat{\mathbf{Q}}_{1y} \\ \hat{\mathbf{Q}}_{2y} \end{bmatrix} = \begin{bmatrix} \frac{1}{n}\mathbf{X}'_1\mathbf{y} \\ \frac{1}{n}\mathbf{X}'_2\mathbf{y} \end{bmatrix}.$$

By the partitioned matrix inversion formula (A.4)

$$\hat{\mathbf{Q}}_{xx}^{-1} = \begin{bmatrix} \hat{\mathbf{Q}}_{11} & \hat{\mathbf{Q}}_{12} \\ \hat{\mathbf{Q}}_{21} & \hat{\mathbf{Q}}_{22} \end{bmatrix}^{-1} \stackrel{def}{=} \begin{bmatrix} \hat{\mathbf{Q}}^{11} & \hat{\mathbf{Q}}^{12} \\ \hat{\mathbf{Q}}^{21} & \hat{\mathbf{Q}}^{22} \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{Q}}_{11 \cdot 2}^{-1} & -\hat{\mathbf{Q}}_{11 \cdot 2}^{-1}\hat{\mathbf{Q}}_{12}\hat{\mathbf{Q}}_{22}^{-1} \\ -\hat{\mathbf{Q}}_{22 \cdot 1}^{-1}\hat{\mathbf{Q}}_{21}\hat{\mathbf{Q}}_{11}^{-1} & \hat{\mathbf{Q}}_{22 \cdot 1}^{-1} \end{bmatrix} \quad (3.39)$$

where  $\hat{\mathbf{Q}}_{11 \cdot 2} = \hat{\mathbf{Q}}_{11} - \hat{\mathbf{Q}}_{12}\hat{\mathbf{Q}}_{22}^{-1}\hat{\mathbf{Q}}_{21}$  and  $\hat{\mathbf{Q}}_{22 \cdot 1} = \hat{\mathbf{Q}}_{22} - \hat{\mathbf{Q}}_{21}\hat{\mathbf{Q}}_{11}^{-1}\hat{\mathbf{Q}}_{12}$ .

Thus

$$\begin{aligned} \hat{\beta} &= \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} \\ &= \begin{bmatrix} \hat{\mathbf{Q}}_{11 \cdot 2}^{-1} & -\hat{\mathbf{Q}}_{11 \cdot 2}^{-1}\hat{\mathbf{Q}}_{12}\hat{\mathbf{Q}}_{22}^{-1} \\ -\hat{\mathbf{Q}}_{22 \cdot 1}^{-1}\hat{\mathbf{Q}}_{21}\hat{\mathbf{Q}}_{11}^{-1} & \hat{\mathbf{Q}}_{22 \cdot 1}^{-1} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{Q}}_{1y} \\ \hat{\mathbf{Q}}_{2y} \end{bmatrix} \\ &= \begin{pmatrix} \hat{\mathbf{Q}}_{11 \cdot 2}^{-1}\hat{\mathbf{Q}}_{1y \cdot 2} \\ \hat{\mathbf{Q}}_{22 \cdot 1}^{-1}\hat{\mathbf{Q}}_{2y \cdot 1} \end{pmatrix}. \end{aligned}$$

Now

$$\begin{aligned} \hat{\mathbf{Q}}_{11 \cdot 2} &= \hat{\mathbf{Q}}_{11} - \hat{\mathbf{Q}}_{12}\hat{\mathbf{Q}}_{22}^{-1}\hat{\mathbf{Q}}_{21} \\ &= \frac{1}{n}\mathbf{X}'_1\mathbf{X}_1 - \frac{1}{n}\mathbf{X}'_1\mathbf{X}_2 \left( \frac{1}{n}\mathbf{X}'_2\mathbf{X}_2 \right)^{-1} \frac{1}{n}\mathbf{X}'_2\mathbf{X}_1 \\ &= \frac{1}{n}\mathbf{X}'_1\mathbf{M}_2\mathbf{X}_1 \end{aligned}$$

where

$$\mathbf{M}_2 = \mathbf{I}_n - \mathbf{X}_2(\mathbf{X}'_2\mathbf{X}_2)^{-1}\mathbf{X}'_2$$

is the orthogonal projection matrix for  $\mathbf{X}_2$ . Similarly  $\hat{\mathbf{Q}}_{22 \cdot 1} = \frac{1}{n}\mathbf{X}'_2\mathbf{M}_1\mathbf{X}_2$  where

$$\mathbf{M}_1 = \mathbf{I}_n - \mathbf{X}_1(\mathbf{X}'_1\mathbf{X}_1)^{-1}\mathbf{X}'_1$$

is the orthogonal projection matrix for  $\mathbf{X}_1$ . Also

$$\begin{aligned}\widehat{\mathbf{Q}}_{1y \cdot 2} &= \widehat{\mathbf{Q}}_{1y} - \widehat{\mathbf{Q}}_{12} \widehat{\mathbf{Q}}_{22}^{-1} \widehat{\mathbf{Q}}_{2y} \\ &= \frac{1}{n} \mathbf{X}'_1 \mathbf{y} - \frac{1}{n} \mathbf{X}'_1 \mathbf{X}_2 \left( \frac{1}{n} \mathbf{X}'_2 \mathbf{X}_2 \right)^{-1} \frac{1}{n} \mathbf{X}'_2 \mathbf{y} \\ &= \frac{1}{n} \mathbf{X}'_1 \mathbf{M}_2 \mathbf{y}\end{aligned}$$

and  $\widehat{\mathbf{Q}}_{2y \cdot 1} = \frac{1}{n} \mathbf{X}'_2 \mathbf{M}_1 \mathbf{y}$ .

Therefore

$$\widehat{\boldsymbol{\beta}}_1 = (\mathbf{X}'_1 \mathbf{M}_2 \mathbf{X}_1)^{-1} (\mathbf{X}'_1 \mathbf{M}_2 \mathbf{y}) \quad (3.40)$$

and

$$\widehat{\boldsymbol{\beta}}_2 = (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{X}_2)^{-1} (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{y}). \quad (3.41)$$

These are algebraic expressions for the sub-coefficient estimates from (3.38).

### 3.16 Residual Regression

As first recognized by Frisch and Waugh (1933), expressions (3.40) and (3.41) can be used to show that the least-squares estimators  $\widehat{\boldsymbol{\beta}}_1$  and  $\widehat{\boldsymbol{\beta}}_2$  can be found by a two-step regression procedure.

Take (3.41). Since  $\mathbf{M}_1$  is idempotent,  $\mathbf{M}_1 = \mathbf{M}_1 \mathbf{M}_1$  and thus

$$\begin{aligned}\widehat{\boldsymbol{\beta}}_2 &= (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{X}_2)^{-1} (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{y}) \\ &= (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{M}_1 \mathbf{X}_2)^{-1} (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{M}_1 \mathbf{y}) \\ &= (\widetilde{\mathbf{X}}'_2 \widetilde{\mathbf{X}}_2)^{-1} (\widetilde{\mathbf{X}}'_2 \widetilde{\mathbf{e}}_1)\end{aligned}$$

where

$$\widetilde{\mathbf{X}}_2 = \mathbf{M}_1 \mathbf{X}_2$$

and

$$\widetilde{\mathbf{e}}_1 = \mathbf{M}_1 \mathbf{y}.$$

Thus the coefficient estimate  $\widehat{\boldsymbol{\beta}}_2$  is algebraically equal to the least-squares regression of  $\widetilde{\mathbf{e}}_1$  on  $\widetilde{\mathbf{X}}_2$ . Notice that these two are  $\mathbf{y}$  and  $\mathbf{X}_2$ , respectively, premultiplied by  $\mathbf{M}_1$ . But we know that multiplication by  $\mathbf{M}_1$  is equivalent to creating least-squares residuals. Therefore  $\widetilde{\mathbf{e}}_1$  is simply the least-squares residual from a regression of  $\mathbf{y}$  on  $\mathbf{X}_1$ , and the columns of  $\widetilde{\mathbf{X}}_2$  are the least-squares residuals from the regressions of the columns of  $\mathbf{X}_2$  on  $\mathbf{X}_1$ .

We have proven the following theorem.

**Theorem 3.16.1 Frisch-Waugh-Lovell (FWL)**

*In the model (3.37), the OLS estimator of  $\boldsymbol{\beta}_2$  and the OLS residuals  $\widehat{\mathbf{e}}$  may be equivalently computed by either the OLS regression (3.38) or via the following algorithm:*

1. Regress  $\mathbf{y}$  on  $\mathbf{X}_1$ , obtain residuals  $\widetilde{\mathbf{e}}_1$ ;
2. Regress  $\mathbf{X}_2$  on  $\mathbf{X}_1$ , obtain residuals  $\widetilde{\mathbf{X}}_2$ ;
3. Regress  $\widetilde{\mathbf{e}}_1$  on  $\widetilde{\mathbf{X}}_2$ , obtain OLS estimates  $\widehat{\boldsymbol{\beta}}_2$  and residuals  $\widehat{\mathbf{e}}$ .

In some contexts, the FWL theorem can be used to speed computation, but in most cases there is little computational advantage to using the two-step algorithm.

This result is a direct analogy of the coefficient representation obtained in Section 2.22. The result obtained in that section concerned the population projection coefficients, the result obtained here concern the least-squares estimates. The key message is the same. In the least-squares regression (3.38), the estimated coefficient  $\hat{\beta}_2$  numerically equals the regression of  $\mathbf{y}$  on the regressors  $\mathbf{X}_2$ , only after the regressors  $\mathbf{X}_1$  have been linearly projected out. Similarly, the coefficient estimate  $\hat{\beta}_1$  numerically equals the regression of  $\mathbf{y}$  on the regressors  $\mathbf{X}_1$ , after the regressors  $\mathbf{X}_2$  have been linearly projected out. This result can be very insightful when interpreting regression coefficients.

A common application of the FWL theorem is the demeaning formula for regression obtained in (3.20). Partition  $\mathbf{X} = [\mathbf{X}_1 \ \mathbf{X}_2]$  where  $\mathbf{X}_1 = \mathbf{1}$  is a vector of ones and  $\mathbf{X}_2$  is a matrix of observed regressors. In this case,

$$\mathbf{M}_1 = \mathbf{I}_n - \mathbf{1} (\mathbf{1}'\mathbf{1})^{-1} \mathbf{1}'.$$

Observe that

$$\widetilde{\mathbf{X}}_2 = \mathbf{M}_1 \mathbf{X}_2 = \mathbf{X}_2 - \overline{\mathbf{X}}_2$$

and

$$\mathbf{M}_1 \mathbf{y} = \mathbf{y} - \bar{y}$$

are the “demeaned” variables. The FWL theorem says that  $\hat{\beta}_2$  is the OLS estimate from a regression of  $y_i - \bar{y}$  on  $x_{2i} - \bar{x}_2$ :

$$\hat{\beta}_2 = \left( \sum_{i=1}^n (x_{2i} - \bar{x}_2) (x_{2i} - \bar{x}_2)' \right)^{-1} \left( \sum_{i=1}^n (x_{2i} - \bar{x}_2) (y_i - \bar{y}) \right).$$

This is (3.20).

### Ragnar Frisch

Ragnar Frisch (1895-1973) was co-winner with Jan Tinbergen of the first Nobel Memorial Prize in Economic Sciences in 1969 for their work in developing and applying dynamic models for the analysis of economic problems. Frisch made a number of foundational contributions to modern economics beyond the Frisch-Waugh-Lovell Theorem, including formalizing consumer theory, production theory, and business cycle theory.

## 3.17 Prediction Errors

The least-squares residual  $\hat{e}_i$  are not true prediction errors, as they are constructed based on the full sample including  $y_i$ . A proper prediction for  $y_i$  should be based on estimates constructed using only the other observations. We can do this by defining the **leave-one-out** OLS estimator of  $\beta$  as that obtained from the sample of  $n - 1$  observations *excluding* the  $i^{th}$  observation:

$$\begin{aligned} \hat{\beta}_{(-i)} &= \left( \frac{1}{n-1} \sum_{j \neq i} \mathbf{x}_j \mathbf{x}_j' \right)^{-1} \left( \frac{1}{n-1} \sum_{j \neq i} \mathbf{x}_j y_j \right) \\ &= \left( \mathbf{X}_{(-i)}' \mathbf{X}_{(-i)} \right)^{-1} \mathbf{X}_{(-i)}' \mathbf{y}_{(-i)}. \end{aligned} \tag{3.42}$$



Here,  $\mathbf{X}_{(-i)}$  and  $\mathbf{y}_{(-i)}$  are the data matrices omitting the  $i^{th}$  row. The leave-one-out predicted value for  $y_i$  is

$$\tilde{y}_i = \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{(-i)},$$

and the **leave-one-out residual** or **prediction error** or **prediction residual** is

$$\tilde{e}_i = y_i - \tilde{y}_i.$$

A convenient alternative expression for  $\hat{\boldsymbol{\beta}}_{(-i)}$  (derived in Section 3.21) is

$$\hat{\boldsymbol{\beta}}_{(-i)} = \hat{\boldsymbol{\beta}} - (1 - h_{ii})^{-1} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i \hat{e}_i \quad (3.43)$$

where  $h_{ii}$  are the leverage values as defined in (3.25).

Using (3.43) we can simplify the expression for the prediction error:

$$\begin{aligned} \tilde{e}_i &= y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{(-i)} \\ &= y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}} + (1 - h_{ii})^{-1} \mathbf{x}_i' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i \hat{e}_i \\ &= \hat{e}_i + (1 - h_{ii})^{-1} h_{ii} \hat{e}_i \\ &= (1 - h_{ii})^{-1} \hat{e}_i. \end{aligned} \quad (3.44)$$

To write this in vector notation, define

$$\begin{aligned} \mathbf{M}^* &= (\mathbf{I}_n - \text{diag}\{h_{11}, \dots, h_{nn}\})^{-1} \\ &= \text{diag}\{(1 - h_{11})^{-1}, \dots, (1 - h_{nn})^{-1}\}. \end{aligned} \quad (3.45)$$

Then (3.44) is equivalent to

$$\tilde{\mathbf{e}} = \mathbf{M}^* \hat{\mathbf{e}}. \quad (3.46)$$

A convenient feature of this expression is that it shows that computation of the full vector of prediction errors  $\tilde{\mathbf{e}}$  is based on a simple linear operation, and does not really require  $n$  separate estimations.

One use of the prediction errors is to estimate the out-of-sample mean squared error

$$\begin{aligned} \tilde{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n \tilde{e}_i^2 \\ &= \frac{1}{n} \sum_{i=1}^n (1 - h_{ii})^{-2} \hat{e}_i^2. \end{aligned} \quad (3.47)$$

This is also known as the **sample mean squared prediction error**. Its square root  $\tilde{\sigma} = \sqrt{\tilde{\sigma}^2}$  is the **prediction standard error**.

### 3.18 Influential Observations

Another use of the leave-one-out estimator is to investigate the impact of **influential observations**, sometimes called **outliers**. We say that observation  $i$  is influential if its omission from the sample induces a substantial change in a parameter estimate of interest.

For illustration, consider Figure 3.2 which shows a scatter plot of random variables  $(y_i, x_i)$ . The 25 observations shown with the open circles are generated by  $x_i \sim U[1, 10]$  and  $y_i \sim N(x_i, 4)$ . The 26<sup>th</sup> observation shown with the filled circle is  $x_{26} = 9$ ,  $y_{26} = 0$ . (Imagine that  $y_{26} = 0$  was incorrectly recorded due to a mistaken key entry.) The Figure shows both the least-squares fitted line from the full sample and that obtained after deletion of the 26<sup>th</sup> observation from the sample. In this example we can see how the 26<sup>th</sup> observation (the “outlier”) greatly tilts the least-squares

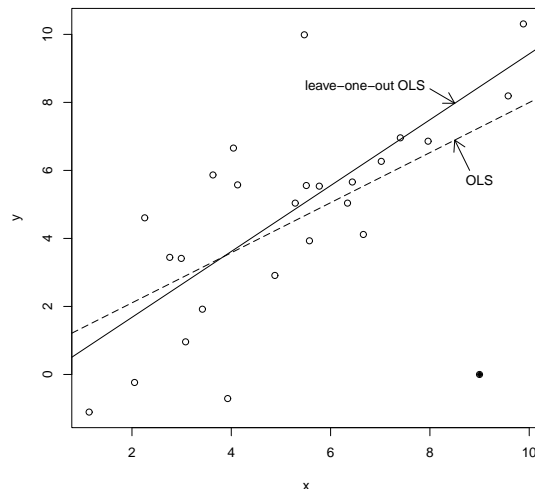


Figure 3.2: Impact of an influential observation on the least-squares estimator

fitted line towards the 26<sup>th</sup> observation. In fact, the slope coefficient decreases from 0.97 (which is close to the true value of 1.00) to 0.56, which is substantially reduced. Neither  $y_{26}$  nor  $x_{26}$  are unusual values relative to their marginal distributions, so this outlier would not have been detected from examination of the marginal distributions of the data. The change in the slope coefficient of  $-0.41$  is meaningful and should raise concern to an applied economist.

From (3.43)-(3.44) we know that

$$\begin{aligned}\hat{\beta} - \hat{\beta}_{(-i)} &= (1 - h_{ii})^{-1} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i \hat{e}_i \\ &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i \tilde{e}_i.\end{aligned}\tag{3.48}$$

By direct calculation of this quantity for each observation  $i$ , we can directly discover if a specific observation  $i$  is influential for a coefficient estimate of interest.

For a general assessment, we can focus on the predicted values. The difference between the full-sample and leave-one-out predicted values is

$$\begin{aligned}\hat{y}_i - \tilde{y}_i &= \mathbf{x}_i' \hat{\beta} - \mathbf{x}_i' \hat{\beta}_{(-i)} \\ &= \mathbf{x}_i' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i \tilde{e}_i \\ &= h_{ii} \tilde{e}_i\end{aligned}$$

which is a simple function of the leverage values  $h_{ii}$  and prediction errors  $\tilde{e}_i$ . Observation  $i$  is influential for the predicted value if  $|h_{ii} \tilde{e}_i|$  is large, which requires that both  $h_{ii}$  and  $|\tilde{e}_i|$  are large.

One way to think about this is that a large leverage value  $h_{ii}$  gives the potential for observation  $i$  to be influential. A large  $h_{ii}$  means that observation  $i$  is unusual in the sense that the regressor  $\mathbf{x}_i$  is far from its sample mean. We call an observation with large  $h_{ii}$  a **leverage point**. A leverage point is not necessarily influential as the latter also requires that the prediction error  $\tilde{e}_i$  is large.

To determine if any individual observations are influential in this sense, several diagnostics have been proposed (some names include DFITS, Cook's Distance, and Welsch Distance). Unfortunately, from a statistical perspective it is difficult to recommend these diagnostics for applications as they are not based on statistical theory. Probably the most relevant measure is the change in the coefficient estimates given in (3.48). The ratio of these changes to the coefficient's standard error is called its DFBETA, and is a postestimation diagnostic available in Stata. While there is no magic threshold, the concern is whether or not an individual observation meaningfully changes an

estimated coefficient of interest. A simple diagnostic for influential observations is to calculate

$$Influence = \max_{1 \leq i \leq n} |\hat{y}_i - \tilde{y}_i| = \max_{1 \leq i \leq n} |h_{ii}\tilde{e}_i|.$$

This is the largest (absolute) change in the predicted value due to a single observation. If this diagnostic is large relative to the distribution of  $y_i$ , it may indicate that that observation is influential.

If an observation is determined to be influential, what should be done? As a common cause of influential observations is data entry error, the influential observations should be examined for evidence that the observation was mis-recorded. Perhaps the observation falls outside of permitted ranges, or some observables are inconsistent (for example, a person is listed as having a job but receives earnings of \$0). If it is determined that an observation is incorrectly recorded, then the observation is typically deleted from the sample. This process is often called “cleaning the data”. The decisions made in this process involve a fair amount of individual judgment. When this is done it is proper empirical practice to document such choices. (It is useful to keep the source data in its original form, a revised data file after cleaning, and a record describing the revision process. This is especially useful when revising empirical work at a later date.)

It is also possible that an observation is correctly measured, but unusual and influential. In this case it is unclear how to proceed. Some researchers will try to alter the specification to properly model the influential observation. Other researchers will delete the observation from the sample. The motivation for this choice is to prevent the results from being skewed or determined by individual observations, but this practice is viewed skeptically by many researchers who believe it reduces the integrity of reported empirical results.

For an empirical illustration, consider the log wage regression (3.15) for single Asian males. This regression, which has 268 observations, has  $Influence = 0.29$ . This means that the most influential observation, when deleted, changes the predicted (fitted) value of the dependent variable  $\log(Wage)$  by 0.29, or equivalently the wage by 29%. This is a meaningful change and suggests further investigation. We examine the influential observation, and find that its leverage  $h_{ii}$  is 0.33, which is disturbingly large. (Recall that the leverage values are all positive and sum to  $k$ . One twelfth of the leverage in this sample of 268 observations is contained in just this single observation!) Examining further, we find that this individual is 65 years old with 8 years education, so that his potential experience is 51 years. This is the highest experience in the subsample – the next highest is 41 years. The large leverage is due to his unusual characteristics (very low education and very high experience) within this sample. Essentially, regression (3.15) is attempting to estimate the conditional mean at  $experience = 51$  with only one observation, so it is not surprising that this observation determines the fit and is thus influential. A reasonable conclusion is the regression function can only be estimated over a smaller range of  $experience$ . We restrict the sample to individuals with less than 45 years experience, re-estimate, and obtain the following estimates.

$$\widehat{\log(Wage)} = 0.144 \text{ education} + 0.043 \text{ experience} - 0.095 \text{ experience}^2/100 + 0.531. \quad (3.49)$$

For this regression, we calculate that  $Influence = 0.11$ , which is greatly reduced relative to the regression (3.15). Comparing (3.49) with (3.15), the slope coefficient for education is essentially unchanged, but the coefficients in experience and its square have slightly increased.

By eliminating the influential observation, equation (3.49) can be viewed as a more robust estimate of the conditional mean for most levels of  $experience$ . Whether to report (3.15) or (3.49) in an application is largely a matter of judgment.

### 3.19 CPS Data Set

In this section we describe the data set used in the empirical illustrations.

The Current Population Survey (CPS) is a monthly survey of about 57,000 U.S. households conducted by the Bureau of the Census of the Bureau of Labor Statistics. The CPS is the primary source of information on the labor force characteristics of the U.S. population. The survey covers employment, earnings, educational attainment, income, poverty, health insurance coverage, job experience, voting and registration, computer usage, veteran status, and other variables. Details can be found at [www.census.gov/cps](http://www.census.gov/cps) and [dataferrett.census.gov](http://dataferrett.census.gov).

From the March 2009 survey we extracted the individuals with non-allocated variables who were full-time employed (defined as those who had worked at least 36 hours per week for at least 48 weeks the past year), and excluded those in the military. This sample has 50,742 individuals. We extracted 14 variables from the CPS on these individuals and created the data files `cps09mar.dta` (Stata format), `cps09mar.xlsx` (Excel format) and `cps09mar.txt` (text format). The variables are described in the file `cps09mar_description.pdf`. All data files are available at <http://www.ssc.wisc.edu/~bhansen/econometrics/>

## 3.20 Programming

Most packages allow both interactive programming (where you enter commands one-by-one) and batch programming (where you run a pre-written sequence of commands from a file). Interactive programming can be useful for exploratory analysis, but eventually all work should be executed in batch mode. This is the best way to control and document your work.

Batch programs are text files where each line executes a single command. For Stata, this file needs to have the filename extension “.do”, and for MATLAB “.m”. For R there is no specific naming requirements, though it is typical to use the extension “.r”.

To execute a program file, you type a command within the program.

Stata: `do chapter3` executes the file `chapter3.do`

MATLAB: `run chapter3` executes the file `chapter3.m`

R: `source("chapter3.r")` executes the file `chapter3.r`

When writing batch files, it is useful to include comments for documentation and readability.

We illustrate programming files for Stata, R, and MATLAB, which execute a portion of the empirical illustrations from Sections 3.7 and 3.18.

### Stata do File

```
*      Clear memory and load the data
clear
use cps09mar.dta
*      Generate transformations
gen wage=ln(earnings/(hours*week))
gen experience = age - education - 6
gen exp2 = (experience^2)/100
*      Create indicator for subsamples
gen mbf = (race == 2) & (marital <= 2) & (female == 1)
gen sam = (race == 4) & (marital == 7) & (female == 0)
*      Regressions
reg wage education if (mbf == 1) & (experience == 12)
reg wage education experience exp2 if sam == 1
*      Leverage and influence
predict leverage,hat
predict e,residual
gen d=e*leverage/(1-leverage)
summarize d if sam ==1
```

**R Program File**

```

#    Load the data and create subsamples
dat <- read.table("cps09mar.txt")
experience <- dat[,1]-dat[,4]-6
mbf <- (dat[,11]==2)&(dat[,12]<=2)&(dat[,2]==1)&(experience==12)
sam <- (dat[,11]==4)&(dat[,12]==7)&(dat[,2]==0)
dat1 <- dat[mbf,]
dat2 <- dat[sam,]
#    First regression
y <- as.matrix(log(dat1[,5]/(dat1[,6]*dat1[,7])))
x <- cbind(dat1[,4],matrix(1,nrow(dat1),1))
beta <- solve(t(x)%*%x,t(x)%*%y)
print(beta)
#    Second regression
y <- as.matrix(log(dat2[,5]/(dat2[,6]*dat2[,7])))
experience <- dat2[,1]-dat2[,4]-6
exp2 <- (experience^2)/100
x <- cbind(dat2[,4],experience,exp2,matrix(1,nrow(dat2),1))
beta <- solve(t(x)%*%x,t(x)%*%y)print(beta)
#    Create leverage and influence
e <- y-x%*%beta
leverage <- rowSums(x*(x%*%solve(t(x)%*%x)))
r <- e/(1-leverage)
d <- leverage*e/(1-leverage)
print(max(abs(d)))

```

**MATLAB Program File**

```

% Load the data and create subsamples
load cps09mar.txt;
dat=cps09mar;
experience=dat(:,1)-dat(:,4)-6;
mbf = (dat(:,11)==2)&(dat(:,12)<=2)&(dat(:,2)==1)&(experience==12);
sam = (dat(:,11)==4)&(dat(:,12)==7)&(dat(:,2)==0);
dat1=dat(mbf,:);
dat2=dat(sam,:);
%    First regression
y=log(dat1(:,5)./(dat1(:,6).*dat1(:,7)));
x=[dat1(:,4),ones(length(dat1),1)];
beta=inv(x'*x)*(x'*y);display(beta);
%    Second regression
y=log(dat2(:,5)./(dat2(:,6).*dat2(:,7)));
experience=dat2(:,1)-dat2(:,4)-6;
exp2 = (experience.^2)/100;
x=[dat2(:,4),experience,exp2,ones(length(dat2),1)];
beta=inv(x'*x)*(x'*y);display(beta);
%    Create leverage and influence
e=y-x*beta;
leverage=sum((x.*(x*inv(x'*x))))';d=leverage.*e./(1-leverage);
influence=max(abs(d));
display(influence);

```

Instead, to load from an excel file, we can replace the first two lines ('load' and 'dat=') with

```

dat=xlsread('cps09mar.xlsx');

```

### 3.21 Technical Proofs\*

**Proof of Theorem 3.11.1, equation (3.27):** First,  $h_{ii} = \mathbf{x}'_i (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i \geq 0$  since it is a quadratic form and  $\mathbf{X}'\mathbf{X} > 0$ . Next, since  $h_{ii}$  is the  $i^{\text{th}}$  diagonal element of the projection matrix  $\mathbf{P} = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$ , then

$$h_{ii} = \mathbf{s}' \mathbf{P} \mathbf{s}$$

where

$$\mathbf{s} = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}$$

is a unit vector with a 1 in the  $i^{\text{th}}$  place (and zeros elsewhere).

By the spectral decomposition of the idempotent matrix  $\mathbf{P}$  (see equation (A.10))

$$\mathbf{P} = \mathbf{B}' \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{B}$$

where  $\mathbf{B}'\mathbf{B} = \mathbf{I}_n$ . Thus letting  $\mathbf{b} = \mathbf{B}\mathbf{s}$  denote the  $i^{\text{th}}$  column of  $\mathbf{B}$ , and partitioning  $\mathbf{b}' = (\mathbf{b}'_1 \quad \mathbf{b}'_2)$  then

$$\begin{aligned} h_{ii} &= \mathbf{s}' \mathbf{B}' \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{B} \mathbf{s} \\ &= \mathbf{b}'_1 \begin{bmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{b}_1 \\ &= \mathbf{b}'_1 \mathbf{b}_1 \\ &\leq \mathbf{b}' \mathbf{b} \\ &= 1 \end{aligned}$$

the final equality since  $\mathbf{b}$  is the  $i^{\text{th}}$  column of  $\mathbf{B}$  and  $\mathbf{B}'\mathbf{B} = \mathbf{I}_n$ . We have shown that  $h_{ii} \leq 1$ , establishing (3.27). ■

**Proof of Equation (3.43).** The Sherman–Morrison formula (A.3) from Appendix A.6 states that for nonsingular  $\mathbf{A}$  and vector  $\mathbf{b}$

$$(\mathbf{A} - \mathbf{b}\mathbf{b}')^{-1} = \mathbf{A}^{-1} + (1 - \mathbf{b}'\mathbf{A}^{-1}\mathbf{b})^{-1} \mathbf{A}^{-1}\mathbf{b}\mathbf{b}'\mathbf{A}^{-1}.$$

This implies

$$(\mathbf{X}'\mathbf{X} - \mathbf{x}_i\mathbf{x}'_i)^{-1} = (\mathbf{X}'\mathbf{X})^{-1} + (1 - h_{ii})^{-1} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i\mathbf{x}'_i (\mathbf{X}'\mathbf{X})^{-1}$$

and thus

$$\begin{aligned} \widehat{\boldsymbol{\beta}}_{(-i)} &= (\mathbf{X}'\mathbf{X} - \mathbf{x}_i\mathbf{x}'_i)^{-1} (\mathbf{X}'\mathbf{y} - \mathbf{x}_iy_i) \\ &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_iy_i \\ &\quad + (1 - h_{ii})^{-1} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i\mathbf{x}'_i (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{y} - \mathbf{x}_iy_i) \\ &= \widehat{\boldsymbol{\beta}} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_iy_i + (1 - h_{ii})^{-1} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i (\mathbf{x}'_i\widehat{\boldsymbol{\beta}} - h_{ii}y_i) \\ &= \widehat{\boldsymbol{\beta}} - (1 - h_{ii})^{-1} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i \left( (1 - h_{ii})y_i - \mathbf{x}'_i\widehat{\boldsymbol{\beta}} + h_{ii}y_i \right) \\ &= \widehat{\boldsymbol{\beta}} - (1 - h_{ii})^{-1} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i \widehat{e}_i \end{aligned}$$

the third equality making the substitutions  $\widehat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$  and  $h_{ii} = \mathbf{x}'_i (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i$ , and the remainder collecting terms. ■

## Exercises

**Exercise 3.1** Let  $y$  be a random variable with  $\mu = \mathbb{E}(y)$  and  $\sigma^2 = \text{var}(y)$ . Define

$$g(y, \mu, \sigma^2) = \begin{pmatrix} y - \mu \\ (y - \mu)^2 - \sigma^2 \end{pmatrix}.$$

Let  $(\hat{\mu}, \hat{\sigma}^2)$  be the values such that  $\bar{g}_n(\hat{\mu}, \hat{\sigma}^2) = \mathbf{0}$  where  $\bar{g}_n(m, s) = n^{-1} \sum_{i=1}^n g(y_i, m, s)$ . Show that  $\hat{\mu}$  and  $\hat{\sigma}^2$  are the sample mean and variance.

**Exercise 3.2** Consider the OLS regression of the  $n \times 1$  vector  $\mathbf{y}$  on the  $n \times k$  matrix  $\mathbf{X}$ . Consider an alternative set of regressors  $\mathbf{Z} = \mathbf{X}\mathbf{C}$ , where  $\mathbf{C}$  is a  $k \times k$  non-singular matrix. Thus, each column of  $\mathbf{Z}$  is a mixture of some of the columns of  $\mathbf{X}$ . Compare the OLS estimates and residuals from the regression of  $\mathbf{y}$  on  $\mathbf{X}$  to the OLS estimates from the regression of  $\mathbf{y}$  on  $\mathbf{Z}$ .

**Exercise 3.3** Using matrix algebra, show  $\mathbf{X}'\hat{\mathbf{e}} = \mathbf{0}$ .

**Exercise 3.4** Let  $\hat{\mathbf{e}}$  be the OLS residual from a regression of  $\mathbf{y}$  on  $\mathbf{X} = [\mathbf{X}_1 \ \mathbf{X}_2]$ . Find  $\mathbf{X}_2'\hat{\mathbf{e}}$ .

**Exercise 3.5** Let  $\hat{\mathbf{e}}$  be the OLS residual from a regression of  $\mathbf{y}$  on  $\mathbf{X}$ . Find the OLS coefficient from a regression of  $\hat{\mathbf{e}}$  on  $\mathbf{X}$ .

**Exercise 3.6** Let  $\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ . Find the OLS coefficient from a regression of  $\hat{\mathbf{y}}$  on  $\mathbf{X}$ .

**Exercise 3.7** Show that if  $\mathbf{X} = [\mathbf{X}_1 \ \mathbf{X}_2]$  then  $\mathbf{P}\mathbf{X}_1 = \mathbf{X}_1$  and  $\mathbf{M}\mathbf{X}_1 = \mathbf{0}$ .

**Exercise 3.8** Show that  $\mathbf{M}$  is idempotent:  $\mathbf{M}\mathbf{M} = \mathbf{M}$ .

**Exercise 3.9** Show that  $\text{tr } \mathbf{M} = n - k$ .

**Exercise 3.10** Show that if  $\mathbf{X} = [\mathbf{X}_1 \ \mathbf{X}_2]$  and  $\mathbf{X}_1'\mathbf{X}_2 = \mathbf{0}$  then  $\mathbf{P} = \mathbf{P}_1 + \mathbf{P}_2$ .

**Exercise 3.11** Show that when  $\mathbf{X}$  contains a constant,  $\frac{1}{n} \sum_{i=1}^n \hat{y}_i = \bar{y}$ .

**Exercise 3.12** A dummy variable takes on only the values 0 and 1. It is used for categorical data, such as an individual's gender. Let  $\mathbf{d}_1$  and  $\mathbf{d}_2$  be vectors of 1's and 0's, with the  $i^{\text{th}}$  element of  $\mathbf{d}_1$  equaling 1 and that of  $\mathbf{d}_2$  equaling 0 if the person is a man, and the reverse if the person is a woman. Suppose that there are  $n_1$  men and  $n_2$  women in the sample. Consider fitting the following three equations by OLS

$$\mathbf{y} = \mu + \mathbf{d}_1\alpha_1 + \mathbf{d}_2\alpha_2 + \mathbf{e} \quad (3.50)$$

$$\mathbf{y} = \mathbf{d}_1\alpha_1 + \mathbf{d}_2\alpha_2 + \mathbf{e} \quad (3.51)$$

$$\mathbf{y} = \mu + \mathbf{d}_1\phi + \mathbf{e} \quad (3.52)$$

Can all three equations (3.50), (3.51), and (3.52) be estimated by OLS? Explain if not.

- Compare regressions (3.51) and (3.52). Is one more general than the other? Explain the relationship between the parameters in (3.51) and (3.52).
- Compute  $\boldsymbol{\iota}'\mathbf{d}_1$  and  $\boldsymbol{\iota}'\mathbf{d}_2$ , where  $\boldsymbol{\iota}$  is an  $n \times 1$  vector of ones.
- Letting  $\boldsymbol{\alpha} = (\alpha_1 \ \alpha_2)'$ , write equation (3.51) as  $\mathbf{y} = \mathbf{X}\boldsymbol{\alpha} + \mathbf{e}$ . Consider the assumption  $\mathbb{E}(\mathbf{x}_i\mathbf{e}_i) = \mathbf{0}$ . Is there any content to this assumption in this setting?



**Exercise 3.13** Let  $\mathbf{d}_1$  and  $\mathbf{d}_2$  be defined as in the previous exercise.

(a) In the OLS regression

$$\mathbf{y} = \mathbf{d}_1\hat{\gamma}_1 + \mathbf{d}_2\hat{\gamma}_2 + \hat{\mathbf{u}},$$

show that  $\hat{\gamma}_1$  is the sample mean of the dependent variable among the men of the sample ( $\bar{y}_1$ ), and that  $\hat{\gamma}_2$  is the sample mean among the women ( $\bar{y}_2$ ).

(b) Let  $\mathbf{X}$  ( $n \times k$ ) be an additional matrix of regressors. Describe in words the transformations

$$\begin{aligned}\mathbf{y}^* &= \mathbf{y} - \mathbf{d}_1\bar{y}_1 - \mathbf{d}_2\bar{y}_2 \\ \mathbf{X}^* &= \mathbf{X} - \mathbf{d}_1\bar{\mathbf{x}}_1' - \mathbf{d}_2\bar{\mathbf{x}}_2'\end{aligned}$$

where  $\bar{\mathbf{x}}_1$  and  $\bar{\mathbf{x}}_2$  are the  $k \times 1$  means of the regressors for men and women, respectively.

(c) Compare  $\tilde{\boldsymbol{\beta}}$  from the OLS regression

$$\mathbf{y}^* = \mathbf{X}^*\tilde{\boldsymbol{\beta}} + \tilde{\mathbf{e}}$$

with  $\hat{\boldsymbol{\beta}}$  from the OLS regression

$$\mathbf{y} = \mathbf{d}_1\hat{\alpha}_1 + \mathbf{d}_2\hat{\alpha}_2 + \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\mathbf{e}}.$$

**Exercise 3.14** Let  $\hat{\boldsymbol{\beta}}_n = (\mathbf{X}_n' \mathbf{X}_n)^{-1} \mathbf{X}_n' \mathbf{y}_n$  denote the OLS estimate when  $\mathbf{y}_n$  is  $n \times 1$  and  $\mathbf{X}_n$  is  $n \times k$ . A new observation ( $y_{n+1}, \mathbf{x}_{n+1}$ ) becomes available. Prove that the OLS estimate computed using this additional observation is

$$\hat{\boldsymbol{\beta}}_{n+1} = \hat{\boldsymbol{\beta}}_n + \frac{1}{1 + \mathbf{x}_{n+1}' (\mathbf{X}_n' \mathbf{X}_n)^{-1} \mathbf{x}_{n+1}} (\mathbf{X}_n' \mathbf{X}_n)^{-1} \mathbf{x}_{n+1} (y_{n+1} - \mathbf{x}_{n+1}' \hat{\boldsymbol{\beta}}_n).$$

**Exercise 3.15** Prove that  $R^2$  is the square of the sample correlation between  $\mathbf{y}$  and  $\hat{\mathbf{y}}$ .

**Exercise 3.16** Consider two least-squares regressions

$$\mathbf{y} = \mathbf{X}_1\tilde{\boldsymbol{\beta}}_1 + \tilde{\mathbf{e}}$$

and

$$\mathbf{y} = \mathbf{X}_1\hat{\boldsymbol{\beta}}_1 + \mathbf{X}_2\hat{\boldsymbol{\beta}}_2 + \hat{\mathbf{e}}.$$

Let  $R_1^2$  and  $R_2^2$  be the  $R$ -squared from the two regressions. Show that  $R_2^2 \geq R_1^2$ . Is there a case (explain) when there is equality  $R_2^2 = R_1^2$ ?

**Exercise 3.17** Show that  $\tilde{\sigma}^2 \geq \hat{\sigma}^2$ . Is equality possible?

**Exercise 3.18** For which observations will  $\hat{\boldsymbol{\beta}}_{(-i)} = \hat{\boldsymbol{\beta}}$ ?

**Exercise 3.19** Consider the least-squares regression estimates

$$y_i = x_{1i}\hat{\beta}_1 + x_{2i}\hat{\beta}_2 + \hat{e}_i$$

and the “one regressor at a time” regression estimates

$$y_i = \tilde{\beta}_1 x_{1i} + \tilde{e}_{1i} \quad y_i = \tilde{\beta}_2 x_{2i} + \tilde{e}_{2i}$$

Under what condition does  $\tilde{\beta}_1 = \hat{\beta}_1$  and  $\tilde{\beta}_2 = \hat{\beta}_2$ ?

**Exercise 3.20** You estimate a least-squares regression

$$y_i = \mathbf{x}'_{1i} \tilde{\boldsymbol{\beta}}_1 + \tilde{u}_i$$

and then regress the residuals on another set of regressors

$$\tilde{u}_i = \mathbf{x}'_{2i} \tilde{\boldsymbol{\beta}}_2 + \tilde{e}_i$$

Does this second regression give you the same estimated coefficients as from estimation of a least-squares regression on both set of regressors?

$$y_i = \mathbf{x}'_{1i} \hat{\boldsymbol{\beta}}_1 + \mathbf{x}'_{2i} \hat{\boldsymbol{\beta}}_2 + \hat{e}_i$$

In other words, is it true that  $\tilde{\boldsymbol{\beta}}_2 = \hat{\boldsymbol{\beta}}_2$ ? Explain your reasoning.

**Exercise 3.21** The data matrix is  $(\mathbf{y}, \mathbf{X})$  with  $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2]$ , and consider the transformed regressor matrix  $\mathbf{Z} = [\mathbf{X}_1, \mathbf{X}_2 - \mathbf{X}_1]$ . Suppose you do a least-squares regression of  $\mathbf{y}$  on  $\mathbf{X}$ , and a least-squares regression of  $\mathbf{y}$  on  $\mathbf{Z}$ . Let  $\hat{\sigma}^2$  and  $\tilde{\sigma}^2$  denote the residual variance estimates from the two regressions. Give a formula relating  $\hat{\sigma}^2$  and  $\tilde{\sigma}^2$ ? (Explain your reasoning.)

**Exercise 3.22** Use the data set from Section 3.19 and the sub-sample used for equation (3.49) (see Section 3.20) for data construction)

- Estimate equation (3.49) and compute the equation  $R^2$  and sum of squared errors.
- Re-estimate the slope on education using the residual regression approach. Regress  $\log(\text{Wage})$  on experience and its square, regress education on experience and its square, and the residuals on the residuals. Report the estimates from this final regression, along with the equation  $R^2$  and sum of squared errors. Does the slope coefficient equal the value in (3.49)? Explain.
- Are the  $R^2$  and sum-of-squared errors from parts (a) and (b) equal? Explain.

**Exercise 3.23** Estimate equation (3.49) as in part (a) of the previous question. Let  $\hat{e}_i$  be the OLS residual,  $\hat{y}_i$  the predicted value from the regression,  $x_{1i}$  be education and  $x_{2i}$  be experience. Numerically calculate the following:

- $\sum_{i=1}^n \hat{e}_i$
- $\sum_{i=1}^n x_{1i} \hat{e}_i$
- $\sum_{i=1}^n x_{2i} \hat{e}_i$
- $\sum_{i=1}^n x_{1i}^2 \hat{e}_i$
- $\sum_{i=1}^n x_{2i}^2 \hat{e}_i$
- $\sum_{i=1}^n \hat{y}_i \hat{e}_i$
- $\sum_{i=1}^n \hat{e}_i^2$

Are these calculations consistent with the theoretical properties of OLS? Explain.

**Exercise 3.24** Use the data set from Section 3.19.

- Estimate a log wage regression for the subsample of white male Hispanics. In addition to education, experience, and its square, include a set of binary variables for regions and marital status. For regions, you create dummy variables for Northeast, South and West so that Midwest is the excluded group. For marital status, create variables for married, widowed or divorced, and separated, so that single (never married) is the excluded group.
- Repeat this estimation using a different econometric package. Compare your results. Do they agree?

## Chapter 4

# Least Squares Regression

### 4.1 Introduction

In this chapter we investigate some finite-sample properties of the least-squares estimator in the linear regression model. In particular, we calculate the finite-sample mean and covariance matrix and propose standard errors for the coefficient estimates.

### 4.2 Random Sampling

Assumption 3.2.1 specified that the observations have identical distributions. To derive the finite-sample properties of the estimators we will need to additionally specify the dependence structure across the observations.

The simplest context is when the observations are mutually independent, in which case we say that they are **independent and identically distributed**, or **i.i.d.** It is also common to describe iid observations as a **random sample**. Traditionally, random sampling has been the default assumption in cross-section (e.g. survey) contexts. It is quite convenient as iid sampling leads to straightforward expressions for estimation variance. The assumption seems appropriate (meaning that it should be approximately valid) when samples are small and relatively dispersed. That is, if you randomly sample 1000 people from a large country such as the United States it seems reasonable to model their responses as mutually independent.

**Assumption 4.2.1** *The observations  $\{(y_1, \mathbf{x}_1), \dots, (y_i, \mathbf{x}_i), \dots, (y_n, \mathbf{x}_n)\}$  are independent and identically distributed.*

For most of this chapter, we will use Assumption 4.2.1 to derive properties of the OLS estimator.

Assumption 4.2.1 means that if you take any two individuals  $i \neq j$  in a sample, the values  $(y_i, \mathbf{x}_i)$  are independent of the values  $(y_j, \mathbf{x}_j)$  yet have the same distribution. Independence means that the decisions and choices of individual  $i$  do not affect the decisions of individual  $j$ , and conversely.

This assumption may be violated if individuals in the sample are connected in some way, for example if they are neighbors, members of the same village, classmates at a school, or even firms within a specific industry. In this case, it seems plausible that decisions may be inter-connected and thus mutually dependent rather than independent. Allowing for such interactions complicates inference and requires specialized treatment. A currently popular approach which allows for mutual dependence is known as **clustered dependence**, which assumes that observations are grouped into “clusters” (for example, schools). We will discuss clustering in more detail in Section 4.20.

### 4.3 Sample Mean

To start with the simplest setting, we first consider the intercept-only model

$$\begin{aligned} y_i &= \mu + e_i \\ \mathbb{E}(e_i) &= 0. \end{aligned}$$

which is equivalent to the regression model with  $k = 1$  and  $x_i = 1$ . In the intercept model,  $\mu = \mathbb{E}(y_i)$  is the mean of  $y_i$ . (See Exercise 2.15.) The least-squares estimator  $\hat{\mu} = \bar{y}$  equals the sample mean as shown in equation (3.10).

We now calculate the mean and variance of the estimator  $\bar{y}$ . Since the sample mean is a linear function of the observations, its expectation is simple to calculate

$$\mathbb{E}(\bar{y}) = \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n y_i\right) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(y_i) = \mu.$$

This shows that the expected value of the least-squares estimator (the sample mean) equals the projection coefficient (the population mean). An estimator with the property that its expectation equals the parameter it is estimating is called **unbiased**.

**Definition 4.3.1** An estimator  $\hat{\theta}$  for  $\theta$  is **unbiased** if  $\mathbb{E}(\hat{\theta}) = \theta$ .

We next calculate the variance of the estimator  $\bar{y}$  under Assumption 4.2.1. Making the substitution  $y_i = \mu + e_i$  we find

$$\bar{y} - \mu = \frac{1}{n} \sum_{i=1}^n e_i.$$

Then

$$\begin{aligned} \text{var}(\bar{y}) &= \mathbb{E}(\bar{y} - \mu)^2 \\ &= \mathbb{E}\left(\left(\frac{1}{n} \sum_{i=1}^n e_i\right) \left(\frac{1}{n} \sum_{j=1}^n e_j\right)\right) \\ &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}(e_i e_j) \\ &= \frac{1}{n^2} \sum_{i=1}^n \sigma^2 \\ &= \frac{1}{n} \sigma^2. \end{aligned}$$

The second-to-last inequality is because  $\mathbb{E}(e_i e_j) = \sigma^2$  for  $i = j$  yet  $\mathbb{E}(e_i e_j) = 0$  for  $i \neq j$  due to independence.

We have shown that  $\text{var}(\bar{y}) = \frac{1}{n} \sigma^2$ . This is the familiar formula for the variance of the sample mean.

## 4.4 Linear Regression Model

We now consider the linear regression model. Throughout this chapter we maintain the following.

**Assumption 4.4.1 *Linear Regression Model***

*The observations  $(y_i, \mathbf{x}_i)$  satisfy the linear regression equation*

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i \quad (4.1)$$

$$\mathbb{E}(e_i | \mathbf{x}_i) = 0. \quad (4.2)$$

*The variables have finite second moments*

$$\mathbb{E}(y_i^2) < \infty,$$

$$\mathbb{E} \|\mathbf{x}_i\|^2 < \infty,$$

*and an invertible design matrix*

$$\mathbf{Q}_{\mathbf{x}\mathbf{x}} = \mathbb{E}(\mathbf{x}_i \mathbf{x}_i') > 0.$$

We will consider both the general case of heteroskedastic regression, where the conditional variance

$$\mathbb{E}(e_i^2 | \mathbf{x}_i) = \sigma^2(\mathbf{x}_i) = \sigma_i^2$$

is unrestricted, and the specialized case of homoskedastic regression, where the conditional variance is constant. In the latter case we add the following assumption.

**Assumption 4.4.2 *Homoskedastic Linear Regression Model***

*In addition to Assumption 4.4.1,*

$$\mathbb{E}(e_i^2 | \mathbf{x}_i) = \sigma^2(\mathbf{x}_i) = \sigma^2 \quad (4.3)$$

*is independent of  $\mathbf{x}_i$ .*

## 4.5 Mean of Least-Squares Estimator

In this section we show that the OLS estimator is unbiased in the linear regression model. This calculation can be done using either summation notation or matrix notation. We will use both.

First take summation notation. Observe that under (4.1)-(4.2)

$$\mathbb{E}(y_i | \mathbf{X}) = \mathbb{E}(y_i | \mathbf{x}_i) = \mathbf{x}_i' \boldsymbol{\beta}. \quad (4.4)$$

The first equality states that the conditional expectation of  $y_i$  given  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  only depends on  $\mathbf{x}_i$ , since the observations are independent across  $i$ . The second equality is the assumption of a linear conditional mean.

Using definition (3.12), the conditioning theorem, the linearity of expectations, (4.4), and properties of the matrix inverse,

$$\begin{aligned}
\mathbb{E}(\hat{\beta} \mid \mathbf{X}) &= \mathbb{E}\left(\left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'\right)^{-1} \left(\sum_{i=1}^n \mathbf{x}_i y_i\right) \mid \mathbf{X}\right) \\
&= \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'\right)^{-1} \mathbb{E}\left(\left(\sum_{i=1}^n \mathbf{x}_i y_i\right) \mid \mathbf{X}\right) \\
&= \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'\right)^{-1} \sum_{i=1}^n \mathbb{E}(\mathbf{x}_i y_i \mid \mathbf{X}) \\
&= \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'\right)^{-1} \sum_{i=1}^n \mathbf{x}_i \mathbb{E}(y_i \mid \mathbf{X}) \\
&= \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'\right)^{-1} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \beta \\
&= \beta.
\end{aligned}$$

Now let's show the same result using matrix notation. (4.4) implies

$$\mathbb{E}(\mathbf{y} \mid \mathbf{X}) = \begin{pmatrix} \vdots \\ \mathbb{E}(y_i \mid \mathbf{X}) \\ \vdots \end{pmatrix} = \begin{pmatrix} \vdots \\ \mathbf{x}_i' \beta \\ \vdots \end{pmatrix} = \mathbf{X} \beta. \quad (4.5)$$

Similarly

$$\mathbb{E}(\mathbf{e} \mid \mathbf{X}) = \begin{pmatrix} \vdots \\ \mathbb{E}(e_i \mid \mathbf{X}) \\ \vdots \end{pmatrix} = \begin{pmatrix} \vdots \\ \mathbb{E}(e_i \mid \mathbf{x}_i) \\ \vdots \end{pmatrix} = \mathbf{0}. \quad (4.6)$$

Using definition (3.22), the conditioning theorem, the linearity of expectations, (4.5), and the properties of the matrix inverse,

$$\begin{aligned}
\mathbb{E}(\hat{\beta} \mid \mathbf{X}) &= \mathbb{E}\left((\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} \mid \mathbf{X}\right) \\
&= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \mathbb{E}(\mathbf{y} \mid \mathbf{X}) \\
&= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \mathbf{X} \beta \\
&= \beta.
\end{aligned}$$

At the risk of belaboring the derivation, another way to calculate the same result is as follows. Insert  $\mathbf{y} = \mathbf{X}\beta + \mathbf{e}$  into the formula (3.22) for  $\hat{\beta}$  to obtain

$$\begin{aligned}
\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'(\mathbf{X}\beta + \mathbf{e})) \\
&= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X}\beta + (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{e}) \\
&= \beta + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{e}.
\end{aligned} \quad (4.7)$$

This is a useful linear decomposition of the estimator  $\hat{\beta}$  into the true parameter  $\beta$  and the stochastic component  $(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{e}$ . Once again, we can calculate that

$$\begin{aligned}
\mathbb{E}(\hat{\beta} - \beta \mid \mathbf{X}) &= \mathbb{E}\left((\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{e} \mid \mathbf{X}\right) \\
&= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \mathbb{E}(\mathbf{e} \mid \mathbf{X}) \\
&= \mathbf{0}.
\end{aligned}$$

Regardless of the method, we have shown that  $\mathbb{E}(\hat{\beta} | \mathbf{X}) = \beta$ .

We have shown the following theorem.

**Theorem 4.5.1 Mean of Least-Squares Estimator**

*In the linear regression model (Assumption 4.4.1) and i.i.d. sampling (Assumption 4.2.1)*

$$\mathbb{E}(\hat{\beta} | \mathbf{X}) = \beta \quad (4.8)$$

Equation (4.8) says that the estimator  $\hat{\beta}$  is unbiased for  $\beta$ , conditional on  $\mathbf{X}$ . This means that the conditional distribution of  $\hat{\beta}$  is centered at  $\beta$ . By “conditional on  $\mathbf{X}$ ” this means that the distribution is unbiased (centered at  $\beta$ ) for any realization of the regressor matrix  $\mathbf{X}$ . In conditional models, we simply refer to this as saying “ $\hat{\beta}$  is unbiased for  $\beta$ ”.

Strictly speaking, “unbiasedness” is a property of the unconditional distribution. Assuming the unconditional mean is well defined, that is  $\mathbb{E}\|\hat{\beta}\| < \infty$ , then applying the law of iterated expectations, we find that the unconditional mean of  $\hat{\beta}$  is also  $\beta$

$$\mathbb{E}(\hat{\beta}) = \mathbb{E}(\mathbb{E}(\hat{\beta} | \mathbf{X})) = \beta. \quad (4.9)$$

## 4.6 Variance of Least Squares Estimator

In this section we calculate the conditional variance of the OLS estimator.

For any  $r \times 1$  random vector  $\mathbf{Z}$  define the  $r \times r$  covariance matrix

$$\begin{aligned} \text{var}(\mathbf{Z}) &= \mathbb{E}((\mathbf{Z} - \mathbb{E}(\mathbf{Z}))(\mathbf{Z} - \mathbb{E}(\mathbf{Z}))') \\ &= \mathbb{E}(\mathbf{Z}\mathbf{Z}') - (\mathbb{E}(\mathbf{Z}))(\mathbb{E}(\mathbf{Z}))' \end{aligned}$$

and for any pair  $(\mathbf{Z}, \mathbf{X})$  define the conditional covariance matrix

$$\text{var}(\mathbf{Z} | \mathbf{X}) = \mathbb{E}((\mathbf{Z} - \mathbb{E}(\mathbf{Z} | \mathbf{X}))(\mathbf{Z} - \mathbb{E}(\mathbf{Z} | \mathbf{X}))' | \mathbf{X}).$$

We define

$$\mathbf{V}_{\hat{\beta}} \stackrel{\text{def}}{=} \text{var}(\hat{\beta} | \mathbf{X})$$

as the conditional covariance matrix of the regression coefficient estimates. We now derive its form.

The conditional covariance matrix of the  $n \times 1$  regression error  $\mathbf{e}$  is the  $n \times n$  matrix

$$\text{var}(\mathbf{e} | \mathbf{X}) = \mathbb{E}(\mathbf{e}\mathbf{e}' | \mathbf{X}) \stackrel{\text{def}}{=} \mathbf{D}.$$

The  $i^{\text{th}}$  diagonal element of  $\mathbf{D}$  is

$$\mathbb{E}(e_i^2 | \mathbf{X}) = \mathbb{E}(e_i^2 | \mathbf{x}_i) = \sigma_i^2$$

while the  $ij^{\text{th}}$  off-diagonal element of  $\mathbf{D}$  is

$$\mathbb{E}(e_i e_j | \mathbf{X}) = \mathbb{E}(e_i | \mathbf{x}_i) \mathbb{E}(e_j | \mathbf{x}_j) = 0.$$

where the first equality uses independence of the observations (Assumption 1.5.2) and the second is (4.2). Thus  $\mathbf{D}$  is a diagonal matrix with  $i^{th}$  diagonal element  $\sigma_i^2$ :

$$\mathbf{D} = \text{diag}(\sigma_1^2, \dots, \sigma_n^2) = \begin{pmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{pmatrix}. \quad (4.10)$$

In the special case of the linear homoskedastic regression model (4.3), then

$$\mathbb{E}(e_i^2 | \mathbf{x}_i) = \sigma_i^2 = \sigma^2$$

and we have the simplification

$$\mathbf{D} = \mathbf{I}_n \sigma^2.$$

In general, however,  $\mathbf{D}$  need not necessarily take this simplified form.

For any  $n \times r$  matrix  $\mathbf{A} = \mathbf{A}(\mathbf{X})$ ,

$$\text{var}(\mathbf{A}'\mathbf{y} | \mathbf{X}) = \text{var}(\mathbf{A}'\mathbf{e} | \mathbf{X}) = \mathbf{A}'\mathbf{D}\mathbf{A}. \quad (4.11)$$

In particular, we can write  $\hat{\boldsymbol{\beta}} = \mathbf{A}'\mathbf{y}$  where  $\mathbf{A} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$  and thus

$$\mathbf{V}_{\hat{\boldsymbol{\beta}}} = \text{var}(\hat{\boldsymbol{\beta}} | \mathbf{X}) = \mathbf{A}'\mathbf{D}\mathbf{A} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{D}\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1}.$$

It is useful to note that

$$\mathbf{X}'\mathbf{D}\mathbf{X} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \sigma_i^2,$$

a weighted version of  $\mathbf{X}'\mathbf{X}$ .

In the special case of the linear homoskedastic regression model,  $\mathbf{D} = \mathbf{I}_n \sigma^2$ , so  $\mathbf{X}'\mathbf{D}\mathbf{X} = \mathbf{X}'\mathbf{X} \sigma^2$ , and the variance matrix simplifies to

$$\mathbf{V}_{\hat{\boldsymbol{\beta}}} = (\mathbf{X}'\mathbf{X})^{-1} \sigma^2.$$

**Theorem 4.6.1 Variance of Least-Squares Estimator**

*In the linear regression model (Assumption 4.4.1) and i.i.d. sampling (Assumption 4.2.1)*

$$\begin{aligned} \mathbf{V}_{\hat{\boldsymbol{\beta}}} &= \text{var}(\hat{\boldsymbol{\beta}} | \mathbf{X}) \\ &= (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{D}\mathbf{X}) (\mathbf{X}'\mathbf{X})^{-1} \end{aligned} \quad (4.12)$$

*where  $\mathbf{D}$  is defined in (4.10).*

*In the homoskedastic linear regression model (Assumption 4.4.2) and i.i.d. sampling (Assumption 4.2.1)*

$$\mathbf{V}_{\hat{\boldsymbol{\beta}}} = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}. \quad (4.13)$$



## 4.7 Gauss-Markov Theorem

Now consider the class of estimators of  $\beta$  which are linear functions of the vector  $\mathbf{y}$ , and thus can be written as

$$\tilde{\beta} = \mathbf{A}'\mathbf{y}$$

where  $\mathbf{A}$  is an  $n \times k$  function of  $\mathbf{X}$ . As noted before, the least-squares estimator is the special case obtained by setting  $\mathbf{A} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$ . What is the best choice of  $\mathbf{A}$ ? The Gauss-Markov theorem, which we now present, says that the least-squares estimator is the best choice among linear unbiased estimators when the errors are homoskedastic, in the sense that the least-squares estimator has the smallest variance among all unbiased linear estimators.

To see this, since  $\mathbb{E}(\mathbf{y} | \mathbf{X}) = \mathbf{X}\beta$ , then for any linear estimator  $\tilde{\beta} = \mathbf{A}'\mathbf{y}$  we have

$$\mathbb{E}(\tilde{\beta} | \mathbf{X}) = \mathbf{A}'\mathbb{E}(\mathbf{y} | \mathbf{X}) = \mathbf{A}'\mathbf{X}\beta,$$

so  $\tilde{\beta}$  is unbiased if (and only if)  $\mathbf{A}'\mathbf{X} = \mathbf{I}_k$ . Furthermore, we saw in (4.11) that

$$\text{var}(\tilde{\beta} | \mathbf{X}) = \text{var}(\mathbf{A}'\mathbf{y} | \mathbf{X}) = \mathbf{A}'\mathbf{D}\mathbf{A} = \mathbf{A}'\mathbf{A}\sigma^2$$

the last equality using the homoskedasticity assumption  $\mathbf{D} = \mathbf{I}_n\sigma^2$ . The “best” unbiased linear estimator is obtained by finding the matrix  $\mathbf{A}_0$  satisfying  $\mathbf{A}_0'\mathbf{X} = \mathbf{I}_k$  such that  $\mathbf{A}_0'\mathbf{A}_0$  is minimized in the positive definite sense, in that for any other matrix  $\mathbf{A}$  satisfying  $\mathbf{A}'\mathbf{X} = \mathbf{I}_k$ , then  $\mathbf{A}'\mathbf{A} - \mathbf{A}_0'\mathbf{A}_0$  is positive semi-definite.

**Theorem 4.7.1 Gauss-Markov.** *In the homoskedastic linear regression model (Assumption 4.4.2) and i.i.d. sampling (Assumption 4.2.1), if  $\tilde{\beta}$  is a linear unbiased estimator of  $\beta$  then*

$$\text{var}(\tilde{\beta} | \mathbf{X}) \geq \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}.$$

The Gauss-Markov theorem provides a lower bound on the variance matrix of unbiased linear estimators under the assumption of homoskedasticity. It says that no unbiased linear estimator can have a variance matrix smaller (in the positive definite sense) than  $\sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$ . Since the variance of the OLS estimator is exactly equal to this bound, this means that the OLS estimator is efficient in the class of linear unbiased estimator. This gives rise to the description of OLS as BLUE, standing for “best linear unbiased estimator”. This is an efficiency justification for the least-squares estimator. The justification is limited because the class of models is restricted to homoskedastic linear regression and the class of potential estimators is restricted to linear unbiased estimators. This latter restriction is particularly unsatisfactory as the theorem leaves open the possibility that a non-linear or biased estimator could have lower mean squared error than the least-squares estimator.

We give a proof of the Gauss-Markov theorem below.

---

**Proof of Theorem 4.7.1.1.** Let  $\mathbf{A}$  be any  $n \times k$  function of  $\mathbf{X}$  such that  $\mathbf{A}'\mathbf{X} = \mathbf{I}_k$ . The variance of the least-squares estimator is  $(\mathbf{X}'\mathbf{X})^{-1}\sigma^2$  and that of  $\mathbf{A}'\mathbf{y}$  is  $\mathbf{A}'\mathbf{A}\sigma^2$ . It is sufficient to show

that the difference  $\mathbf{A}'\mathbf{A} - (\mathbf{X}'\mathbf{X})^{-1}$  is positive semi-definite. Set  $\mathbf{C} = \mathbf{A} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$ . Note that  $\mathbf{X}'\mathbf{C} = \mathbf{0}$ . Then we calculate that

$$\begin{aligned}\mathbf{A}'\mathbf{A} - (\mathbf{X}'\mathbf{X})^{-1} &= \left(\mathbf{C} + \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\right)' \left(\mathbf{C} + \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\right) - (\mathbf{X}'\mathbf{X})^{-1} \\ &= \mathbf{C}'\mathbf{C} + \mathbf{C}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{C} \\ &\quad + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} - (\mathbf{X}'\mathbf{X})^{-1} \\ &= \mathbf{C}'\mathbf{C}.\end{aligned}$$

The matrix  $\mathbf{C}'\mathbf{C}$  is positive semi-definite (see Appendix A.9) as required.

---

## 4.8 Generalized Least Squares

Take the linear regression model in matrix format

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}. \quad (4.14)$$

Consider a generalized situation where the observation errors are possibly correlated and/or heteroskedastic. Specifically, suppose that

$$\mathbb{E}(\mathbf{e} \mid \mathbf{X}) = \mathbf{0} \quad (4.15)$$

$$\text{var}(\mathbf{e} \mid \mathbf{X}) = \boldsymbol{\Omega} \quad (4.16)$$

for some  $n \times n$  covariance matrix  $\boldsymbol{\Omega}$ , possibly a function of  $\mathbf{X}$ . This includes the iid sampling framework where  $\boldsymbol{\Omega} = \mathbf{D}$  but allows for non-diagonal covariance matrices as well.

Under these assumptions, by similar arguments we can calculate the mean and variance of the OLS estimator:

$$\mathbb{E}(\hat{\boldsymbol{\beta}} \mid \mathbf{X}) = \boldsymbol{\beta} \quad (4.17)$$

$$\text{var}(\hat{\boldsymbol{\beta}} \mid \mathbf{X}) = (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\boldsymbol{\Omega}\mathbf{X}) (\mathbf{X}'\mathbf{X})^{-1} \quad (4.18)$$

(see Exercise 4.5).

We have an analog of the Gauss-Markov Theorem.

**Theorem 4.8.1** *If (4.15)-(4.16) hold and if  $\tilde{\boldsymbol{\beta}}$  is a linear unbiased estimator of  $\boldsymbol{\beta}$  then*

$$\text{var}(\tilde{\boldsymbol{\beta}} \mid \mathbf{X}) \geq (\mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{X})^{-1}.$$

We leave the proof for Exercise 4.6.

The theorem provides a lower bound on the variance matrix of unbiased linear estimators. The bound is different from the variance matrix of the OLS estimator except when  $\boldsymbol{\Omega} = \mathbf{I}_n\sigma^2$ . This suggests that we may be able to improve on the OLS estimator.

This is indeed the case when  $\boldsymbol{\Omega}$  is known up to scale. That is, suppose that  $\boldsymbol{\Omega} = c^2\boldsymbol{\Sigma}$  where  $c^2 > 0$  is real and  $\boldsymbol{\Sigma}$  is  $n \times n$  and known. Take the linear model (4.14) and pre-multiply by  $\boldsymbol{\Sigma}^{-1/2}$ . This produces the equation

$$\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\boldsymbol{\beta} + \tilde{\mathbf{e}}$$

where  $\tilde{\mathbf{y}} = \boldsymbol{\Sigma}^{-1/2}\mathbf{y}$ ,  $\tilde{\mathbf{X}} = \boldsymbol{\Sigma}^{-1/2}\mathbf{X}$ , and  $\tilde{\mathbf{e}} = \boldsymbol{\Sigma}^{-1/2}\mathbf{e}$ . Consider OLS estimation of  $\boldsymbol{\beta}$  in this equation

$$\begin{aligned}\tilde{\boldsymbol{\beta}}_{\text{ols}} &= (\tilde{\mathbf{X}}'\tilde{\mathbf{X}})^{-1}\tilde{\mathbf{X}}'\tilde{\mathbf{y}} \\ &= \left( \left( \boldsymbol{\Sigma}^{-1/2}\mathbf{X} \right)' \left( \boldsymbol{\Sigma}^{-1/2}\mathbf{X} \right) \right)^{-1} \left( \boldsymbol{\Sigma}^{-1/2}\mathbf{X} \right)' \left( \boldsymbol{\Sigma}^{-1/2}\mathbf{y} \right) \\ &= (\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{y}.\end{aligned}\tag{4.19}$$

This is called the **Generalized Least Squares** (GLS) estimator of  $\boldsymbol{\beta}$ .

You can calculate that

$$\mathbb{E}(\tilde{\boldsymbol{\beta}}_{\text{ols}} \mid \mathbf{X}) = \boldsymbol{\beta}\tag{4.20}$$

$$\text{var}(\tilde{\boldsymbol{\beta}}_{\text{ols}} \mid \mathbf{X}) = (\mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{X})^{-1}.\tag{4.21}$$

This shows that the GLS estimator is unbiased, and has a covariance matrix which equals the lower bound from Theorem 4.8.1. This shows that the lower bound is sharp when  $\boldsymbol{\Sigma}$  is known and the GLS is efficient in the class of linear unbiased estimators.

In the linear regression model with independent observations and known conditional variances, where  $\boldsymbol{\Omega} = \boldsymbol{\Sigma} = \mathbf{D} = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$ , the GLS estimator takes the form

$$\begin{aligned}\tilde{\boldsymbol{\beta}}_{\text{ols}} &= (\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{D}^{-1}\mathbf{y} \\ &= \left( \sum_{i=1}^n \sigma_i^{-2} \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left( \sum_{i=1}^n \sigma_i^{-2} \mathbf{x}_i y_i \right).\end{aligned}$$

In practice, the covariance matrix  $\boldsymbol{\Omega}$  is unknown, so the GLS estimator as presented here is not feasible. However, the form of the GLS estimator motivates feasible versions, effectively by replacing  $\boldsymbol{\Omega}$  with an estimate. We return to this issue in Section 20.2.

## 4.9 Residuals

What are some properties of the residuals  $\hat{e}_i = y_i - \mathbf{x}_i'\hat{\boldsymbol{\beta}}$  and prediction errors  $\tilde{e}_i = y_i - \mathbf{x}_i'\hat{\boldsymbol{\beta}}_{(-i)}$ , at least in the context of the linear regression model?

Recall from (3.31) that we can write the residuals in vector notation as

$$\hat{\mathbf{e}} = \mathbf{M}\mathbf{e}$$

where  $\mathbf{M} = \mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$  is the orthogonal projection matrix. Using the properties of conditional expectation

$$\mathbb{E}(\hat{\mathbf{e}} \mid \mathbf{X}) = \mathbb{E}(\mathbf{M}\mathbf{e} \mid \mathbf{X}) = \mathbf{M}\mathbb{E}(\mathbf{e} \mid \mathbf{X}) = \mathbf{0}$$

and

$$\text{var}(\hat{\mathbf{e}} \mid \mathbf{X}) = \text{var}(\mathbf{M}\mathbf{e} \mid \mathbf{X}) = \mathbf{M} \text{var}(\mathbf{e} \mid \mathbf{X}) \mathbf{M} = \mathbf{M}\mathbf{D}\mathbf{M}\tag{4.22}$$

where  $\mathbf{D}$  is defined in (4.10).

We can simplify this expression under the assumption of conditional homoskedasticity

$$\mathbb{E}(e_i^2 \mid \mathbf{x}_i) = \sigma^2.$$

In this case (4.22) simplifies to

$$\text{var}(\hat{\mathbf{e}} \mid \mathbf{X}) = \mathbf{M}\sigma^2.\tag{4.23}$$

In particular, for a single observation  $i$ , we can find the (conditional) variance of  $\hat{e}_i$  by taking the  $i^{th}$  diagonal element of (4.23). Since the  $i^{th}$  diagonal element of  $\mathbf{M}$  is  $1 - h_{ii}$  as defined in (3.25) we obtain

$$\text{var}(\hat{e}_i | \mathbf{X}) = \mathbb{E}(\hat{e}_i^2 | \mathbf{X}) = (1 - h_{ii})\sigma^2. \quad (4.24)$$

As this variance is a function of  $h_{ii}$  and hence  $\mathbf{x}_i$ , the residuals  $\hat{e}_i$  are heteroskedastic even if the errors  $e_i$  are homoskedastic. Notice as well that this implies  $\hat{e}_i^2$  is a biased estimator of  $\sigma^2$ .

Similarly, recall from (3.46) that the prediction errors  $\tilde{e}_i = (1 - h_{ii})^{-1}\hat{e}_i$  can be written in vector notation as  $\tilde{\mathbf{e}} = \mathbf{M}^*\hat{\mathbf{e}}$  where  $\mathbf{M}^*$  is a diagonal matrix with  $i^{th}$  diagonal element  $(1 - h_{ii})^{-1}$ . Thus  $\tilde{\mathbf{e}} = \mathbf{M}^*\mathbf{M}\mathbf{e}$ . We can calculate that

$$\mathbb{E}(\tilde{\mathbf{e}} | \mathbf{X}) = \mathbf{M}^*\mathbf{M}\mathbb{E}(\mathbf{e} | \mathbf{X}) = \mathbf{0}$$

and

$$\text{var}(\tilde{\mathbf{e}} | \mathbf{X}) = \mathbf{M}^*\mathbf{M} \text{var}(\mathbf{e} | \mathbf{X}) \mathbf{M}\mathbf{M}^* = \mathbf{M}^*\mathbf{M}\mathbf{D}\mathbf{M}\mathbf{M}^*$$

which simplifies under homoskedasticity to

$$\begin{aligned} \text{var}(\tilde{\mathbf{e}} | \mathbf{X}) &= \mathbf{M}^*\mathbf{M}\mathbf{M}\mathbf{M}^*\sigma^2 \\ &= \mathbf{M}^*\mathbf{M}\mathbf{M}^*\sigma^2. \end{aligned}$$

The variance of the  $i^{th}$  prediction error is then

$$\begin{aligned} \text{var}(\tilde{e}_i | \mathbf{X}) &= \mathbb{E}(\tilde{e}_i^2 | \mathbf{X}) \\ &= (1 - h_{ii})^{-1} (1 - h_{ii}) (1 - h_{ii})^{-1} \sigma^2 \\ &= (1 - h_{ii})^{-1} \sigma^2. \end{aligned}$$

A residual with constant conditional variance can be obtained by rescaling. The **standardized residuals** are

$$\bar{e}_i = (1 - h_{ii})^{-1/2} \hat{e}_i, \quad (4.25)$$

and in vector notation

$$\bar{\mathbf{e}} = (\bar{e}_1, \dots, \bar{e}_n)' = \mathbf{M}^{*1/2}\mathbf{M}\mathbf{e}. \quad (4.26)$$

From our above calculations, under homoskedasticity,

$$\text{var}(\bar{\mathbf{e}} | \mathbf{X}) = \mathbf{M}^{*1/2}\mathbf{M}\mathbf{M}^{*1/2}\sigma^2$$

and

$$\text{var}(\bar{e}_i | \mathbf{X}) = \mathbb{E}(\bar{e}_i^2 | \mathbf{X}) = \sigma^2 \quad (4.27)$$

and thus these standardized residuals have the same bias and variance as the original errors when the latter are homoskedastic.

## 4.10 Estimation of Error Variance

The error variance  $\sigma^2 = \mathbb{E}(e_i^2)$  can be a parameter of interest even in a heteroskedastic regression or a projection model.  $\sigma^2$  measures the variation in the “unexplained” part of the regression. Its method of moments estimator (MME) is the sample average of the squared residuals:

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{e}_i^2.$$

In the linear regression model we can calculate the mean of  $\hat{\sigma}^2$ . From (3.35) and the properties of the trace operator, observe that

$$\hat{\sigma}^2 = \frac{1}{n} \mathbf{e}'\mathbf{M}\mathbf{e} = \frac{1}{n} \text{tr}(\mathbf{e}'\mathbf{M}\mathbf{e}) = \frac{1}{n} \text{tr}(\mathbf{M}\mathbf{e}\mathbf{e}').$$

Then

$$\begin{aligned}\mathbb{E}(\hat{\sigma}^2 \mid \mathbf{X}) &= \frac{1}{n} \operatorname{tr}(\mathbb{E}(\mathbf{M} \mathbf{e} \mathbf{e}' \mid \mathbf{X})) \\ &= \frac{1}{n} \operatorname{tr}(\mathbf{M} \mathbb{E}(\mathbf{e} \mathbf{e}' \mid \mathbf{X})) \\ &= \frac{1}{n} \operatorname{tr}(\mathbf{M} \mathbf{D}).\end{aligned}\tag{4.28}$$

Adding the assumption of conditional homoskedasticity  $\mathbb{E}(e_i^2 \mid \mathbf{x}_i) = \sigma^2$ , so that  $\mathbf{D} = \mathbf{I}_n \sigma^2$ , then (4.28) simplifies to

$$\begin{aligned}\mathbb{E}(\hat{\sigma}^2 \mid \mathbf{X}) &= \frac{1}{n} \operatorname{tr}(\mathbf{M} \sigma^2) \\ &= \sigma^2 \left( \frac{n-k}{n} \right),\end{aligned}$$

the final equality by (3.29). This calculation shows that  $\hat{\sigma}^2$  is biased towards zero. The order of the bias depends on  $k/n$ , the ratio of the number of estimated coefficients to the sample size.

Another way to see this is to use (4.24). Note that

$$\begin{aligned}\mathbb{E}(\hat{\sigma}^2 \mid \mathbf{X}) &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}(\hat{e}_i^2 \mid \mathbf{X}) \\ &= \frac{1}{n} \sum_{i=1}^n (1 - h_{ii}) \sigma^2 \\ &= \left( \frac{n-k}{n} \right) \sigma^2\end{aligned}\tag{4.29}$$

the last equality using Theorem 3.11.1.

Since the bias takes a scale form, a classic method to obtain an unbiased estimator is by rescaling the estimator. Define

$$s^2 = \frac{1}{n-k} \sum_{i=1}^n \hat{e}_i^2.\tag{4.30}$$

By the above calculation,

$$\mathbb{E}(s^2 \mid \mathbf{X}) = \sigma^2\tag{4.31}$$

and

$$\mathbb{E}(s^2) = \sigma^2.$$

Hence the estimator  $s^2$  is unbiased for  $\sigma^2$ . Consequently,  $s^2$  is known as the “bias-corrected estimator” for  $\sigma^2$  and in empirical practice  $s^2$  is the most widely used estimator for  $\sigma^2$ .

Interestingly, this is not the only method to construct an unbiased estimator for  $\sigma^2$ . An estimator constructed with the standardized residuals  $\bar{e}_i$  from (4.25) is

$$\bar{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \bar{e}_i^2 = \frac{1}{n} \sum_{i=1}^n (1 - h_{ii})^{-1} \hat{e}_i^2.\tag{4.32}$$

You can show (see Exercise 4.9) that

$$\mathbb{E}(\bar{\sigma}^2 \mid \mathbf{X}) = \sigma^2\tag{4.33}$$

and thus  $\bar{\sigma}^2$  is unbiased for  $\sigma^2$  (in the homoskedastic linear regression model).

When  $k/n$  is small (typically, this occurs when  $n$  is large), the estimators  $\hat{\sigma}^2$ ,  $s^2$  and  $\bar{\sigma}^2$  are likely to be close. However, if not then  $s^2$  and  $\bar{\sigma}^2$  are generally preferred to  $\hat{\sigma}^2$ . Consequently it is best to use one of the bias-corrected variance estimators in applications.

## 4.11 Mean-Square Forecast Error

A major purpose of estimated regressions is to predict out-of-sample values. Consider an out-of-sample observation  $(y_{n+1}, \mathbf{x}_{n+1})$  where  $\mathbf{x}_{n+1}$  is observed but not  $y_{n+1}$ . Given the coefficient estimate  $\hat{\beta}$  the standard point estimate of  $\mathbb{E}(y_{n+1} | \mathbf{x}_{n+1}) = \mathbf{x}_{n+1}'\beta$  is  $\tilde{y}_{n+1} = \mathbf{x}_{n+1}'\hat{\beta}$ . The forecast error is the difference between the actual value  $y_{n+1}$  and the point forecast  $\tilde{y}_{n+1}$ . This is the forecast error  $\tilde{e}_{n+1} = y_{n+1} - \tilde{y}_{n+1}$ . The mean-squared forecast error (MSFE) is its expected squared value

$$MSFE_n = \mathbb{E}(\tilde{e}_{n+1}^2).$$

In the linear regression model,  $\tilde{e}_{n+1} = e_{n+1} - \mathbf{x}_{n+1}'(\hat{\beta} - \beta)$ , so

$$\begin{aligned} MSFE_n &= \mathbb{E}(e_{n+1}^2) - 2\mathbb{E}\left(e_{n+1}\mathbf{x}_{n+1}'(\hat{\beta} - \beta)\right) \\ &\quad + \mathbb{E}\left(\mathbf{x}_{n+1}'(\hat{\beta} - \beta)(\hat{\beta} - \beta)' \mathbf{x}_{n+1}\right). \end{aligned} \quad (4.34)$$

The first term in (4.34) is  $\sigma^2$ . The second term in (4.34) is zero since  $e_{n+1}\mathbf{x}_{n+1}'$  is independent of  $\hat{\beta} - \beta$  and both are mean zero. Using the properties of the trace operator, the third term in (4.34) is

$$\begin{aligned} &\text{tr}\left(\mathbb{E}(\mathbf{x}_{n+1}\mathbf{x}_{n+1}') \mathbb{E}\left((\hat{\beta} - \beta)(\hat{\beta} - \beta)'\right)\right) \\ &= \text{tr}\left(\mathbb{E}(\mathbf{x}_{n+1}\mathbf{x}_{n+1}') \mathbb{E}\left(\mathbb{E}\left((\hat{\beta} - \beta)(\hat{\beta} - \beta)' | \mathbf{X}\right)\right)\right) \\ &= \text{tr}\left(\mathbb{E}(\mathbf{x}_{n+1}\mathbf{x}_{n+1}') \mathbb{E}(\mathbf{V}_{\hat{\beta}})\right) \\ &= \mathbb{E}\text{tr}\left((\mathbf{x}_{n+1}\mathbf{x}_{n+1}') \mathbf{V}_{\hat{\beta}}\right) \\ &= \mathbb{E}\left(\mathbf{x}_{n+1}' \mathbf{V}_{\hat{\beta}} \mathbf{x}_{n+1}\right) \end{aligned} \quad (4.35)$$

where we use the fact that  $\mathbf{x}_{n+1}$  is independent of  $\hat{\beta}$ , the definition  $\mathbf{V}_{\hat{\beta}} = \mathbb{E}\left((\hat{\beta} - \beta)(\hat{\beta} - \beta)' | \mathbf{X}\right)$  and the fact that  $\mathbf{x}_{n+1}$  is independent of  $\mathbf{V}_{\hat{\beta}}$ . Thus

$$MSFE_n = \sigma^2 + \mathbb{E}\left(\mathbf{x}_{n+1}' \mathbf{V}_{\hat{\beta}} \mathbf{x}_{n+1}\right).$$

Under conditional homoskedasticity, this simplifies to

$$MSFE_n = \sigma^2 \left(1 + \mathbb{E}\left(\mathbf{x}_{n+1}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_{n+1}\right)\right).$$

A simple estimator for the MSFE is obtained by averaging the squared prediction errors (3.47)

$$\tilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \tilde{e}_i^2$$

where  $\tilde{e}_i = y_i - \mathbf{x}_i' \hat{\beta}_{(-i)} = \hat{e}_i(1 - h_{ii})^{-1}$ . Indeed, we can calculate that

$$\begin{aligned} \mathbb{E}(\tilde{\sigma}^2) &= \mathbb{E}(\tilde{e}_i^2) \\ &= \mathbb{E}\left(e_i - \mathbf{x}_i'(\hat{\beta}_{(-i)} - \beta)\right)^2 \\ &= \sigma^2 + \mathbb{E}\left(\mathbf{x}_i'(\hat{\beta}_{(-i)} - \beta)(\hat{\beta}_{(-i)} - \beta)' \mathbf{x}_i\right). \end{aligned}$$

By a similar calculation as in (4.35) we find

$$\mathbb{E}(\tilde{\sigma}^2) = \sigma^2 + \mathbb{E}\left(\mathbf{x}'_i \mathbf{V}_{\hat{\beta}_{(-i)}} \mathbf{x}_i\right) = MSFE_{n-1}.$$

This is the MSFE based on a sample of size  $n-1$ , rather than size  $n$ . The difference arises because the in-sample prediction errors  $\tilde{e}_i$  for  $i \leq n$  are calculated using an effective sample size of  $n-1$ , while the out-of sample prediction error  $\tilde{e}_{n+1}$  is calculated from a sample with the full  $n$  observations. Unless  $n$  is very small we should expect  $MSFE_{n-1}$  (the MSFE based on  $n-1$  observations) to be close to  $MSFE_n$  (the MSFE based on  $n$  observations). Thus  $\tilde{\sigma}^2$  is a reasonable estimator for  $MSFE_n$ .

**Theorem 4.11.1 MSFE**

*In the linear regression model (Assumption 4.4.1) and i.i.d. sampling (Assumption 4.2.1)*

$$MSFE_n = \mathbb{E}(\tilde{e}_{n+1}^2) = \sigma^2 + \mathbb{E}\left(\mathbf{x}'_{n+1} \mathbf{V}_{\hat{\beta}} \mathbf{x}_{n+1}\right)$$

where  $\mathbf{V}_{\hat{\beta}} = \text{var}(\hat{\beta} \mid \mathbf{X})$ . Furthermore,  $\tilde{\sigma}^2$  defined in (3.47) is an unbiased estimator of  $MSFE_{n-1}$ :

$$\mathbb{E}(\tilde{\sigma}^2) = MSFE_{n-1}.$$

## 4.12 Covariance Matrix Estimation Under Homoskedasticity

For inference, we need an estimate of the covariance matrix  $\mathbf{V}_{\hat{\beta}}$  of the least-squares estimator. In this section we consider the homoskedastic regression model (Assumption 4.4.2).

Under homoskedasticity, the covariance matrix takes the relatively simple form

$$\mathbf{V}_{\hat{\beta}}^0 = (\mathbf{X}'\mathbf{X})^{-1} \sigma^2$$

which is known up to the unknown scale  $\sigma^2$ . In Section 4.10 we discussed three estimators of  $\sigma^2$ . The most commonly used choice is  $s^2$ , leading to the classic covariance matrix estimator

$$\hat{\mathbf{V}}_{\hat{\beta}}^0 = (\mathbf{X}'\mathbf{X})^{-1} s^2. \quad (4.36)$$

Since  $s^2$  is conditionally unbiased for  $\sigma^2$ , it is simple to calculate that  $\hat{\mathbf{V}}_{\hat{\beta}}^0$  is conditionally unbiased for  $\mathbf{V}_{\hat{\beta}}$  under the assumption of homoskedasticity:

$$\begin{aligned} \mathbb{E}(\hat{\mathbf{V}}_{\hat{\beta}}^0 \mid \mathbf{X}) &= (\mathbf{X}'\mathbf{X})^{-1} \mathbb{E}(s^2 \mid \mathbf{X}) \\ &= (\mathbf{X}'\mathbf{X})^{-1} \sigma^2 \\ &= \mathbf{V}_{\hat{\beta}}. \end{aligned}$$

This was the dominant covariance matrix estimator in applied econometrics for many years, and is still the default method in most regression packages. For example, Stata uses the covariance matrix estimator (4.36) by default in linear regression unless an alternative is specified.

If the estimator (4.36) is used, but the regression error is heteroskedastic, it is possible for  $\hat{\mathbf{V}}_{\hat{\beta}}^0$  to be quite biased for the correct covariance matrix  $\mathbf{V}_{\hat{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{D}\mathbf{X}) (\mathbf{X}'\mathbf{X})^{-1}$ . For example,

suppose  $k = 1$  and  $\sigma_i^2 = x_i^2$  with  $\mathbb{E}(x_i) = 0$ . The ratio of the true variance of the least-squares estimator to the expectation of the variance estimator is

$$\frac{\mathbf{V}_{\hat{\beta}}}{\mathbb{E}(\hat{\mathbf{V}}_{\hat{\beta}}^0 | \mathbf{X})} = \frac{\sum_{i=1}^n x_i^4}{\sigma^2 \sum_{i=1}^n x_i^2} \simeq \frac{\mathbb{E}(x_i^4)}{(\mathbb{E}(x_i^2))^2} \stackrel{def}{=} \kappa.$$

(Notice that we use the fact that  $\sigma_i^2 = x_i^2$  implies  $\sigma^2 = \mathbb{E}(\sigma_i^2) = \mathbb{E}(x_i^2)$ .) The constant  $\kappa$  is the standardized fourth moment (or kurtosis) of the regressor  $x_i$ , and can be any number greater than one. For example, if  $x_i \sim N(0, \sigma^2)$  then  $\kappa = 3$ , so the true variance  $\mathbf{V}_{\hat{\beta}}$  is three times larger than the expected homoskedastic estimator  $\hat{\mathbf{V}}_{\hat{\beta}}^0$ . But  $\kappa$  can be much larger. Suppose, for example, that  $x_i \sim \chi_1^2 - 1$ . In this case  $\kappa = 15$ , so that the true variance  $\mathbf{V}_{\hat{\beta}}$  is fifteen times larger than the expected homoskedastic estimator  $\hat{\mathbf{V}}_{\hat{\beta}}^0$ . While this is an extreme and constructed example, the point is that the classic covariance matrix estimator (4.36) may be quite biased when the homoskedasticity assumption fails.

### 4.13 Covariance Matrix Estimation Under Heteroskedasticity

In the previous section we showed that the classic covariance matrix estimator can be highly biased if homoskedasticity fails. In this section we show how to construct covariance matrix estimators which do not require homoskedasticity.

Recall that the general form for the covariance matrix is

$$\mathbf{V}_{\hat{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{D}\mathbf{X}) (\mathbf{X}'\mathbf{X})^{-1}.$$

This depends on the unknown matrix  $\mathbf{D}$  which we can write as

$$\begin{aligned} \mathbf{D} &= \text{diag}(\sigma_1^2, \dots, \sigma_n^2) \\ &= \mathbb{E}(\mathbf{e}\mathbf{e}' | \mathbf{X}) \\ &= \mathbb{E}(\mathbf{D}_0 | \mathbf{X}) \end{aligned}$$

where  $\mathbf{D}_0 = \text{diag}(e_1^2, \dots, e_n^2)$ . Thus  $\mathbf{D}_0$  is a conditionally unbiased estimator for  $\mathbf{D}$ . If the squared errors  $e_i^2$  were observable, we could construct the unbiased estimator

$$\begin{aligned} \hat{\mathbf{V}}_{\hat{\beta}}^{ideal} &= (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{D}_0\mathbf{X}) (\mathbf{X}'\mathbf{X})^{-1} \\ &= (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' e_i^2 \right) (\mathbf{X}'\mathbf{X})^{-1}. \end{aligned}$$

Indeed,

$$\begin{aligned} \mathbb{E}(\hat{\mathbf{V}}_{\hat{\beta}}^{ideal} | \mathbf{X}) &= (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \mathbb{E}(e_i^2 | \mathbf{X}) \right) (\mathbf{X}'\mathbf{X})^{-1} \\ &= (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \sigma_i^2 \right) (\mathbf{X}'\mathbf{X})^{-1} \\ &= (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{D}\mathbf{X}) (\mathbf{X}'\mathbf{X})^{-1} \\ &= \mathbf{V}_{\hat{\beta}} \end{aligned}$$

verifying that  $\hat{\mathbf{V}}_{\hat{\beta}}^{ideal}$  is unbiased for  $\mathbf{V}_{\hat{\beta}}$ .



Since the errors  $e_i^2$  are unobserved,  $\widehat{\mathbf{V}}_{\beta}^{ideal}$  is not a feasible estimator. However, we can replace the errors  $e_i$  with the least-squares residuals  $\widehat{e}_i$ . Making this substitution we obtain the estimator

$$\widehat{\mathbf{V}}_{\beta}^W = (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \widehat{e}_i^2 \right) (\mathbf{X}'\mathbf{X})^{-1}. \quad (4.37)$$

We know, however, that  $\widehat{e}_i^2$  is biased towards zero (recall equation (4.24)). To estimate the variance  $\sigma^2$  the unbiased estimator  $s^2$  scales the moment estimator  $\widehat{\sigma}^2$  by  $n/(n-k)$ . Making the same adjustment we obtain the estimator

$$\widehat{\mathbf{V}}_{\beta} = \left( \frac{n}{n-k} \right) (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \widehat{e}_i^2 \right) (\mathbf{X}'\mathbf{X})^{-1}. \quad (4.38)$$

While the scaling by  $n/(n-k)$  is *ad hoc*, it is recommended over the unscaled estimator (4.37).

Alternatively, we could use the prediction errors  $\widetilde{e}_i$  or the standardized residuals  $\bar{e}_i$ , yielding the estimators

$$\begin{aligned} \widetilde{\mathbf{V}}_{\beta} &= (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \widetilde{e}_i^2 \right) (\mathbf{X}'\mathbf{X})^{-1} \\ &= (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n (1 - h_{ii})^{-2} \mathbf{x}_i \mathbf{x}_i' \widehat{e}_i^2 \right) (\mathbf{X}'\mathbf{X})^{-1} \end{aligned} \quad (4.39)$$

and

$$\begin{aligned} \overline{\mathbf{V}}_{\beta} &= (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \bar{e}_i^2 \right) (\mathbf{X}'\mathbf{X})^{-1} \\ &= (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n (1 - h_{ii})^{-1} \mathbf{x}_i \mathbf{x}_i' \widehat{e}_i^2 \right) (\mathbf{X}'\mathbf{X})^{-1}. \end{aligned} \quad (4.40)$$

The four estimators  $\widehat{\mathbf{V}}_{\beta}^W$ ,  $\widehat{\mathbf{V}}_{\beta}$ ,  $\widetilde{\mathbf{V}}_{\beta}$ , and  $\overline{\mathbf{V}}_{\beta}$  are collectively called **robust, heteroskedasticity-consistent**, or **heteroskedasticity-robust** covariance matrix estimators. The estimator  $\widehat{\mathbf{V}}_{\beta}^W$  was first developed by Eicker (1963) and introduced to econometrics by White (1980), and is sometimes called the **Eicker-White** or **White** covariance matrix estimator. The degree-of-freedom adjustment in  $\widehat{\mathbf{V}}_{\beta}$  was recommended by Hinkley (1977), and is the default robust covariance matrix estimator implemented in Stata. (It is implemented by the “,r” option, for example by a regression executed with the command “reg y x, r”. In current applied econometric practice, this is the method used by most users.) The estimator  $\overline{\mathbf{V}}_{\beta}$  was introduced by Horn, Horn and Duncan (1975) (and is implemented using the vce(hc2) option in Stata). The estimator  $\widetilde{\mathbf{V}}_{\beta}$  was derived by MacKinnon and White from the jackknife principle, and by Andrews (1991) based on the principle of leave-one-out cross-validation (and is implemented using the vce(hc3) option in Stata).

Since  $(1 - h_{ii})^{-2} > (1 - h_{ii})^{-1} > 1$  it is straightforward to show that

$$\widehat{\mathbf{V}}_{\beta}^W < \overline{\mathbf{V}}_{\beta} < \widetilde{\mathbf{V}}_{\beta} \quad (4.41)$$

(See Exercise 4.10). The inequality  $\mathbf{A} < \mathbf{B}$  when applied to matrices means that the matrix  $\mathbf{B} - \mathbf{A}$  is positive definite.

In general, the bias of the covariance matrix estimators is quite complicated, but they greatly simplify under the assumption of homoskedasticity (4.3). For example, using (4.24),

$$\begin{aligned}
\mathbb{E}(\widehat{\mathbf{V}}_{\widehat{\beta}}^W | \mathbf{X}) &= (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \mathbb{E}(\widehat{e}_i^2 | \mathbf{X}) \right) (\mathbf{X}'\mathbf{X})^{-1} \\
&= (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' (1 - h_{ii}) \sigma^2 \right) (\mathbf{X}'\mathbf{X})^{-1} \\
&= (\mathbf{X}'\mathbf{X})^{-1} \sigma^2 - (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' h_{ii} \right) (\mathbf{X}'\mathbf{X})^{-1} \sigma^2 \\
&< (\mathbf{X}'\mathbf{X})^{-1} \sigma^2 \\
&= \mathbf{V}_{\widehat{\beta}}.
\end{aligned}$$

This calculation shows that  $\widehat{\mathbf{V}}_{\widehat{\beta}}^W$  is biased towards zero.

By a similarly calculation (again under homoskedasticity) we can calculate that the estimator  $\overline{\mathbf{V}}_{\widehat{\beta}}$  is unbiased

$$\mathbb{E}(\overline{\mathbf{V}}_{\widehat{\beta}} | \mathbf{X}) = (\mathbf{X}'\mathbf{X})^{-1} \sigma^2. \quad (4.42)$$

(See Exercise 4.11.)

It might seem rather odd to compare the bias of heteroskedasticity-robust estimators under the assumption of homoskedasticity, but it does give us a baseline for comparison.

Another interesting calculation shows that in general (that is, without assuming homoskedasticity)  $\widehat{\mathbf{V}}_{\widehat{\beta}}$  is biased away from zero. Indeed, using the definition of the prediction errors (3.44)

$$\widetilde{e}_i = y_i - \mathbf{x}_i' \widehat{\beta}_{(-i)} = e_i - \mathbf{x}_i' (\widehat{\beta}_{(-i)} - \beta)$$

so

$$\widetilde{e}_i^2 = e_i^2 - 2\mathbf{x}_i' (\widehat{\beta}_{(-i)} - \beta) e_i + \left( \mathbf{x}_i' (\widehat{\beta}_{(-i)} - \beta) \right)^2.$$

Note that  $e_i$  and  $\widehat{\beta}_{(-i)}$  are functions of non-overlapping observations and are thus independent. Hence  $\mathbb{E}(\left( \widehat{\beta}_{(-i)} - \beta \right) e_i | \mathbf{X}) = 0$  and

$$\begin{aligned}
\mathbb{E}(\widetilde{e}_i^2 | \mathbf{X}) &= \mathbb{E}(e_i^2 | \mathbf{X}) - 2\mathbf{x}_i' \mathbb{E}(\left( \widehat{\beta}_{(-i)} - \beta \right) e_i | \mathbf{X}) + \mathbb{E}\left( \left( \mathbf{x}_i' (\widehat{\beta}_{(-i)} - \beta) \right)^2 | \mathbf{X} \right) \\
&= \sigma_i^2 + \mathbb{E}\left( \left( \mathbf{x}_i' (\widehat{\beta}_{(-i)} - \beta) \right)^2 | \mathbf{X} \right) \\
&\geq \sigma_i^2.
\end{aligned}$$

It follows that

$$\begin{aligned}
\mathbb{E}(\widetilde{\mathbf{V}}_{\widehat{\beta}} | \mathbf{X}) &= (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \mathbb{E}(\widetilde{e}_i^2 | \mathbf{X}) \right) (\mathbf{X}'\mathbf{X})^{-1} \\
&\geq (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \sigma_i^2 \right) (\mathbf{X}'\mathbf{X})^{-1} \\
&= \mathbf{V}_{\widehat{\beta}}.
\end{aligned}$$

This means that  $\widetilde{\mathbf{V}}_{\widehat{\beta}}$  is conservative in the sense that it is weakly larger (in expectation) than the correct variance for any realization of  $\mathbf{X}$ .

We have introduced five covariance matrix estimators,  $\widehat{\mathbf{V}}_{\hat{\beta}}^0$ ,  $\widehat{\mathbf{V}}_{\hat{\beta}}^W$ ,  $\widehat{\mathbf{V}}_{\hat{\beta}}$ ,  $\widetilde{\mathbf{V}}_{\hat{\beta}}$ , and  $\overline{\mathbf{V}}_{\hat{\beta}}$ . Which should you use? The classic estimator  $\widehat{\mathbf{V}}_{\hat{\beta}}^0$  is typically a poor choice, as it is only valid under the unlikely homoskedasticity restriction. For this reason it is not typically used in contemporary econometric research. Unfortunately, standard regression packages set their default choice as  $\widehat{\mathbf{V}}_{\hat{\beta}}^0$ , so users must intentionally select a robust covariance matrix estimator.

Of the four robust estimators,  $\widehat{\mathbf{V}}_{\hat{\beta}}$  is the most commonly used as it is the default robust covariance matrix option in Stata. However,  $\widetilde{\mathbf{V}}_{\hat{\beta}}$  may be the preferred choice since it is conservative for any  $\mathbf{X}$ . As  $\widetilde{\mathbf{V}}_{\hat{\beta}}$  is simple to implement, this should not be a barrier.

### Halbert L. White

Hal White (1950-2012) of the United States was an influential econometrician of recent years. His 1980 paper on heteroskedasticity-consistent covariance matrix estimation for many years has been the most cited paper in economics. His research was central to the movement to view econometric models as approximations, and to the drive for increased mathematical rigor in the discipline. In addition to being a highly prolific and influential scholar, he also co-founded the economic consulting firm Bates White.

## 4.14 Standard Errors

A variance estimator such as  $\widehat{\mathbf{V}}_{\hat{\beta}}$  is an estimate of the variance of the distribution of  $\hat{\beta}$ . A more easily interpretable measure of spread is its square root – the standard deviation. This is so important when discussing the distribution of parameter estimates, we have a special name for estimates of their standard deviation.

**Definition 4.14.1** A *standard error*  $s(\hat{\beta})$  for a real-valued estimator  $\hat{\beta}$  is an estimate of the standard deviation of the distribution of  $\hat{\beta}$ .

When  $\beta$  is a vector with estimate  $\hat{\beta}$  and covariance matrix estimate  $\widehat{\mathbf{V}}_{\hat{\beta}}$ , standard errors for individual elements are the square roots of the diagonal elements of  $\widehat{\mathbf{V}}_{\hat{\beta}}$ . That is,

$$s(\hat{\beta}_j) = \sqrt{\widehat{\mathbf{V}}_{\hat{\beta}_j}} = \sqrt{[\widehat{\mathbf{V}}_{\hat{\beta}}]_{jj}}.$$

When the classical covariance matrix estimate (4.36) is used, the standard error takes the particularly simple form

$$s(\hat{\beta}_j) = s \sqrt{[(\mathbf{X}'\mathbf{X})^{-1}]_{jj}}. \quad (4.43)$$

As we discussed in the previous section, there are multiple possible covariance matrix estimators, so standard errors are not unique. It is therefore important to understand what formula and method is used by an author when studying their work. It is also important to understand that a particular standard error may be relevant under one set of model assumptions, but not under another set of assumptions.

To illustrate, we return to the log wage regression (3.14) of Section 3.7. We calculate that  $s^2 = 0.160$ . Therefore the homoskedastic covariance matrix estimate is

$$\hat{\mathbf{V}}_{\hat{\beta}}^0 = \begin{pmatrix} 5010 & 314 \\ 314 & 20 \end{pmatrix}^{-1} 0.160 = \begin{pmatrix} 0.002 & -0.031 \\ -0.031 & 0.499 \end{pmatrix}.$$

We also calculate that

$$\sum_{i=1}^n (1 - h_{ii})^{-1} \mathbf{x}_i \mathbf{x}_i' \hat{e}_i^2 = \begin{pmatrix} 763.26 & 48.513 \\ 48.513 & 3.1078 \end{pmatrix}.$$

Therefore the Horn-Horn-Duncan covariance matrix estimate is

$$\begin{aligned} \bar{\mathbf{V}}_{\hat{\beta}} &= \begin{pmatrix} 5010 & 314 \\ 314 & 20 \end{pmatrix}^{-1} \begin{pmatrix} 763.26 & 48.513 \\ 48.513 & 3.1078 \end{pmatrix} \begin{pmatrix} 5010 & 314 \\ 314 & 20 \end{pmatrix}^{-1} \\ &= \begin{pmatrix} 0.001 & -0.015 \\ -0.015 & 0.243 \end{pmatrix}. \end{aligned} \quad (4.44)$$

The standard errors are the square roots of the diagonal elements of these matrices. A conventional format to write the estimated equation with standard errors is

$$\widehat{\log(\text{Wage})} = \begin{matrix} 0.155 & \text{Education} + & 0.698 \\ (0.031) & & (0.493) \end{matrix}.$$

Alternatively, standard errors could be calculated using the other formulae. We report the different standard errors in the following table.

	Education	Intercept
Homoskedastic (4.36)	0.045	0.707
White (4.37)	0.029	0.461
Scaled White (4.38)	0.030	0.486
Andrews (4.39)	0.033	0.527
Horn-Horn-Duncan (4.40)	0.031	0.493

The homoskedastic standard errors are noticeably different (larger, in this case) than the others, but the four robust standard errors are quite close to one another.

## 4.15 Covariance Matrix Estimation with Sparse Dummy Variables

The heteroskedasticity-robust covariance matrix estimators can be quite imprecise in some contexts. One is in the presence of **sparse dummy variables** – when a dummy variable only takes the value 1 or 0 for very few observations. In these contexts one component of the variance matrix is estimated on just those few observations and thus will be imprecise. This is effectively hidden from the user.

To see the problem, let  $d_{1i}$  be a dummy variable (takes on the values 1 and 0) for “group 1” and let  $d_{2i} = 1 - d_{1i}$  be the complement for “group 2”. Consider the dummy-only regression

$$y_i = \beta_1 d_{1i} + \beta_2 d_{2i} + e_i$$

which excludes the intercept for identification. The number of observations in the two “groups” are  $n_1 = \sum_{i=1}^n d_{1i}$  and  $n_2 = \sum_{i=1}^n d_{2i}$ . The least-squares estimates for  $\beta_1$  and  $\beta_2$  are the averages

within the two groups. We say the design is sparse if either  $n_1$  or  $n_2$  is small. One implication is that the coefficient for the small group will be imprecisely estimated.

An extreme situation is when  $n_1 = 1$ , thus group 1 has only a single observation. This would be unlikely to occur intentionally, but is actually remarkably likely when a large number of interactions are included in a regression. In this context, the least-squares estimate for  $\beta_1$  is  $\hat{\beta}_1 = y_1$ , where for simplicity we have assumed that the first observation is the one for which  $d_{1i} = 1$ . This means that the corresponding residual is  $\hat{e}_1 = 0$ .

The implication for covariance matrix estimation is rather unpleasant. The White estimator is

$$\hat{\mathbf{V}}_{\hat{\beta}}^W = \begin{pmatrix} 0 & 0 \\ 0 & \frac{\hat{\sigma}^2}{n-1} \end{pmatrix}$$

where  $\hat{\sigma}^2$  is a variance estimator computed with all observations excluding the first. The covariance matrix  $\hat{\mathbf{V}}_{\hat{\beta}}^W$  is singular, and in particular produces the standard error  $s(\hat{\beta}_1) = 0$ ! That is, the standard regression package will print out a standard error of 0 for the least-precisely estimated coefficient!

The reason is that the estimator is effectively estimating the variance of  $\hat{\beta}_1$  from a single observation. The point estimate of a variance from a single observation is 0. Essentially, while it is impossible to estimate a variance from a single observation the standard formula gives a misleadingly precise answer.

In most practical regressions, estimated standard errors will not be zero as we typically estimate models with an omitted dummy category and an intercept. What are the implications? In this case, while the reported “standard errors” are non-zero, the covariance matrix estimator itself is singular. This means that there is a linear combination of the estimates with a zero estimated variance. This is generally troubling as this situation is largely hidden from the user.

This problem does not arise if the homoskedastic form of the covariance matrix estimate is used. In the above example, the estimate is

$$\hat{\mathbf{V}}_{\hat{\beta}}^0 = \begin{pmatrix} s^2 & 0 \\ 0 & \frac{s^2}{n-1} \end{pmatrix}.$$

Consequently, in models with sparse dummy variable designs, it may be prudent to use (or at least check) the homoskedastic standard error formulae.

In general, users should be cautious about regression results when dummy variables (and interactions of dummy variables) are sparse.

## 4.16 Computation

We illustrate methods to compute standard errors for equation (3.15) extending the code of Section 3.20.

### Stata do File (continued)

```
*      Homoskedastic formula (4.36):
reg wage education experience exp2 if (mnwf == 1)
*      Scaled White formula (4.38):
reg wage education experience exp2 if (mnwf == 1), r
*      Andrews formula (4.39):
reg wage education experience exp2 if (mnwf == 1), vce(hc3)
*      Horn-Horn-Duncan formula (4.40):
reg wage education experience exp2 if (mnwf == 1), vce(hc2)
```

**R Program File (continued)**

```

n <- nrow(y)
k <- ncol(x)
a <- n/(n-k)
sig2 <- (t(e) %*% e)/(n-k)
u1 <- x*(e%*%matrix(1,1,k))
u2 <- x*((e/(1-leverage))%*%matrix(1,1,k))
u3 <- x*((e/sqrt(1-leverage))%*%matrix(1,1,k))
v0 <- xx*sig2
xx <- solve(t(x)%*%x)
v1 <- xx %*% (t(u1)%*%u1) %*% xx
v1a <- a * xx %*% (t(u1)%*%u1) %*% xx
v2 <- xx %*% (t(u2)%*%u2) %*% xx
v3 <- xx %*% (t(u3)%*%u3) %*% xx
s0 <- sqrt(diag(v0))          # Homoskedastic formula
s1 <- sqrt(diag(v1))          # White formula
s1a <- sqrt(diag(v1a))        # Scaled White formula
s2 <- sqrt(diag(v2))          # Andrews formula
s3 <- sqrt(diag(v3))          # Horn-Horn-Duncan formula

```

**MATLAB Program File (continued)**

```

[n,k]=size(x);
a=n/(n-k);
sig2=(e'*e)/(n-k);
u1=x.*(e*ones(1,k));
u2=x.*((e./(1-leverage))*ones(1,k));u3=x.*((e./sqrt(1-
leverage))*ones(1,k));
xx=inv(x'*x);
v0=xx*sig2;
v1=xx*(u1'*u1)*xx;
v1a=a*xx*(u1'*u1)*xx;
v2=xx*(u2'*u2)*xx;
v3=xx*(u3'*u3)*xx;
s0=sqrt(diag(v0));           # Homoskedastic formula
s1=sqrt(diag(v1));           # White formula
s1a=sqrt(diag(v1a));         # Scaled White formula
s2=sqrt(diag(v2));           # Andrews formula
s3=sqrt(diag(v3));           # Horn-Horn-Duncan formula

```

**4.17 Measures of Fit**

As we described in the previous chapter, a commonly reported measure of regression fit is the regression  $R^2$  defined as

$$R^2 = 1 - \frac{\sum_{i=1}^n \hat{e}_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\hat{\sigma}^2}{\hat{\sigma}_y^2}.$$

where  $\hat{\sigma}_y^2 = n^{-1} \sum_{i=1}^n (y_i - \bar{y})^2$ .  $R^2$  can be viewed as an estimator of the population parameter

$$\rho^2 = \frac{\text{var}(\mathbf{x}'_i \boldsymbol{\beta})}{\text{var}(y_i)} = 1 - \frac{\sigma^2}{\sigma_y^2}.$$

However,  $\hat{\sigma}^2$  and  $\hat{\sigma}_y^2$  are biased estimators. Theil (1961) proposed replacing these by the unbiased versions  $s^2$  and  $\tilde{\sigma}_y^2 = (n-1)^{-1} \sum_{i=1}^n (y_i - \bar{y})^2$  yielding what is known as **R-bar-squared** or **adjusted R-squared**:

$$\bar{R}^2 = 1 - \frac{s^2}{\tilde{\sigma}_y^2} = 1 - \frac{(n-1) \sum_{i=1}^n \tilde{e}_i^2}{(n-k) \sum_{i=1}^n (y_i - \bar{y})^2}.$$

While  $\bar{R}^2$  is an improvement on  $R^2$ , a much better improvement is

$$\tilde{R}^2 = 1 - \frac{\sum_{i=1}^n \tilde{e}_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\tilde{\sigma}^2}{\tilde{\sigma}_y^2}$$

where  $\tilde{e}_i$  are the prediction errors (3.44) and  $\tilde{\sigma}^2$  is the MSPE from (3.47). As described in Section (4.11),  $\tilde{\sigma}^2$  is a good estimator of the out-of-sample mean-squared forecast error, so  $\tilde{R}^2$  is a good estimator of the percentage of the forecast variance which is explained by the regression forecast. In this sense,  $\tilde{R}^2$  is a good measure of fit.

One problem with  $R^2$ , which is partially corrected by  $\bar{R}^2$  and fully corrected by  $\tilde{R}^2$ , is that  $R^2$  necessarily increases when regressors are added to a regression model. This occurs because  $R^2$  is a negative function of the sum of squared residuals which cannot increase when a regressor is added. In contrast,  $\bar{R}^2$  and  $\tilde{R}^2$  are non-monotonic in the number of regressors.  $\tilde{R}^2$  can even be negative, which occurs when an estimated model predicts worse than a constant-only model.

In the statistical literature the MSPE  $\tilde{\sigma}^2$  is known as the **leave-one-out cross validation** criterion, and is popular for model comparison and selection, especially in high-dimensional (non-parametric) contexts. It is equivalent to use  $\tilde{R}^2$  or  $\tilde{\sigma}^2$  to compare and select models. Models with high  $\tilde{R}^2$  (or low  $\tilde{\sigma}^2$ ) are better models in terms of expected out of sample squared error. In contrast,  $R^2$  cannot be used for model selection, as it necessarily increases when regressors are added to a regression model.  $\bar{R}^2$  is also an inappropriate choice for model selection (it tends to select models with too many parameters), though a justification of this assertion requires a study of the theory of model selection. Unfortunately,  $\bar{R}^2$  is routinely used by some economists, possibly as a hold-over from previous generations.

In summary, it is recommended to calculate and report  $\tilde{R}^2$  and/or  $\tilde{\sigma}^2$  in regression analysis, and omit  $R^2$  and  $\bar{R}^2$ .

### Henri Theil

Henri Theil (1924-2000) of the Netherlands invented  $\bar{R}^2$  and two-stage least squares, both of which are routinely seen in applied econometrics. He also wrote an early influential advanced textbook on econometrics (Theil, 1971).

## 4.18 Empirical Example

We again return to our wage equation, but use a much larger sample of all individuals with at least 12 years of education. For regressors we include years of education, potential work experience, experience squared, and dummy variable indicators for the following: female, female union member,

male union member, married female<sup>1</sup>, married male, formerly married female<sup>2</sup>, formerly married male, Hispanic, black, American Indian, Asian, and mixed race<sup>3</sup>. The available sample is 46,943 so the parameter estimates are quite precise and reported in Table 4.1. For standard errors we use the unbiased Horn-Horn-Duncan formula.

Table 4.1 displays the parameter estimates in a standard tabular format. The table clearly states the estimation method (OLS), the dependent variable ( $\log(\text{Wage})$ ), and the regressors are clearly labeled. Both parameter estimates and standard errors are reported for all coefficients. In addition to the coefficient estimates, the table also reports the estimated error standard deviation and the sample size. These are useful summary measures of fit which aid readers.

Table 4.1  
OLS Estimates of Linear Equation for  $\log(\text{Wage})$

	$\hat{\beta}$	$s(\hat{\beta})$
Education	0.117	0.001
Experience	0.033	0.001
Experience <sup>2</sup> /100	-0.056	0.002
Female	-0.098	0.011
Female Union Member	0.023	0.020
Male Union Member	0.095	0.020
Married Female	0.016	0.010
Married Male	0.211	0.010
Formerly Married Female	-0.006	0.012
Formerly Married Male	0.083	0.015
Hispanic	-0.108	0.008
Black	-0.096	0.008
American Indian	-0.137	0.027
Asian	-0.038	0.013
Mixed Race	-0.041	0.021
Intercept	0.909	0.021
$\hat{\sigma}$	0.565	
Sample Size	46,943	

Note: Standard errors are heteroskedasticity-consistent (Horn-Horn-Duncan formula)

As a general rule, it is advisable to always report standard errors along with parameter estimates. This allows readers to assess the precision of the parameter estimates, and as we will discuss in later chapters, form confidence intervals and t-tests for individual coefficients if desired.

The results in Table 4.1 confirm our earlier findings that the return to a year of education is approximately 12%, the return to experience is concave, that single women earn approximately 10% less than single men, and blacks earn about 10% less than whites. In addition, we see that Hispanics earn about 11% less than whites, American Indians 14% less, and Asians and Mixed races about 4% less. We also see there are wage premiums for men who are members of a labor union (about 10%), married (about 21%) or formerly married (about 8%), but no similar premiums are apparent for women.

<sup>1</sup>Defining “married” as marital code 1, 2, or 3.

<sup>2</sup>Defining “formerly married” as marital code 4, 5, or 6.

<sup>3</sup>Race code 6 or higher.



## 4.19 Multicollinearity

If  $\mathbf{X}'\mathbf{X}$  is singular, then  $(\mathbf{X}'\mathbf{X})^{-1}$  and  $\hat{\boldsymbol{\beta}}$  are not defined. This situation is called **strict multicollinearity**, as the columns of  $\mathbf{X}$  are linearly dependent, i.e., there is some  $\boldsymbol{\alpha} \neq \mathbf{0}$  such that  $\mathbf{X}\boldsymbol{\alpha} = \mathbf{0}$ . Most commonly, this arises when sets of regressors are included which are identically related. For example, if  $\mathbf{X}$  includes both the logs of two prices and the log of the relative prices,  $\log(p_1)$ ,  $\log(p_2)$  and  $\log(p_1/p_2)$ , then  $\mathbf{X}'\mathbf{X}$  will necessarily be singular. When this happens, the applied researcher quickly discovers the error as the statistical software will be unable to construct  $(\mathbf{X}'\mathbf{X})^{-1}$ . Since the error is discovered quickly, this is rarely a problem for applied econometric practice.

The more relevant situation is **near multicollinearity**, which is often called “multicollinearity” for brevity. This is the situation when the  $\mathbf{X}'\mathbf{X}$  matrix is near singular, when the columns of  $\mathbf{X}$  are close to linearly dependent. This definition is not precise, because we have not said what it means for a matrix to be “near singular”. This is one difficulty with the definition and interpretation of multicollinearity.

One potential complication of near singularity of matrices is that the numerical reliability of the calculations may be reduced. In practice this is rarely an important concern, except when the number of regressors is very large.

A more relevant implication of near multicollinearity is that individual coefficient estimates will be imprecise. We can see this most simply in a homoskedastic linear regression model with two regressors

$$y_i = x_{1i}\beta_1 + x_{2i}\beta_2 + e_i,$$

and

$$\frac{1}{n}\mathbf{X}'\mathbf{X} = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}.$$

In this case

$$\text{var}(\hat{\boldsymbol{\beta}} | \mathbf{X}) = \frac{\sigma^2}{n} \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}^{-1} = \frac{\sigma^2}{n(1-\rho^2)} \begin{pmatrix} 1 & -\rho \\ -\rho & 1 \end{pmatrix}.$$

The correlation  $\rho$  indexes collinearity, since as  $\rho$  approaches 1 the matrix becomes singular. We can see the effect of collinearity on precision by observing that the variance of a coefficient estimate  $\sigma^2 [n(1-\rho^2)]^{-1}$  approaches infinity as  $\rho$  approaches 1. Thus the more “collinear” are the regressors, the worse the precision of the individual coefficient estimates.

What is happening is that when the regressors are highly dependent, it is statistically difficult to disentangle the impact of  $\beta_1$  from that of  $\beta_2$ . As a consequence, the precision of individual estimates are reduced. The imprecision, however, will be reflected by large standard errors, so there is no distortion in inference.

Some earlier textbooks overemphasized a concern about multicollinearity. A very amusing parody of these texts appeared in Chapter 23.3 of Goldberger’s *A Course in Econometrics* (1991), which is reprinted below. To understand his basic point, you should notice how the estimation variance  $\sigma^2 [n(1-\rho^2)]^{-1}$  depends equally and symmetrically on the correlation  $\rho$  and the sample size  $n$ .

**Arthur S. Goldberger**

Art Goldberger (1930-2009) was one of the most distinguished members of the Department of Economics at the University of Wisconsin. His PhD thesis developed an early macroeconometric forecasting model (known as the Klein-Goldberger model) but most of his career focused on microeconomic issues. He was the leading pioneer of what has been called the Wisconsin Tradition of empirical work – a combination of formal econometric theory with a careful critical analysis of empirical work. Goldberger wrote a series of highly regarded and influential graduate econometric textbooks, including *Econometric Theory* (1964), *Topics in Regression Analysis* (1968), and *A Course in Econometrics* (1991).

**Micronumerosity**

Arthur S. Goldberger

*A Course in Econometrics* (1991), Chapter 23.3

Econometrics texts devote many pages to the problem of multicollinearity in multiple regression, but they say little about the closely analogous problem of small sample size in estimating a univariate mean. Perhaps that imbalance is attributable to the lack of an exotic polysyllabic name for “small sample size.” If so, we can remove that impediment by introducing the term *micronumerosity*.

Suppose an econometrician set out to write a chapter about small sample size in sampling from a univariate population. Judging from what is now written about multicollinearity, the chapter might look like this:

1. *Micronumerosity*

The extreme case, “exact micronumerosity,” arises when  $n = 0$ , in which case the sample estimate of  $\mu$  is not unique. (Technically, there is a violation of the rank condition  $n > 0$  : the matrix  $0$  is singular.) The extreme case is easy enough to recognize. “Near micronumerosity” is more subtle, and yet very serious. It arises when the rank condition  $n > 0$  is barely satisfied. Near micronumerosity is very prevalent in empirical economics.

2. *Consequences of micronumerosity*

The consequences of micronumerosity are serious. Precision of estimation is reduced. There are two aspects of this reduction: estimates of  $\mu$  may have large errors, and not only that, but  $V_{\bar{y}}$  will be large.

Investigators will sometimes be led to accept the hypothesis  $\mu = 0$  because  $\bar{y}/\hat{\sigma}_{\bar{y}}$  is small, even though the true situation may be not that  $\mu = 0$  but simply that the sample data have not enabled us to pick  $\mu$  up.

The estimate of  $\mu$  will be very sensitive to sample data, and the addition of a few more observations can sometimes produce drastic shifts in the sample mean.

The true  $\mu$  may be sufficiently large for the null hypothesis  $\mu = 0$  to be rejected, even though  $V_{\bar{y}} = \sigma^2/n$  is large because of micronumerosity. But if the true  $\mu$  is small (although nonzero) the hypothesis  $\mu = 0$  may mistakenly be accepted.

### 3. *Testing for micronumerosity*

Tests for the presence of micronumerosity require the judicious use of various fingers. Some researchers prefer a single finger, others use their toes, still others let their thumbs rule.

A generally reliable guide may be obtained by counting the number of observations. Most of the time in econometric analysis, when  $n$  is close to zero, it is also far from infinity.

Several test procedures develop critical values  $n^*$ , such that micronumerosity is a problem only if  $n$  is smaller than  $n^*$ . But those procedures are questionable.

### 4. *Remedies for micronumerosity*

If micronumerosity proves serious in the sense that the estimate of  $\mu$  has an unsatisfactorily low degree of precision, we are in the statistical position of not being able to make bricks without straw. The remedy lies essentially in the acquisition, if possible, of larger samples from the same population.

But more data are no remedy for micronumerosity if the additional data are simply “more of the same.” So obtaining lots of small samples from the same population will not help.

## 4.20 Clustered Sampling

In Section 4.2 we briefly mentioned clustered sampling as an alternative to the assumption of random sampling. We now introduce the framework in more detail and extend the primary results of this Chapter to encompass clustered dependence.

It might be easiest to understand the idea of clusters by considering a concrete example. Duflo, Dupas and Kremer (2011) investigate the impact of tracking (assigning students based on initial test score) on educational attainment in a randomized experiment. An extract of their data set is available on the textbook webpage in the file DDK2011.

In 2005, 140 primary schools in Kenya received funding to hire an extra first grade teacher to reduce class sizes. In half of the schools (selected randomly), students were assigned to classrooms based on an initial test score (“tracking”); in the remaining schools the students were randomly assigned to classrooms. For their analysis, the authors restricted attention to the 121 schools which initially had a single first-grade class, and if we further restrict attention to those with full data availability the resulting sample has 111 schools.

The key regression in the paper takes the form

$$TestScore_{ig} = -0.082 + 0.147Tracking_g + e_{ig} \quad (4.45)$$

where  $TestScore_{ig}$  is the standardized test score (normalized to have mean 0 and variance 1) of student  $i$  in school  $g$ , and  $Tracking_g$  is a dummy equal to 1 if school  $g$  was tracking. The OLS estimates indicate that schools which tracked the students had an overall increase in test scores by 0.15 standard deviations, which is quite meaningful. More general versions of this regression are estimated, many of which take the form

$$TestScore_{ig} = \alpha + \gamma Tracking_g + \mathbf{x}'_{ig} \boldsymbol{\beta} + e_{ig} \quad (4.46)$$

where  $\mathbf{x}_{ig}$  is a set of controls specific to the student (including age, sex and initial test score).

A difficulty with applying the classical regression framework is that student achievement is likely to be dependent within a given school. Student achievement may be affected by local demographics, individual teachers, and classmates, all of which imply dependence within a school. These concerns, however, do not suggest that achievement will be correlated across schools, so it seems reasonable to model achievement across schools as mutually independent.

In clustering contexts it is convenient to double index the observations as  $(y_{ig}, \mathbf{x}_{ig})$  where  $g = 1, \dots, G$  indexes the cluster and  $i = 1, \dots, n_g$  indexes the individual within the  $g^{\text{th}}$  cluster. The number of observations per cluster  $n_g$  may vary across clusters. The number of clusters is  $G$ . The total number of observations is  $n = \sum_{g=1}^G n_g$ . In the Kenyan schooling example, the number of clusters (schools) in the estimation sample is  $G = 111$ , the number of students per school varies from 19 to 62, and the total number of observations is  $n = 5269$ .

While it is typical to write the observations using the double index notation  $(y_{ig}, \mathbf{x}_{ig})$ , it is also useful to use cluster-level notation. Let  $\mathbf{y}_g = (y_{1g}, \dots, y_{n_g g})'$  and  $\mathbf{X}_g = (\mathbf{x}_{1g}, \dots, \mathbf{x}_{n_g g})'$  denote the  $n_g \times 1$  vector of dependent variables and  $n_g \times k$  matrix of regressors for the  $g^{\text{th}}$  cluster. A linear regression model can be written for the individual observations as

$$y_{ig} = \mathbf{x}_{ig}'\boldsymbol{\beta} + e_{ig}$$

and using cluster notation as

$$\mathbf{y}_g = \mathbf{X}_g\boldsymbol{\beta} + \mathbf{e}_g \quad (4.47)$$

where  $\mathbf{e}_g = (e_{1g}, \dots, e_{n_g g})'$  is a  $n_g \times 1$  error vector.

Using this notation we can write the sums over the observations using the double sum  $\sum_{g=1}^G \sum_{i=1}^{n_g}$ . This is the sum across clusters of the sum across observations within each cluster. The OLS estimator can be written as

$$\hat{\boldsymbol{\beta}} = \left( \sum_{g=1}^G \sum_{i=1}^{n_g} \mathbf{x}_{ig} \mathbf{x}_{ig}' \right)^{-1} \left( \sum_{g=1}^G \sum_{i=1}^{n_g} \mathbf{x}_{ig} y_{ig} \right)$$

or

$$\hat{\boldsymbol{\beta}} = \left( \sum_{g=1}^G \mathbf{X}_g' \mathbf{X}_g \right)^{-1} \left( \sum_{g=1}^G \mathbf{X}_g' \mathbf{y}_g \right). \quad (4.48)$$

The OLS residuals are  $\hat{e}_{ig} = y_{ig} - \mathbf{x}_{ig}'\hat{\boldsymbol{\beta}}$  in individual level notation and  $\hat{\mathbf{e}}_g = \mathbf{y}_g - \mathbf{X}_g\hat{\boldsymbol{\beta}}$  in cluster level notation.

The standard clustering assumption is that the clusters are known to the researcher and that the observations are independent across clusters.

**Assumption 4.20.1** *The clusters  $(\mathbf{y}_g, \mathbf{X}_g)$  are mutually independent across clusters  $g$ .*

In our example, clusters are schools. In other common applications, cluster dependence has been assumed within individual classrooms, families, villages, regions, and within larger units such as industries and states. This choice is up to the researcher, though the justification will depend on the context, the nature of the data, and will reflect information and assumptions on the dependence structure across observations.

The model is a linear regression under the assumption

$$\mathbb{E}(\mathbf{e}_g \mid \mathbf{X}_g) = 0. \quad (4.49)$$

This is the same as assuming that the individual errors are conditionally mean zero

$$\mathbb{E}(e_{ig} \mid \mathbf{X}_g) = 0$$

or that the conditional mean of  $\mathbf{y}_g$  given  $\mathbf{X}_g$  is linear. As in the independent case, equation (4.49) means that the linear regression model is correctly specified. In the clustered regression model, this requires that all interaction effects within clusters have been accounted for in the specification of the individual regressors  $\mathbf{x}_{ig}$ .

In the regression (4.45), the conditional mean is necessarily linear and satisfies (4.49) since the regressor  $Tracking_g$  is a dummy variable at the cluster level. In the regression (4.46) with individual controls, (4.49) requires that the achievement of any student is unaffected by the individual controls (e.g. age, sex and initial test score) of other students within the same school.

Given (4.49), we can calculate the mean of the OLS estimator. Substituting (4.47) into (4.48) we find

$$\hat{\beta} - \beta = \left( \sum_{g=1}^G \mathbf{X}'_g \mathbf{X}_g \right)^{-1} \left( \sum_{g=1}^G \mathbf{X}'_g e_g \right).$$

The mean of  $\hat{\beta} - \beta$  conditioning on all the regressors is

$$\begin{aligned} \mathbb{E}(\hat{\beta} - \beta \mid \mathbf{X}) &= \left( \sum_{g=1}^G \mathbf{X}'_g \mathbf{X}_g \right)^{-1} \left( \sum_{g=1}^G \mathbf{X}'_g \mathbb{E}(e_g \mid \mathbf{X}) \right) \\ &= \left( \sum_{g=1}^G \mathbf{X}'_g \mathbf{X}_g \right)^{-1} \left( \sum_{g=1}^G \mathbf{X}'_g \mathbb{E}(e_g \mid \mathbf{X}_g) \right) \\ &= \mathbf{0}. \end{aligned}$$

The first equality holds by linearity, the second by Assumption 4.20.1 and the third by (4.49).

This shows that OLS is unbiased under clustering if the conditional mean is linear.

**Theorem 4.20.1** *In the clustered linear regression model (Assumption 4.20.1 and (4.49))*

$$\mathbb{E}(\hat{\beta} \mid \mathbf{X}) = \beta$$

Now consider the covariance matrix of  $\hat{\beta}$ . Let

$$\Sigma_g = \mathbb{E}(e_g e'_g \mid \mathbf{X}_g)$$

denote the  $n_g \times n_g$  conditional covariance matrix of the errors within the  $g^{th}$  cluster. Since the observations are independent across clusters,

$$\begin{aligned} \text{var} \left( \left( \sum_{g=1}^G \mathbf{X}'_g e_g \right) \mid \mathbf{X} \right) &= \sum_{g=1}^G \text{var}(\mathbf{X}'_g e_g \mid \mathbf{X}_g) \\ &= \sum_{g=1}^G \mathbf{X}'_g \mathbb{E}(e_g e'_g \mid \mathbf{X}_g) \mathbf{X}_g \\ &= \sum_{g=1}^G \mathbf{X}'_g \Sigma_g \mathbf{X}_g \\ &\stackrel{\text{def}}{=} \Omega_n \end{aligned} \tag{4.50}$$

It follows that

$$\begin{aligned} \mathbf{V}_{\hat{\beta}} &= \text{var}(\hat{\beta} \mid \mathbf{X}) \\ &= (\mathbf{X}'\mathbf{X})^{-1} \boldsymbol{\Omega}_n (\mathbf{X}'\mathbf{X})^{-1} \end{aligned} \quad (4.51)$$

where we write  $\mathbf{X}'\mathbf{X} = \sum_{g=1}^G \mathbf{X}'_g \mathbf{X}_g = \sum_{g=1}^G \sum_{i=1}^{n_g} \mathbf{x}_{ig} \mathbf{x}'_{ig}$ .

This differs from the formula in the independent case due to the correlation between observations within clusters. The magnitude of the difference depends on the degree of correlation between observations within clusters and the number of observations within clusters. To see this, suppose that all clusters have the same number of observations  $n_g = N$ ,  $\mathbb{E}(e_{ig}^2 \mid \mathbf{x}_g) = \sigma^2$ ,  $\mathbb{E}(e_{ig}e_{\ell g} \mid \mathbf{x}_g) = \sigma^2\rho$  for  $i \neq \ell$ , and the regressors  $\mathbf{x}_{ig}$  do not vary within a cluster. In this case the exact variance of the OLS estimator equals

$$\mathbf{V}_{\hat{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \sigma^2 (1 + \rho(N - 1)).$$

If  $\rho > 0$ , this shows that the actual variance is appropriately a multiple  $\rho N$  of the conventional formula. In the Kenyan school example, the average cluster size is 48, so if the correlation between students is  $\rho = 0.25$  the actual variance exceeds the conventional formula by a factor of about twelve. In this case the correct standard errors (the square root of the variance) should be a multiple of about three times the conventional formula. This is a substantial difference, and should not be neglected.

The typical solution is to use a covariance matrix estimate which extends the robust White formula to allow for general correlation within clusters. Recall that the insight of the White covariance estimator is that the squared error  $e_i^2$  is unbiased for  $\mathbb{E}(e_i^2 \mid \mathbf{x}_i) = \sigma_i^2$ . Similarly with cluster dependence the matrix  $\mathbf{e}_g \mathbf{e}'_g$  is unbiased for  $\mathbb{E}(\mathbf{e}_g \mathbf{e}'_g \mid \mathbf{X}_g) = \boldsymbol{\Sigma}_g$ . This means that an unbiased estimate for (4.50) is  $\hat{\boldsymbol{\Omega}}_n = \sum_{g=1}^G \mathbf{X}'_g \mathbf{e}_g \mathbf{e}'_g \mathbf{X}_g$ . This is not feasible, but we can replace the unknown errors by the OLS residuals to obtain the estimator

$$\begin{aligned} \hat{\boldsymbol{\Omega}}_n &= \sum_{g=1}^G \mathbf{X}'_g \hat{\mathbf{e}}_g \hat{\mathbf{e}}'_g \mathbf{X}_g \\ &= \sum_{g=1}^G \sum_{i=1}^{n_g} \sum_{\ell=1}^{n_g} \mathbf{x}_{ig} \mathbf{x}'_{\ell g} \hat{e}_{ig} \hat{e}_{\ell g} \\ &= \sum_{g=1}^G \left( \sum_{i=1}^{n_g} \mathbf{x}_{ig} \hat{e}_{ig} \right) \left( \sum_{\ell=1}^{n_g} \mathbf{x}_{\ell g} \hat{e}_{\ell g} \right)'. \end{aligned} \quad (4.52)$$

The three expressions in (4.50) give three equivalent formula which could be used to calculate  $\hat{\boldsymbol{\Omega}}_n$ . The final expression writes  $\hat{\boldsymbol{\Omega}}_n$  in terms of the cluster sums  $\sum_{\ell=1}^{n_g} \mathbf{x}_{\ell g} \hat{e}_{\ell g}$  which is basis for our example R and MATLAB codes shown below.

Given the expressions (4.50)-(4.51), a natural cluster covariance matrix estimator takes the form

$$\hat{\mathbf{V}}_{\hat{\beta}} = a_n (\mathbf{X}'\mathbf{X})^{-1} \hat{\boldsymbol{\Omega}}_n (\mathbf{X}'\mathbf{X})^{-1} \quad (4.53)$$

where the term  $a_n$  is a possible finite-sample adjustment. The Stata cluster command uses

$$a_n = \left( \frac{n-1}{n-k} \right) \left( \frac{G}{G-1} \right). \quad (4.54)$$

The factor  $G/(G-1)$  was derived by Chris Hansen (2007) in the context of equal-sized clusters to improve performance when the number of clusters  $G$  is small. The factor  $(n-1)/(n-k)$  is an

ad hoc generalization which nests the adjustment used in (4.38), since when  $G = n$  we have the simplification  $a_n = n/(n - k)$ .

Alternative cluster-robust covariance matrix estimators can be constructed using cluster-level prediction errors such as

$$\tilde{e}_g = \mathbf{y}_g - \mathbf{X}_g \hat{\boldsymbol{\beta}}_{-g}$$

where  $\hat{\boldsymbol{\beta}}_{-g}$  is the least-squares estimator omitting cluster  $g$ . We then have the robust covariance matrix estimator

$$\tilde{\mathbf{V}}_{\hat{\boldsymbol{\beta}}} = (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{g=1}^G \mathbf{X}'_g \tilde{e}_g \tilde{e}'_g \mathbf{X}_g \right) (\mathbf{X}'\mathbf{X})^{-1}.$$

Similarly to the heteroskedastic-robust case, you can show that  $\tilde{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}$  is a conservative estimator for  $\mathbf{V}_{\hat{\boldsymbol{\beta}}}$  in the sense that the conditional expectation of  $\tilde{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}$  exceeds  $\mathbf{V}_{\hat{\boldsymbol{\beta}}}$ . This covariance matrix estimator is more cumbersome to implement, however, as the cluster-level prediction errors do not have a simple computational form so require a loop to estimate.

To illustrate in the context of the Kenyan schooling example, we present the regression of student test scores on the school-level tracking dummy, with two standard errors displayed. The first (in parenthesis) is the conventional robust standard error. The second [in square brackets] is the clustered standard error, where clustering is at the level of the school.

$$\begin{array}{rcccl} \textit{TestScore}_{ig} = & - & 0.082 & + & 0.147 & \textit{Tracking}_g + e_{ig} & (4.55) \\ & & (0.020) & & (0.028) & & \\ & & [0.054] & & [0.077] & & \end{array}$$

We can see that the cluster-robust standard errors are roughly three times the conventional robust standard errors. Consequently, confidence intervals for the coefficients are greatly affected by the choice.

For illustration, we list here the commands needed to produce the regression results with clustered standard errors in Stata, R, and MATLAB.

<b>Stata do File</b>	
*	Load data:
	use "DDK2011.dta"
*	Standard the test score variable to have mean zero and unit variance:
	egen testscore = std(totalscore)
*	Regression with standard errors clustered at the school level:
	reg testscore tracking, cluster(schoolid)

You can see that clustered standard errors are simple to calculate in Stata.



**R Program File**

```

# Load the data and create variables
data <- read.table("DDK2011.txt",header=TRUE,sep="\t")
y <- scale(as.matrix(data$totalscore))
n <- nrow(y)
x <- cbind(as.matrix(data$tracking),matrix(1,n,1))
schoolid <- as.matrix(data$schoolid)
k <- ncol(x)
invx <- solve(t(x)%*%x)
beta <- invx%*%t(x)%*%y
xe <- x*rep(y-x%*%beta,times=k)
# Clustered robust standard error
xe_sum <- rowsum(xe,schoolid)
G <- nrow(xe_sum)
omega <- t(xe_sum)%*%xe_sum
scale <- G/(G-1)*(n-1)/(n-k)
V_clustered = scale*invx%*%omega%*%invx
se_clustered <- sqrt(diag(V_clustered))
print(beta)
print(se_clustered)

```

Programming clustered standard errors in R is also relatively easy due to the convenient **rowsum** command, which sums variables within clusters.

**MATLAB Program File**

```

% Load the data and create variables
data = xlsread('DDK2011.xlsx');
schoolid = data(:,2);
tracking = data(:,7);
totalscore = data(:,62);
y = (totalscore - mean(totalscore))./std(totalscore);
x = [tracking,ones(size(y,1),1)];
[n,k] = size(x);
invx = inv(x'*x);
beta = invx*(x'*y);
e = y - x*beta;
% Clustered robust standard error
[schools,~,schoolidx] = unique(schoolid);
G = size(schools,1);
cluster_sums = zeros(G,k);
for j = 1:k
    cluster_sums(:,j) = accumarray(schoolidx,x(:,j).*e);end
omega = cluster_sums'*cluster_sums;
scale = G/(G-1)*(n-1)/(n-k);
V_clustered = scale*invx*omega*invx;
se_clustered = sqrt(diag(V_clustered));
display(beta);
display(se_clustered);

```

Here we see that programming clustered standard errors in MATLAB is less convenient than the other packages, but still can be executed with just a few lines of code. This example uses the `accumarray` command, which is similar to the `rowsum` command in R, but only can be applied to vectors (hence the loop across the regressors) and works best if the clusterid variable are indices (which is why the original *schoolid* variable is transformed into indices in *schoolidx*. Application of these commands requires considerable care and attention.

## 4.21 Inference with Clustered Samples

In this section we give some cautionary remarks and general advice about cluster-robust inference in econometric practice. There has been remarkably little theoretical research about the properties of cluster-robust methods – until quite recently – so these remarks may become dated rather quickly.

In many respects cluster-robust inference should be viewed similarly to heteroskedasticity-robust inference, with where a “cluster” in the cluster-robust case is interpreted similarly to an “observation” in the heteroskedasticity-robust case. In particular, the effective sample size should be viewed as the number of clusters, not the “sample size”  $n$ . This is because the cluster-robust covariance matrix estimator effectively treats each cluster as a single observation, and estimates the covariance matrix based on the variation across cluster means. Hence if there are only  $G = 50$  clusters, inference should be viewed as (at best) similar to heteroskedasticity-robust inference with  $n = 50$  observations. This is a bit unsettling, for if the number of regressors is large (say  $k = 20$ ), then the covariance matrix will be estimated quite imprecisely.

Furthermore, most cluster-robust theory (for example, the work of Chris Hansen (2007)) assumes that the clusters are homogeneous, including the assumption that the cluster sizes are all identical. This turns out to be a very important simplification. When this is violated – when, for example, cluster sizes are highly heterogeneous – this should be viewed as roughly equivalent to the heteroskedasticity-robust case with an extremely high degree of heteroskedasticity. If observations themselves are i.i.d. then cluster sums have variances which are proportional to the cluster sizes, so if the latter is heterogeneous so will be the variances of the cluster sums. This also has a large effect on finite sample inference. When clusters are heterogeneous then cluster-robust inference is similar to heteroskedasticity-robust inference with highly heteroskedastic observations.

Put together, if the number of clusters  $G$  is small and the number of observations per cluster is highly varied, then we should interpret inferential statements with a great degree of caution. Unfortunately, this is the norm. Many empirical studies on U.S. data cluster at the “state” level, meaning that there are 50 or 51 clusters (the District of Columbia is typically treated as a state). The number of observations vary considerably across states, since the populations are highly unequal. Thus when you read empirical papers with individual-level data but clustered at the “state” level you should be very cautious, and recognize that this is equivalent to inference with a small number of extremely heterogeneous observations.

A further complication occurs when we are interested in treatment, as in the tracking example given in the previous section. In many cases (including Duflo, Dupas and Kremer (2011)) the interest is in the effect of a specific treatment which is applied at the cluster level (in their case, treatment applies to schools). In many cases (not, however, Duflo, Dupas and Kremer (2011)), the number of treated clusters is small relative to the total number of clusters, in an extreme case there is just a single treated cluster. Based on the reasoning given above, these applications should be interpreted as equivalent to heteroskedasticity-robust inference with a sparse dummy variable, as discussed in Section 4.15. As discussed there, standard error estimates can be erroneously small. In the extreme of a single treated cluster (in the example, if only a single school was tracked) then if the regression is estimated using the pure dummy (no intercept) design, the estimated *tracking* coefficient will have a cluster standard error of 0. In general, reported standard errors will understate the imprecision of parameter estimates.

A practical question which arises in the context of cluster-robust inference is “At what level should we cluster?” In some examples you could cluster at a very fine level, such as families or classrooms, or at higher levels of aggregation, such as neighborhoods, schools, towns, counties, or states. What is the correct level at which to cluster? Rules of thumb have been advocated by practitioners, but at present there is little formal analysis to provide useful guidance. What do we know? If cluster dependence is ignored or imposed at too fine a level, then variance estimators will be biased and inference will be inaccurate. Typically this means that standard errors will be too small, giving rise to spurious indications of significance and precision. On the other hand when cluster-robust inference is based on higher levels of dependence, then the precision of the covariance matrix estimators will decrease, meaning that standard errors will be very imprecise estimates of the actual sampling uncertainty. This means that there is a trade-off between bias and variance in the estimation of the covariance matrix by cluster-robust methods. It is not at all clear – based on current theory – what to do. I state this emphatically. We really do not know what is the “correct” level at which to do cluster-robust inference. This is a very interesting question and should certainly be explored by econometric research.

## Exercises

**Exercise 4.1** For some integer  $k$ , set  $\mu_k = \mathbb{E}(y^k)$ .

- (a) Construct an estimator  $\hat{\mu}_k$  for  $\mu_k$ .
- (b) Show that  $\hat{\mu}_k$  is unbiased for  $\mu_k$ .
- (c) Calculate the variance of  $\hat{\mu}_k$ , say  $\text{var}(\hat{\mu}_k)$ . What assumption is needed for  $\text{var}(\hat{\mu}_k)$  to be finite?
- (d) Propose an estimator of  $\text{var}(\hat{\mu}_k)$ .

**Exercise 4.2** Calculate  $E((\bar{y} - \mu)^3)$ , the skewness of  $\bar{y}$ . Under what condition is it zero?

**Exercise 4.3** Explain the difference between  $\bar{y}$  and  $\mu$ . Explain the difference between  $n^{-1} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'$  and  $\mathbb{E}(\mathbf{x}_i \mathbf{x}_i')$ .

**Exercise 4.4** True or False. If  $y_i = x_i \beta + e_i$ ,  $x_i \in \mathbb{R}$ ,  $\mathbb{E}(e_i | x_i) = 0$ , and  $\hat{e}_i$  is the OLS residual from the regression of  $y_i$  on  $x_i$ , then  $\sum_{i=1}^n x_i^2 \hat{e}_i = 0$ .

**Exercise 4.5** Prove (4.17) and (4.18)

**Exercise 4.6** Prove Theorem 4.8.1.

**Exercise 4.7** Let  $\tilde{\beta}$  be the GLS estimator (4.19) under the assumptions (4.15) and (4.16). Assume that  $\Omega = c^2 \Sigma$  with  $\Sigma$  known and  $c^2$  unknown. Define the residual vector  $\tilde{\mathbf{e}} = \mathbf{y} - \mathbf{X}\tilde{\beta}$ , and an estimator for  $c^2$

$$\tilde{c}^2 = \frac{1}{n-k} \tilde{\mathbf{e}}' \Sigma^{-1} \tilde{\mathbf{e}}.$$

- (a) Show (4.20).
- (b) Show (4.21).
- (c) Prove that  $\tilde{\mathbf{e}} = \mathbf{M}_1 \mathbf{e}$ , where  $\mathbf{M}_1 = \mathbf{I} - \mathbf{X} (\mathbf{X}' \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}' \Sigma^{-1}$ .
- (d) Prove that  $\mathbf{M}_1' \Sigma^{-1} \mathbf{M}_1 = \Sigma^{-1} - \Sigma^{-1} \mathbf{X} (\mathbf{X}' \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}' \Sigma^{-1}$ .
- (e) Find  $\mathbb{E}(\tilde{c}^2 | \mathbf{X})$ .
- (f) Is  $\tilde{c}^2$  a reasonable estimator for  $c^2$ ?

**Exercise 4.8** Let  $(y_i, \mathbf{x}_i)$  be a random sample with  $\mathbb{E}(\mathbf{y} | \mathbf{X}) = \mathbf{X}\beta$ . Consider the **Weighted Least Squares** (WLS) estimator of  $\beta$

$$\tilde{\beta}_{\text{wls}} = (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} (\mathbf{X}' \mathbf{W} \mathbf{y})$$

where  $\mathbf{W} = \text{diag}(w_1, \dots, w_n)$  and  $w_i = x_{ji}^{-2}$ , where  $x_{ji}$  is one of the  $\mathbf{x}_i$ .

- (a) In which contexts would  $\tilde{\beta}_{\text{wls}}$  be a good estimator?
- (b) Using your intuition, in which situations would you expect that  $\tilde{\beta}_{\text{wls}}$  would perform better than OLS?

**Exercise 4.9** Show (4.33) in the homoskedastic regression model.

**Exercise 4.10** Prove (4.41).

**Exercise 4.11** Show (4.42) in the homoskedastic regression model.

**Exercise 4.12** Let  $\mu = \mathbb{E}(y_i)$ ,  $\sigma^2 = \mathbb{E}\left((y_i - \mu)^2\right)$  and  $\mu_3 = \mathbb{E}\left((y_i - \mu)^3\right)$  and consider the sample mean  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ . Find  $\mathbb{E}\left((\bar{y} - \mu)^3\right)$  as a function of  $\mu$ ,  $\sigma^2$ ,  $\mu_3$  and  $n$ .

**Exercise 4.13** Take the simple regression model  $y_i = x_i\beta + e_i$ ,  $x_i \in \mathbb{R}$ ,  $\mathbb{E}(e_i | x_i) = 0$ . Define  $\sigma_i^2 = \mathbb{E}(e_i^2 | x_i)$  and  $\mu_{3i} = \mathbb{E}(e_i^3 | x_i)$  and consider the OLS coefficient  $\hat{\beta}$ . Find  $\mathbb{E}\left(\left(\hat{\beta} - \beta\right)^3 | \mathbf{X}\right)$ .

**Exercise 4.14** Take a regression model with i.i.d. observations  $(y_i, x_i)$  and scalar  $x_i$

$$y_i = x_i\beta + e_i$$

$$\mathbb{E}(e_i | x_i) = 0$$

The parameter of interest is  $\theta = \beta^2$ . Consider the OLS estimates  $\hat{\beta}$  and  $\hat{\theta} = \hat{\beta}^2$ .

- (a) Find  $\mathbb{E}(\hat{\theta} | \mathbf{X})$  using our knowledge of  $\mathbb{E}(\hat{\beta} | \mathbf{X})$  and  $V_{\hat{\beta}} = \text{var}(\hat{\beta} | \mathbf{X})$ . Is  $\hat{\theta}$  biased for  $\theta$ ?
- (b) Suggest an (approximate) biased-corrected estimator  $\hat{\theta}^*$  using an estimate  $\hat{V}_{\hat{\beta}}$  for  $V_{\hat{\beta}}$ .
- (c) For  $\hat{\theta}^*$  to be potentially unbiased, which estimate of  $V_{\hat{\beta}}$  is most appropriate?

Under which conditions is  $\hat{\theta}^*$  unbiased?

**Exercise 4.15** Consider an iid sample  $\{y_i, \mathbf{x}_i\}$   $i = 1, \dots, n$  where  $\mathbf{x}_i$  is  $k \times 1$ . Assume the linear conditional expectation model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(e_i | \mathbf{x}_i) = 0$$

Assume that  $n^{-1} \mathbf{X}' \mathbf{X} = \mathbf{I}_k$  (orthonormal regressors). Consider the OLS estimator  $\hat{\boldsymbol{\beta}}$  for  $\boldsymbol{\beta}$ .

- (a) Find  $\mathbf{V}_{\hat{\boldsymbol{\beta}}} = \text{var}(\hat{\boldsymbol{\beta}})$
- (b) In general, are  $\hat{\beta}_j$  and  $\hat{\beta}_\ell$  for  $j \neq \ell$  correlated or uncorrelated?
- (c) Find a sufficient condition so that  $\hat{\beta}_j$  and  $\hat{\beta}_\ell$  for  $j \neq \ell$  are uncorrelated.

**Exercise 4.16** Take the linear homoskedastic CEF

$$y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + e_i \tag{4.56}$$

$$\mathbb{E}(e_i | \mathbf{x}_i) = 0$$

$$\mathbb{E}(e_i^2 | \mathbf{x}_i) = \sigma^2$$

and suppose that  $y_i^*$  is measured with error. Instead of  $y_i^*$ , we observe  $y_i$  which satisfies

$$y_i = y_i^* + u_i$$

where  $u_i$  is measurement error. Suppose that  $e_i$  and  $u_i$  are independent and

$$\mathbb{E}(u_i | \mathbf{x}_i) = 0$$

$$\mathbb{E}(u_i^2 | \mathbf{x}_i) = \sigma_u^2(\mathbf{x}_i)$$

- (a) Derive an equation for  $y_i$  as a function of  $\mathbf{x}_i$ . Be explicit to write the error term as a function of the structural errors  $e_i$  and  $u_i$ . What is the effect of this measurement error on the model (4.56)?
- (b) Describe the effect of this measurement error on OLS estimation of  $\beta$  in the feasible regression of the observed  $y_i$  on  $\mathbf{x}_i$ .
- (c) Describe the effect (if any) of this measurement error on appropriate standard error calculation for  $\hat{\beta}$ .

**Exercise 4.17** Suppose that for a pair of observables  $(y_i, x_i)$  with  $x_i > 0$  that an economic model implies

$$\mathbb{E}(y_i | x_i) = (\gamma + \theta x_i)^{1/2}. \quad (4.57)$$

A friend suggests that (given an iid sample) you estimate  $\gamma$  and  $\theta$  by the linear regression of  $y_i^2$  on  $x_i$ , that is, to estimate the equation

$$y_i^2 = \alpha + \beta x_i + e_i. \quad (4.58)$$

- (a) Investigate your friend's suggestion. Define  $u_i = y_i - (\gamma + \theta x_i)^{1/2}$ . Show that  $\mathbb{E}(u_i | x_i) = 0$  is implied by (4.57).
- (b) Use  $y_i = (\gamma + \theta x_i)^{1/2} + u_i$  to calculate  $\mathbb{E}(y_i^2 | x_i)$ . What does this tell you about the implied equation (4.58)?
- (c) Can you recover either  $\gamma$  and/or  $\theta$  from estimation of (4.58)? Are additional assumptions required?
- (d) Is this a reasonable suggestion?

**Exercise 4.18** Take the model

$$\begin{aligned} y_i &= \mathbf{x}'_{1i} \beta_1 + \mathbf{x}'_{2i} \beta_2 + e_i \\ \mathbb{E}(e_i | \mathbf{x}_i) &= 0 \\ \mathbb{E}(e_i^2 | \mathbf{x}_i) &= \sigma^2 \end{aligned}$$

where  $\mathbf{x}_i = (\mathbf{x}_{1i}, \mathbf{x}_{2i})$ , with  $\mathbf{x}_{1i}$   $k_1 \times 1$  and  $\mathbf{x}_{2i}$   $k_2 \times 1$ . Consider the short regression

$$y_i = \mathbf{x}'_{1i} \hat{\beta}_1 + \hat{e}_i$$

and define the error variance estimator

$$s^2 = \frac{1}{n - k_1} \sum_{i=1}^n \hat{e}_i^2.$$

Find  $\mathbb{E}(s^2 | \mathbf{X})$

**Exercise 4.19** Let  $\mathbf{y}$  be  $n \times 1$ ,  $\mathbf{X}$  be  $n \times k$ , and  $\mathbf{X}^* = \mathbf{X}\mathbf{C}$  where  $\mathbf{C}$  is  $k \times k$  and full-rank. Let  $\hat{\beta}$  be the least-squares estimator from the regression of  $\mathbf{y}$  on  $\mathbf{X}$ , and let  $\hat{\mathbf{V}}$  be the estimate of its asymptotic covariance matrix. Let  $\hat{\beta}^*$  and  $\hat{\mathbf{V}}^*$  be those from the regression of  $\mathbf{y}$  on  $\mathbf{X}^*$ . Derive an expression for  $\hat{\mathbf{V}}^*$  as a function of  $\hat{\mathbf{V}}$ .

**Exercise 4.20** Take the model

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{e} \\ \mathbb{E}(\mathbf{e} \mid \mathbf{X}) &= \mathbf{0} \\ \mathbb{E}(\mathbf{e}\mathbf{e}' \mid \mathbf{X}) &= \boldsymbol{\Omega} \end{aligned}$$

Assume for simplicity that  $\boldsymbol{\Omega}$  is known. Consider the OLS and GLS estimators  $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{y})$  and  $\tilde{\boldsymbol{\beta}} = (\mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{X})^{-1}(\mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{y})$ . Compute the (conditional) covariance between  $\hat{\boldsymbol{\beta}}$  and  $\tilde{\boldsymbol{\beta}}$ :

$$\mathbb{E}\left((\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mid \mathbf{X}\right)$$

Find the (conditional) covariance matrix for  $\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}$ :

$$\mathbb{E}\left((\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}})(\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}})' \mid \mathbf{X}\right)$$

**Exercise 4.21** The model is

$$\begin{aligned} y_i &= \mathbf{x}_i'\boldsymbol{\beta} + e_i \\ \mathbb{E}(e_i \mid \mathbf{x}_i) &= 0 \\ \mathbb{E}(e_i^2 \mid \mathbf{x}_i) &= \sigma_i^2 \\ \boldsymbol{\Omega} &= \text{diag}(\sigma_1^2, \dots, \sigma_n^2). \end{aligned}$$

The parameter  $\boldsymbol{\beta}$  is estimated both by OLS  $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$  and GLS  $\tilde{\boldsymbol{\beta}} = (\mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{y}$ . Let  $\hat{\mathbf{e}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  and  $\tilde{\mathbf{e}} = \mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}$  denote the residuals. Let  $\hat{R}^2 = 1 - \hat{\mathbf{e}}'\hat{\mathbf{e}}/(\mathbf{y}^*\mathbf{y}^*)$  and  $\tilde{R}^2 = 1 - \tilde{\mathbf{e}}'\tilde{\mathbf{e}}/(\mathbf{y}^*\mathbf{y}^*)$  denote the equation  $R^2$  where  $\mathbf{y}^* = \mathbf{y} - \bar{\mathbf{y}}$ . If the error  $e_i$  is truly heteroskedastic will  $\hat{R}^2$  or  $\tilde{R}^2$  be smaller?

**Exercise 4.22** An economist friend tells you that the assumption that the observations  $(y_i, \mathbf{x}_i)$  are i.i.d. implies that the regression  $y_i = \mathbf{x}_i'\boldsymbol{\beta} + e_i$  is homoskedastic. Do you agree with your friend? How would you explain your position?

**Exercise 4.23** Take the linear regression model with  $\mathbb{E}(\mathbf{y} \mid \mathbf{X}) = \mathbf{X}\boldsymbol{\beta}$ . Define the *ridge regression* estimator

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X} + \mathbf{I}_k\lambda)^{-1}\mathbf{X}'\mathbf{y}$$

where  $\lambda > 0$  is a fixed constant. Find  $E(\hat{\boldsymbol{\beta}} \mid \mathbf{X})$ . Is  $\hat{\boldsymbol{\beta}}$  biased for  $\boldsymbol{\beta}$ ?

**Exercise 4.24** Continue the empirical analysis in Exercise 3.22.

- Calculate standard errors using the homoskedasticity formula and using the four covariance matrices from Section 4.13.
- Repeat in your second programming language. Are they identical?

**Exercise 4.25** Continue the empirical analysis in Exercise 3.24. Calculate standard errors using the Horn-Horn-Duncan method. Repeat in your second programming language. Are they identical?

**Exercise 4.26** Extend the empirical analysis reported in Section 4.20. Do a regression of standardized test score (*totalscore* normalized to have zero mean and variance 1) on tracking, age, sex, being assigned to the contract teacher, and student's percentile in the initial distribution. Calculate standard errors using both the conventional robust formula, and clustering based on the school.

- (a) Compare the two sets of standard errors. Which standard error changes the most by clustering? Which changes the least?
- (b) How does the coefficient on *tracking* change by inclusion of the individual controls (in comparison to the results from (4.55))?



## Chapter 5

# Normal Regression and Maximum Likelihood

### 5.1 Introduction

This chapter introduces the normal regression model and the method of maximum likelihood. The normal regression model is a special case of the linear regression model. It is important as normality allows precise distributional characterizations and sharp inferences. It also provides a baseline for comparison with alternative inference methods, such as asymptotic approximations and the bootstrap.

The method of maximum likelihood is a powerful statistical method for parametric models (such as the normal regression model) and is widely used in econometric practice.

### 5.2 The Normal Distribution

We say that a random variable  $X$  has the **standard normal distribution**, or **Gaussian**, written  $X \sim N(0, 1)$ , if it has the density

$$\phi(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right), \quad -\infty < x < \infty. \quad (5.1)$$

The standard normal density is typically written with the symbol  $\phi(x)$  and the corresponding distribution function by  $\Phi(x)$ . It is a valid density function by the following result.

**Theorem 5.2.1**

$$\int_0^\infty \exp(-x^2/2) dx = \sqrt{\frac{\pi}{2}}. \quad (5.2)$$

All moments of the normal distribution are finite. Since the density is symmetric about zero all odd moments are zero. By integration by parts, you can show (see Exercises 5.2 and 5.3) that  $\mathbb{E}(X^2) = 1$  and  $\mathbb{E}(X^4) = 3$ . In fact, for any positive integer  $m$ ,

$$\mathbb{E}(X^{2m}) = (2m-1)!! = (2m-1) \cdot (2m-3) \cdots 1.$$

The notation  $k!! = k \cdot (k-2) \cdots 1$  is known as the **double factorial**. For example,  $\mathbb{E}(X^6) = 15$ ,  $\mathbb{E}(X^8) = 105$ , and  $\mathbb{E}(X^{10}) = 945$ .

We say that  $X$  has a **univariate normal distribution**, written  $X \sim N(\mu, \sigma^2)$ , if it has the density

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad -\infty < x < \infty.$$

The mean and variance of  $X$  are  $\mu$  and  $\sigma^2$ , respectively.

We say that the  $k$ -vector  $\mathbf{X}$  has a **multivariate normal distribution**, written  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , if it has the joint density

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{k/2} \det(\boldsymbol{\Sigma})^{1/2}} \exp\left(-\frac{(\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}{2}\right), \quad \mathbf{x} \in \mathbb{R}^k.$$

The mean and covariance matrix of  $\mathbf{X}$  are  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ , respectively. By setting  $k = 1$  you can check that the multivariate normal simplifies to the univariate normal.

For technical purposes it is useful to know the form of the moment generating and characteristic functions.

**Theorem 5.2.2** *If  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  then its moment generating function is*

$$M(\mathbf{t}) = \mathbb{E}(\exp(\mathbf{t}'\mathbf{X})) = \exp\left(\mathbf{t}'\boldsymbol{\mu} + \frac{1}{2}\mathbf{t}'\boldsymbol{\Sigma}\mathbf{t}\right)$$

*(see Exercise 5.8) and its characteristic function is*

$$C(\mathbf{t}) = \mathbb{E}(\exp(i\mathbf{t}'\mathbf{X})) = \exp\left(i\mathbf{t}'\boldsymbol{\mu} - \frac{1}{2}\mathbf{t}'\boldsymbol{\Sigma}\mathbf{t}\right)$$

*(see Exercise 5.9).*

An important property of normal random vectors is that affine functions are also multivariate normal.

**Theorem 5.2.3** *If  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and  $\mathbf{Y} = \mathbf{a} + \mathbf{B}\mathbf{X}$ , then  $\mathbf{Y} \sim N(\mathbf{a} + \mathbf{B}\boldsymbol{\mu}, \mathbf{B}\boldsymbol{\Sigma}\mathbf{B}')$ .*

One simple implication of Theorem 5.2.3 is that if  $\mathbf{X}$  is multivariate normal, then each component of  $\mathbf{X}$  is univariate normal.

Another useful property of the multivariate normal distribution is that uncorrelatedness is the same as independence. That is, if a vector is multivariate normal, subsets of variables are independent if and only if they are uncorrelated.

**Theorem 5.2.4** *If  $\mathbf{X} = (\mathbf{X}'_1, \mathbf{X}'_2)'$  is multivariate normal,  $\mathbf{X}_1$  and  $\mathbf{X}_2$  are uncorrelated if and only if they are independent.*

The normal distribution is frequently used for inference to calculate critical values and p-values. This involves evaluating the normal cdf  $\Phi(x)$  and its inverse. Since the cdf  $\Phi(x)$  is not available in closed form, statistical textbooks have traditionally provided tables for this purpose. Such tables are not used currently as now these calculations are embedded in statistical software. For convenience, we list the appropriate commands in MATLAB and R to compute the cumulative distribution function of commonly used statistical distributions.

Numerical Cumulative Distribution Function Calculation			
To calculate $\Pr(X \leq x)$ for given $x$			
	MATLAB	R	Stata
$N(0, 1)$	<code>normcdf(x)</code>	<code>pnorm(x)</code>	<code>normal(x)</code>
$\chi_r^2$	<code>chi2cdf(x,r)</code>	<code>pchisq(x,r)</code>	<code>chi2(r,x)</code>
$t_r$	<code>tcdf(x,r)</code>	<code>pt(x,r)</code>	<code>1-ttail(r,x)</code>
$F_{r,k}$	<code>fcdf(x,r,k)</code>	<code>pf(x,r,k)</code>	<code>F(r,k,x)</code>
$\chi_r^2(d)$	<code>ncx2cdf(x,r,d)</code>	<code>pchisq(x,r,d)</code>	<code>nchi2(r,d,x)</code>
$F_{r,k}(d)$	<code>ncfcdf(x,r,k,d)</code>	<code>pf(x,r,k,d)</code>	<code>1-nFtail(r,k,d,x)</code>

Here we list the appropriate commands to compute the inverse probabilities (quantiles) of the same distributions.

Numerical Quantile Calculation			
To calculate $x$ which solves $p = \Pr(X \leq x)$ for given $p$			
	MATLAB	R	Stata
$N(0, 1)$	<code>norminv(p)</code>	<code>qnorm(p)</code>	<code>invnormal(p)</code>
$\chi_r^2$	<code>chi2inv(p,r)</code>	<code>qchisq(p,r)</code>	<code>invchi2(r,p)</code>
$t_r$	<code>tinv(p,r)</code>	<code>qt(p,r)</code>	<code>invttail(r,1-p)</code>
$F_{r,k}$	<code>finv(p,r,k)</code>	<code>qf(p,r,k)</code>	<code>invF(r,k,p)</code>
$\chi_r^2(d)$	<code>ncx2inv(p,r,d)</code>	<code>qchisq(p,r,d)</code>	<code>invnchi2(r,d,p)</code>
$F_{r,k}(d)$	<code>ncfinv(p,r,k,d)</code>	<code>qf(p,r,k,d)</code>	<code>invnFtail(r,k,d,1-p)</code>

### 5.3 Chi-Square Distribution

Many important distributions can be derived as transformation of multivariate normal random vectors, including the chi-square, the student  $t$ , and the  $F$ . In this section we introduce the chi-square distribution.

Let  $\mathbf{X} \sim N(\mathbf{0}, \mathbf{I}_r)$  be multivariate standard normal and define  $Q = \mathbf{X}'\mathbf{X}$ . The distribution of  $Q$  is called **chi-square** with  $r$  degrees of freedom, written as  $Q \sim \chi_r^2$ .

The mean and variance of  $Q \sim \chi_r^2$  are  $r$  and  $2r$ , respectively. (See Exercise 5.10.)

The chi-square distribution function is frequently used for inference (critical values and p-values). In practice these calculations are performed numerically by statistical software, but for completeness we provide the density function.

**Theorem 5.3.1** *The density of  $\chi_r^2$  is*

$$f(x) = \frac{1}{2^{r/2}\Gamma\left(\frac{r}{2}\right)} x^{r/2-1} e^{-x/2}, \quad x > 0 \quad (5.3)$$

where  $\Gamma(t) = \int_0^\infty u^{t-1} e^{-u} du$  is the gamma function (Section 5.18).

For some theoretical applications, including the study of the power of statistical tests, it is useful to define a non-central version of the chi-square distribution. When  $\mathbf{X} \sim N(\boldsymbol{\mu}, \mathbf{I}_r)$  is multivariate normal, we say that  $Q = \mathbf{X}'\mathbf{X}$  has a **non-central chi-square** distribution, with  $r$  degrees of freedom and non-centrality parameter  $\lambda = \boldsymbol{\mu}'\boldsymbol{\mu}$ , and is written as  $Q \sim \chi_r^2(\lambda)$ . The non-central chi-square simplifies to the central (conventional) chi-square when  $\lambda = 0$ , so that  $\chi_r^2(0) = \chi_r^2$ .

**Theorem 5.3.2** *The density of  $\chi_r^2(\lambda)$  is*

$$f(x) = \sum_{i=0}^{\infty} \frac{e^{-\lambda/2}}{i!} \left(\frac{\lambda}{2}\right)^i f_{r+2i}(x), \quad x > 0 \quad (5.4)$$

where  $f_{r+2i}(x)$  is the  $\chi_{r+2i}^2$  density function (5.3).

Interestingly, as can be seen from the formula (5.4), the distribution of  $\chi_r^2(\lambda)$  only depends on the scalar non-centrality parameter  $\lambda$ , not the entire mean vector  $\boldsymbol{\mu}$ .

A useful fact about the central and non-central chi-square distributions is that they also can be derived from multivariate normal distributions with general covariance matrices.

**Theorem 5.3.3** *If  $\mathbf{X} \sim N(\boldsymbol{\mu}, \mathbf{A})$  with  $\mathbf{A} > 0$ ,  $r \times r$ , then  $\mathbf{X}'\mathbf{A}^{-1}\mathbf{X} \sim \chi_r^2(\lambda)$ , where  $\lambda = \boldsymbol{\mu}'\mathbf{A}^{-1}\boldsymbol{\mu}$ .*

In particular, Theorem 5.3.3 applies to the central chi-squared distribution, so if  $\mathbf{X} \sim N(0, \mathbf{A})$  then  $\mathbf{X}'\mathbf{A}^{-1}\mathbf{X} \sim \chi_r^2$ .

## 5.4 Student t Distribution

Let  $Z \sim N(0, 1)$  and  $Q \sim \chi_r^2$  be independent, and define  $T = Z/\sqrt{Q/r}$ . The distribution of  $T$  is called the **student t** with  $r$  degrees of freedom, and is written  $T \sim t_r$ . Like the chi-square, the distribution only depends on the degree of freedom parameter  $r$ .

**Theorem 5.4.1** *The density of  $T$  is*

$$f(x) = \frac{\Gamma\left(\frac{r+1}{2}\right)}{\sqrt{r\pi}\Gamma\left(\frac{r}{2}\right)} \left(1 + \frac{x^2}{r}\right)^{-\left(\frac{r+1}{2}\right)}, \quad -\infty < x < \infty.$$

The density function of the student  $t$  is bell-shaped like the normal density function, but the  $t$  has thicker tails. The  $t$  distribution has the property that moments below  $r$  are finite, but absolute moments greater than or equal to  $r$  are infinite.

The student  $t$  can also be seen as a generalization of the standard normal, for the latter is obtained as the limiting case where  $r$  is taken to infinity.

**Theorem 5.4.2** *Let  $f_r(x)$  be the student  $t$  density. As  $r \rightarrow \infty$ ,  $f_r(x) \rightarrow \phi(x)$ .*

Another special case of the student  $t$  distribution occurs when  $r = 1$  and is known as the **Cauchy** distribution. The Cauchy density function is

$$f(x) = \frac{1}{\pi(1+x^2)}, \quad -\infty < x < \infty.$$

A Cauchy random variable  $T = Z_1/Z_2$  can also be derived as the ratio of two independent  $N(0, 1)$  variables. The Cauchy has the property that it has no finite integer moments.

### William Gosset

William S. Gosset (1876-1937) of England is most famous for his derivation of the student's  $t$  distribution, published in the paper "The probable error of a mean" in 1908. At the time, Gosset worked at Guinness Brewery, which prohibited its employees from publishing in order to prevent the possible loss of trade secrets. To circumvent this barrier, Gosset published under the pseudonym "Student". Consequently, this famous distribution is known as the student  $t$  rather than Gosset's  $t$ !

## 5.5 F Distribution

Let  $Q_m \sim \chi_m^2$  and  $Q_r \sim \chi_r^2$  be independent. The distribution of  $F = (Q_m/m) / (Q_r/r)$  is called the  $F$  distribution with degree of freedom parameters  $m$  and  $r$ , and we write  $F \sim F_{m,r}$ .

**Theorem 5.5.1** *The density of  $F$  is*

$$f(x) = \frac{\left(\frac{m}{r}\right)^{m/2} x^{m/2-1} \Gamma\left(\frac{m+r}{2}\right)}{\Gamma\left(\frac{m}{2}\right) \Gamma\left(\frac{r}{2}\right) \left(1 + \frac{m}{r}x\right)^{(m+r)/2}}, \quad x > 0.$$

If  $m = 1$  then we can write  $Q_1 = Z^2$  where  $Z \sim N(0, 1)$ , and  $F = Z^2 / (Q_r/r) = \left(Z / \sqrt{Q_r/r}\right)^2 = T^2$ , the square of a student  $t$  with  $r$  degree of freedom. Thus the  $F$  distribution with  $m = 1$  is equal to the squared student  $t$  distribution. In this sense the  $F$  distribution is a generalization of the student  $t$ .

As a limiting case, as  $r \rightarrow \infty$  the  $F$  distribution simplifies to  $F \rightarrow Q_m/m$ , a normalized  $\chi_m^2$ . Thus the  $F$  distribution is also a generalization of the  $\chi_m^2$  distribution.

**Theorem 5.5.2** *Let  $f_{m,r}(x)$  be the density of  $mF$ . As  $r \rightarrow \infty$ ,  $f_{m,r}(x) \rightarrow f_m(x)$ , the density of  $\chi_m^2$ .*

Similarly with the non-central chi-square we define the non-central  $F$  distribution. If  $Q_m \sim \chi_m^2(\lambda)$  and  $Q_r \sim \chi_r^2$  are independent, then  $F = (Q_m/m) / (Q_r/r)$  is called a **non-central**  $F$  with degree of freedom parameters  $m$  and  $r$  and non-centrality parameter  $\lambda$ .

## 5.6 Joint Normality and Linear Regression

Suppose the variables  $(y, \mathbf{x})$  are jointly normally distributed. Consider the best linear predictor of  $y$  given  $\mathbf{x}$

$$y = \mathbf{x}'\boldsymbol{\beta} + \alpha + e.$$

By the properties of the best linear predictor,  $\mathbb{E}(\mathbf{x}e) = 0$  and  $\mathbb{E}(e) = 0$ , so  $\mathbf{x}$  and  $e$  are uncorrelated. Since  $(e, \mathbf{x})$  is an affine transformation of the normal vector  $(y, \mathbf{x})$ , it follows that  $(e, \mathbf{x})$  is jointly normal (Theorem 5.2.3). Since  $(e, \mathbf{x})$  is jointly normal and uncorrelated they are independent (Theorem 5.2.4). Independence implies that

$$\mathbb{E}(e | \mathbf{x}) = \mathbb{E}(e) = 0$$

and

$$\mathbb{E}(e^2 | \mathbf{x}) = \mathbb{E}(e^2) = \sigma^2$$

which are properties of a homoskedastic linear CEF.

We have shown that when  $(y, \mathbf{x})$  are jointly normally distributed, they satisfy a normal linear CEF

$$y = \mathbf{x}'\boldsymbol{\beta} + \alpha + e$$

where

$$e \sim N(0, \sigma^2)$$

is independent of  $\mathbf{x}$ .

This is a classical motivation for the linear regression model.

## 5.7 Normal Regression Model

The normal regression model is the linear regression model with an independent normal error

$$\begin{aligned} y &= \mathbf{x}'\boldsymbol{\beta} + e \\ e &\sim N(0, \sigma^2). \end{aligned} \tag{5.5}$$

As we learned in Section 5.6, the normal regression model holds when  $(y, \mathbf{x})$  are jointly normally distributed. Normal regression, however, does not require joint normality. All that is required is that the conditional distribution of  $y$  given  $\mathbf{x}$  is normal (the marginal distribution of  $\mathbf{x}$  is unrestricted). In this sense the normal regression model is broader than joint normality. Notice that for notational convenience we have written (5.5) so that  $\mathbf{x}$  contains the intercept.

Normal regression is a parametric model, where likelihood methods can be used for estimation, testing, and distribution theory. The **likelihood** is the name for the joint probability density of the data, evaluated at the observed sample, and viewed as a function of the parameters. The maximum likelihood estimator is the value which maximizes this likelihood function. Let us now derive the likelihood of the normal regression model.

First, observe that model (5.5) is equivalent to the statement that the conditional density of  $y$  given  $\mathbf{x}$  takes the form

$$f(y | \mathbf{x}) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{1}{2\sigma^2} (y - \mathbf{x}'\boldsymbol{\beta})^2\right).$$

Under the assumption that the observations are mutually independent, this implies that the conditional density of  $(y_1, \dots, y_n)$  given  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$  is

$$\begin{aligned} f(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n) &= \prod_{i=1}^n f(y_i | \mathbf{x}_i) \\ &= \prod_{i=1}^n \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{1}{2\sigma^2} (y_i - \mathbf{x}_i'\boldsymbol{\beta})^2\right) \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i'\boldsymbol{\beta})^2\right) \\ &\stackrel{\text{def}}{=} L(\boldsymbol{\beta}, \sigma^2) \end{aligned}$$

and is called the **likelihood function**.

For convenience, it is typical to work with the natural logarithm

$$\begin{aligned} \log f(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n) &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i'\boldsymbol{\beta})^2 \\ &\stackrel{\text{def}}{=} \log L(\boldsymbol{\beta}, \sigma^2) \end{aligned} \tag{5.6}$$

which is called the **log-likelihood function**.

The **maximum likelihood estimator (MLE)**  $(\hat{\boldsymbol{\beta}}_{\text{mle}}, \hat{\sigma}_{\text{mle}}^2)$  is the value which maximizes the log-likelihood. (It is equivalent to maximize the likelihood or the log-likelihood. See Exercise 5.15.) We can write the maximization problem as

$$(\hat{\boldsymbol{\beta}}_{\text{mle}}, \hat{\sigma}_{\text{mle}}^2) = \underset{\boldsymbol{\beta} \in \mathbb{R}^k, \sigma^2 > 0}{\operatorname{argmax}} \log L(\boldsymbol{\beta}, \sigma^2). \tag{5.7}$$

In most applications of maximum likelihood, the MLE must be found by numerical methods. However, in the case of the normal regression model we can find an explicit expression for  $\hat{\boldsymbol{\beta}}_{\text{mle}}$  and  $\hat{\sigma}_{\text{mle}}^2$  as functions of the data.

The maximizers  $(\hat{\boldsymbol{\beta}}_{\text{mle}}, \hat{\sigma}_{\text{mle}}^2)$  of (5.7) jointly solve the first-order conditions (FOC)

$$0 = \frac{\partial}{\partial \boldsymbol{\beta}} \log L(\boldsymbol{\beta}, \sigma^2) \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}_{\text{mle}}, \sigma^2=\hat{\sigma}_{\text{mle}}^2} = \frac{1}{\hat{\sigma}_{\text{mle}}^2} \sum_{i=1}^n \mathbf{x}_i (y_i - \mathbf{x}_i'\hat{\boldsymbol{\beta}}_{\text{mle}}) \tag{5.8}$$

$$0 = \frac{\partial}{\partial \sigma^2} \log L(\boldsymbol{\beta}, \sigma^2) \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}_{\text{mle}}, \sigma^2=\hat{\sigma}_{\text{mle}}^2} = -\frac{n}{2\hat{\sigma}_{\text{mle}}^2} + \frac{1}{\hat{\sigma}_{\text{mle}}^4} \sum_{i=1}^n (y_i - \mathbf{x}_i'\hat{\boldsymbol{\beta}}_{\text{mle}})^2. \tag{5.9}$$

The first FOC (5.8) is proportional to the first-order conditions for the least-squares minimization problem of Section 3.6. It follows that the MLE satisfies

$$\hat{\boldsymbol{\beta}}_{\text{mle}} = \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left( \sum_{i=1}^n \mathbf{x}_i y_i \right) = \hat{\boldsymbol{\beta}}_{\text{ols}}.$$

That is, the MLE for  $\boldsymbol{\beta}$  is algebraically identical to the OLS estimator.

Solving the second FOC (5.9) for  $\hat{\sigma}_{\text{mle}}^2$  we find

$$\hat{\sigma}_{\text{mle}}^2 = \frac{1}{n} \sum_{i=1}^n \left( y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{\text{mle}} \right)^2 = \frac{1}{n} \sum_{i=1}^n \left( y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{\text{ols}} \right)^2 = \frac{1}{n} \sum_{i=1}^n \hat{e}_i^2 = \hat{\sigma}_{\text{ols}}^2.$$

Thus the MLE for  $\sigma^2$  is identical to the OLS/moment estimator from (3.33).

Since the OLS estimate and MLE under normality are equivalent,  $\hat{\boldsymbol{\beta}}$  is described by some authors as the maximum likelihood estimator, and by other authors as the least-squares estimator. It is important to remember, however, that  $\hat{\boldsymbol{\beta}}$  is only the MLE when the error  $e$  has a known normal distribution, and not otherwise.

Plugging the estimators into (5.6) we obtain the maximized log-likelihood

$$\log L \left( \hat{\boldsymbol{\beta}}_{\text{mle}}, \hat{\sigma}_{\text{mle}}^2 \right) = -\frac{n}{2} \log (2\pi \hat{\sigma}_{\text{mle}}^2) - \frac{n}{2}. \quad (5.10)$$

The log-likelihood is typically reported as a measure of fit.

It may seem surprising that the MLE  $\hat{\boldsymbol{\beta}}_{\text{mle}}$  is numerically equal to the OLS estimator, despite emerging from quite different motivations. It is not completely accidental. The least-squares estimator minimizes a particular sample loss function – the sum of squared error criterion – and most loss functions are equivalent to the likelihood of a specific parametric distribution, in this case the normal regression model. In this sense it is not surprising that the least-squares estimator can be motivated as either the minimizer of a sample loss function or as the maximizer of a likelihood function.

### Carl Friedrich Gauss

The mathematician Carl Friedrich Gauss (1777-1855) proposed the normal regression model, and derived the least squares estimator as the maximum likelihood estimator for this model. He claimed to have discovered the method in 1795 at the age of eighteen, but did not publish the result until 1809. Interest in Gauss's approach was reinforced by Laplace's simultaneous discovery of the central limit theorem, which provided a justification for viewing random disturbances as approximately normal.

## 5.8 Distribution of OLS Coefficient Vector

In the normal linear regression model we can derive exact sampling distributions for the OLS/MLE estimates, residuals, and variance estimate. In this section we derive the distribution of the OLS coefficient estimate.

The normality assumption  $e_i \mid \mathbf{x}_i \sim N(0, \sigma^2)$  combined with independence of the observations has the multivariate implication

$$\mathbf{e} \mid \mathbf{X} \sim N(\mathbf{0}, \mathbf{I}_n \sigma^2).$$

That is, the error vector  $\mathbf{e}$  is independent of  $\mathbf{X}$  and is normally distributed.

Recall that the OLS estimator satisfies

$$\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{e}$$



which is a linear function of  $\mathbf{e}$ . Since linear functions of normals are also normal (Theorem 5.2.3), this implies that conditional on  $\mathbf{X}$ ,

$$\begin{aligned}\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta} \Big|_{\mathbf{X}} &\sim (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{N}(0, \mathbf{I}_n \sigma^2) \\ &\sim \mathbf{N}\left(0, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1}\right) \\ &= \mathbf{N}\left(0, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}\right).\end{aligned}$$

An alternative way of writing this is

$$\widehat{\boldsymbol{\beta}} \Big|_{\mathbf{X}} \sim \mathbf{N}\left(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}\right).$$

This shows that under the assumption of normal errors, the OLS estimate has an exact normal distribution.

**Theorem 5.8.1** *In the linear regression model,*

$$\widehat{\boldsymbol{\beta}} \Big|_{\mathbf{X}} \sim \mathbf{N}\left(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}\right).$$

Theorems 5.2.3 and 5.8.1 imply that any affine function of the OLS estimate is also normally distributed, including individual estimates. Letting  $\beta_j$  and  $\widehat{\beta}_j$  denote the  $j^{\text{th}}$  elements of  $\boldsymbol{\beta}$  and  $\widehat{\boldsymbol{\beta}}$ , we have

$$\widehat{\beta}_j \Big|_{\mathbf{X}} \sim \mathbf{N}\left(\beta_j, \sigma^2 \left[(\mathbf{X}'\mathbf{X})^{-1}\right]_{jj}\right). \quad (5.11)$$

## 5.9 Distribution of OLS Residual Vector

Now consider the OLS residual vector. Recall from (3.31) that  $\widehat{\mathbf{e}} = \mathbf{M}\mathbf{e}$  where  $\mathbf{M} = \mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ . This shows that  $\widehat{\mathbf{e}}$  is linear in  $\mathbf{e}$ . So conditional on  $\mathbf{X}$ ,

$$\widehat{\mathbf{e}} = \mathbf{M}\mathbf{e} \Big|_{\mathbf{X}} \sim \mathbf{N}(0, \sigma^2 \mathbf{M}\mathbf{M}) = \mathbf{N}(0, \sigma^2 \mathbf{M})$$

the final equality since  $\mathbf{M}$  is idempotent (see Section 3.12). This shows that the residual vector has an exact normal distribution.

Furthermore, it is useful to understand the joint distribution of  $\widehat{\boldsymbol{\beta}}$  and  $\widehat{\mathbf{e}}$ . This is easiest done by writing the two as a stacked linear function of the error  $\mathbf{e}$ . Indeed,

$$\begin{pmatrix} \widehat{\boldsymbol{\beta}} - \boldsymbol{\beta} \\ \widehat{\mathbf{e}} \end{pmatrix} = \begin{pmatrix} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{e} \\ \mathbf{M}\mathbf{e} \end{pmatrix} = \begin{pmatrix} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \\ \mathbf{M} \end{pmatrix} \mathbf{e}$$

which is a linear function of  $\mathbf{e}$ . The vector thus has a joint normal distribution with covariance matrix

$$\begin{pmatrix} \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} & 0 \\ 0 & \sigma^2 \mathbf{M} \end{pmatrix}.$$

The covariance is zero because  $(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{M} = 0$  as  $\mathbf{X}'\mathbf{M} = 0$  from (3.28). Since the covariance is zero, it follows that  $\widehat{\boldsymbol{\beta}}$  and  $\widehat{\mathbf{e}}$  are statistically independent (Theorem 5.2.4).

**Theorem 5.9.1** *In the linear regression model,  $\widehat{\mathbf{e}} \Big|_{\mathbf{X}} \sim \mathbf{N}(0, \sigma^2 \mathbf{M})$  and is independent of  $\widehat{\boldsymbol{\beta}}$ .*

The fact that  $\widehat{\boldsymbol{\beta}}$  and  $\widehat{\mathbf{e}}$  are independent implies that  $\widehat{\boldsymbol{\beta}}$  is independent of any function of the residual vector, including individual residuals  $\widehat{e}_i$  and the variance estimate  $s^2$  and  $\widehat{\sigma}^2$ .

## 5.10 Distribution of Variance Estimate

Next, consider the variance estimator  $s^2$  from (4.30). Using (3.35), it satisfies  $(n - k) s^2 = \widehat{\mathbf{e}}' \widehat{\mathbf{e}} = \mathbf{e}' \mathbf{M} \mathbf{e}$ . The spectral decomposition of  $\mathbf{M}$  (see equation (A.10)) is  $\mathbf{M} = \mathbf{H} \mathbf{\Lambda} \mathbf{H}'$  where  $\mathbf{H}' \mathbf{H} = \mathbf{I}_n$  and  $\mathbf{\Lambda}$  is diagonal with the eigenvalues of  $\mathbf{M}$  on the diagonal. Since  $\mathbf{M}$  is idempotent with rank  $n - k$  (see Section 3.12) it has  $n - k$  eigenvalues equalling 1 and  $k$  eigenvalues equalling 0, so

$$\mathbf{\Lambda} = \begin{bmatrix} \mathbf{I}_{n-k} & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_k \end{bmatrix}.$$

Let  $\mathbf{u} = \mathbf{H}' \mathbf{e} \sim N(\mathbf{0}, \mathbf{I}_n \sigma^2)$  (see Exercise 5.13) and partition  $\mathbf{u} = (\mathbf{u}'_1, \mathbf{u}'_2)'$  where  $\mathbf{u}_1 \sim N(\mathbf{0}, \mathbf{I}_{n-k} \sigma^2)$ . Then

$$\begin{aligned} (n - k) s^2 &= \mathbf{e}' \mathbf{M} \mathbf{e} \\ &= \mathbf{e}' \mathbf{H} \begin{bmatrix} \mathbf{I}_{n-k} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{H}' \mathbf{e} \\ &= \mathbf{u}' \begin{bmatrix} \mathbf{I}_{n-k} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{u} \\ &= \mathbf{u}'_1 \mathbf{u}_1 \\ &\sim \sigma^2 \chi^2_{n-k}. \end{aligned}$$

We see that in the normal regression model the exact distribution of  $s^2$  is a scaled chi-square.

Since  $\widehat{\mathbf{e}}$  is independent of  $\widehat{\boldsymbol{\beta}}$  it follows that  $s^2$  is independent of  $\widehat{\boldsymbol{\beta}}$  as well.

**Theorem 5.10.1** *In the linear regression model,*

$$\frac{(n - k) s^2}{\sigma^2} \sim \chi^2_{n-k}$$

*and is independent of  $\widehat{\boldsymbol{\beta}}$ .*

## 5.11 t-statistic

An alternative way of writing (5.11) is

$$\frac{\widehat{\beta}_j - \beta_j}{\sqrt{\sigma^2 [(\mathbf{X}' \mathbf{X})^{-1}]_{jj}}} \sim N(0, 1).$$

This is sometimes called a **standardized** statistic, as the distribution is the standard normal.

Now take the standardized statistic and replace the unknown variance  $\sigma^2$  with its estimate  $s^2$ . We call this a **t-ratio** or **t-statistic**

$$T = \frac{\widehat{\beta}_j - \beta_j}{\sqrt{s^2 [(\mathbf{X}' \mathbf{X})^{-1}]_{jj}}} = \frac{\widehat{\beta}_j - \beta_j}{s(\widehat{\beta}_j)}$$

where  $s(\widehat{\beta}_j)$  is the classical (homoskedastic) standard error for  $\widehat{\beta}_j$  from (4.43). We will sometimes write the t-statistic as  $T(\beta_j)$  to explicitly indicate its dependence on the parameter value  $\beta_j$ , and

sometimes will simplify notation and write the t-statistic as  $T$  when the dependence is clear from the context.

By some algebraic re-scaling we can write the t-statistic as the ratio of the standardized statistic and the square root of the scaled variance estimate. Since the distributions of these two components are normal and chi-square, respectively, and independent, then we can deduce that the t-statistic has the distribution

$$\begin{aligned} T &= \frac{\hat{\beta}_j - \beta_j}{\sqrt{\sigma^2 [(\mathbf{X}'\mathbf{X})^{-1}]_{jj}}} \bigg/ \sqrt{\frac{(n-k)s^2}{\sigma^2}} \bigg/ (n-k) \\ &\sim \frac{N(0,1)}{\sqrt{\chi_{n-k}^2 / (n-k)}} \\ &\sim t_{n-k} \end{aligned}$$

a student  $t$  distribution with  $n - k$  degrees of freedom.

This derivation shows that the t-ratio has a sampling distribution which depends only on the quantity  $n - k$ . The distribution does not depend on any other features of the data. In this context, we say that the distribution of the t-ratio is **pivotal**, meaning that it does not depend on unknowns.

The trick behind this result is scaling the centered coefficient by its standard error, and recognizing that each depends on the unknown  $\sigma$  only through scale. Thus the ratio of the two does not depend on  $\sigma$ . This trick (scaling to eliminate dependence on unknowns) is known as **studentization**.

**Theorem 5.11.1** *In the normal regression model,  $T \sim t_{n-k}$ .*

An important caveat about Theorem 5.11.1 is that it only applies to the t-statistic constructed with the homoskedastic (old-fashioned) standard error estimate. It does not apply to a t-statistic constructed with any of the robust standard error estimates. In fact, the robust t-statistics can have finite sample distributions which deviate considerably from  $t_{n-k}$  even when the regression errors are independent  $N(0, \sigma^2)$ . Thus the distributional result in Theorem 5.11.1, and the use of the t distribution in finite samples, should only be applied to classical t-statistics.

## 5.12 Confidence Intervals for Regression Coefficients

An OLS estimate  $\hat{\beta}$  is a **point estimate** for a coefficient  $\beta$ . A broader concept is a **set or interval estimate** which takes the form  $\hat{C} = [\hat{L}, \hat{U}]$ . The goal of an interval estimate  $\hat{C}$  is to contain the true value, e.g.  $\beta \in \hat{C}$ , with high probability.

The interval estimate  $\hat{C}$  is a function of the data and hence is random.

An interval estimate  $\hat{C}$  is called a  $1 - \alpha$  **confidence interval** when  $\Pr(\beta \in \hat{C}) = 1 - \alpha$  for a selected value of  $\alpha$ . The value  $1 - \alpha$  is called the **coverage probability**. Typical choices for the coverage probability  $1 - \alpha$  are 0.95 or 0.90.

The probability calculation  $\Pr(\beta \in \hat{C})$  is easily mis-interpreted as treating  $\beta$  as random and  $\hat{C}$  as fixed. (The probability that  $\beta$  is in  $\hat{C}$ .) This is not the appropriate interpretation. Instead, the correct interpretation is that the probability  $\Pr(\beta \in \hat{C})$  treats the point  $\beta$  as fixed and the set  $\hat{C}$  as random. It is the probability that the random set  $\hat{C}$  covers (or contains) the fixed true coefficient  $\beta$ .

There is not a unique method to construct confidence intervals. For example, one simple (yet silly) interval is

$$\hat{C} = \begin{cases} \mathbb{R} & \text{with probability } 1 - \alpha \\ \{\hat{\beta}\} & \text{with probability } \alpha \end{cases}.$$

If  $\hat{\beta}$  has a continuous distribution, then by construction  $\Pr(\beta \in \hat{C}) = 1 - \alpha$ , so this confidence interval has perfect coverage. However,  $\hat{C}$  is uninformative about  $\hat{\beta}$  and is therefore not useful.

Instead, a good choice for a confidence interval for the regression coefficient  $\beta$  is obtained by adding and subtracting from the estimate  $\hat{\beta}$  a fixed multiple of its standard error:

$$\hat{C} = [\hat{\beta} - c \cdot s(\hat{\beta}), \quad \hat{\beta} + c \cdot s(\hat{\beta})] \quad (5.12)$$

where  $c > 0$  is a pre-specified constant. This confidence interval is symmetric about the point estimate  $\hat{\beta}$ , and its length is proportional to the standard error  $s(\hat{\beta})$ .

Equivalently,  $\hat{C}$  is the set of parameter values for  $\beta$  such that the t-statistic  $T(\beta)$  is smaller (in absolute value) than  $c$ , that is

$$\hat{C} = \{\beta : |T(\beta)| \leq c\} = \left\{ \beta : -c \leq \frac{\hat{\beta} - \beta}{s(\hat{\beta})} \leq c \right\}.$$

The coverage probability of this confidence interval is

$$\begin{aligned} \Pr(\beta \in \hat{C}) &= \Pr(|T(\beta)| \leq c) \\ &= \Pr(-c \leq T(\beta) \leq c). \end{aligned} \quad (5.13)$$

Since the t-statistic  $T(\beta)$  has the  $t_{n-k}$  distribution, (5.13) equals  $F(c) - F(-c)$ , where  $F(u)$  is the student  $t$  distribution function with  $n - k$  degrees of freedom. Since  $F(-c) = 1 - F(c)$  (see Exercise 5.19) we can write (5.13) as

$$\Pr(\beta \in \hat{C}) = 2F(c) - 1.$$

This is the **coverage probability** of the interval  $\hat{C}$ , and only depends on the constant  $c$ .

As we mentioned before, a confidence interval has the coverage probability  $1 - \alpha$ . This requires selecting the constant  $c$  so that  $F(c) = 1 - \alpha/2$ . This holds if  $c$  equals the  $1 - \alpha/2$  quantile of the  $t_{n-k}$  distribution. As there is no closed form expression for these quantiles, we compute their values numerically. For example, by `tinvs(1-alpha/2,n-k)` in MATLAB. With this choice the confidence interval (5.12) has exact coverage probability  $1 - \alpha$ . By default, Stata reports 95% confidence intervals  $\hat{C}$  for each estimated regression coefficient using the same formula.

**Theorem 5.12.1** *In the normal regression model, (5.12) with  $c = F^{-1}(1 - \alpha/2)$  has coverage probability  $\Pr(\beta \in \hat{C}) = 1 - \alpha$ .*

When the degree of freedom is large the distinction between the student  $t$  and the normal distribution is negligible. In particular, for  $n - k \geq 61$  we have  $c \leq 2.00$  for a 95% interval. Using this value we obtain the most commonly used confidence interval in applied econometric practice:

$$\hat{C} = [\hat{\beta} - 2s(\hat{\beta}), \quad \hat{\beta} + 2s(\hat{\beta})]. \quad (5.14)$$

This is a useful rule-of-thumb. This 95% confidence interval  $\hat{C}$  is simple to compute and can be easily calculated from coefficient estimates and standard errors.

**Theorem 5.12.2** *In the normal regression model, if  $n-k \geq 61$  then (5.14) has coverage probability  $\Pr(\beta \in \hat{C}) \geq 0.95$ .*

Confidence intervals are a simple yet effective tool to assess estimation uncertainty. When reading a set of empirical results, look at the estimated coefficient estimates and the standard errors. For a parameter of interest, compute the confidence interval  $\hat{C}$  and consider the meaning of the spread of the suggested values. If the range of values in the confidence interval are too wide to learn about  $\beta$ , then do not jump to a conclusion about  $\beta$  based on the point estimate alone.

### 5.13 Confidence Intervals for Error Variance

We can also construct a confidence interval for the regression error variance  $\sigma^2$  using the sampling distribution of  $s^2$  from Theorem 5.10.1, which states that in the normal regression model

$$\frac{(n-k)s^2}{\sigma^2} \sim \chi_{n-k}^2. \quad (5.15)$$

Let  $F(u)$  denote the  $\chi_{n-k}^2$  distribution function, and for some  $\alpha$  set  $c_1 = F^{-1}(\alpha/2)$  and  $c_2 = F^{-1}(1 - \alpha/2)$  (the  $\alpha/2$  and  $1 - \alpha/2$  quantiles of the  $\chi_{n-k}^2$  distribution). Equation (5.15) implies that

$$\Pr\left(c_1 \leq \frac{(n-k)s^2}{\sigma^2} \leq c_2\right) = F(c_2) - F(c_1) = 1 - \alpha.$$

Rewriting the inequalities we find

$$\Pr\left((n-k)s^2/c_2 \leq \sigma^2 \leq (n-k)s^2/c_1\right) = 1 - \alpha.$$

This shows that an exact  $1 - \alpha$  confidence interval for  $\sigma^2$  is

$$C = \left[ \frac{(n-k)s^2}{c_2}, \quad \frac{(n-k)s^2}{c_1} \right]. \quad (5.16)$$

**Theorem 5.13.1** *In the normal regression model, (5.16) has coverage probability  $\Pr(\sigma^2 \in C) = 1 - \alpha$ .*

The confidence interval (5.16) for  $\sigma^2$  is asymmetric about the point estimate  $s^2$ , due to the latter's asymmetric sampling distribution.

### 5.14 t Test

A typical goal in an econometric exercise is to assess whether or not coefficient  $\beta$  equals a specific value  $\beta_0$ . Often the specific value to be tested is  $\beta_0 = 0$  but this is not essential. This is called **hypothesis testing**, a subject which will be explored in detail in Chapter 9. In this section and the following we give a short introduction specific to the normal regression model.

For simplicity write the coefficient to be tested as  $\beta$ . The **null hypothesis** is

$$\mathbb{H}_0 : \beta = \beta_0. \quad (5.17)$$

This states that the hypothesis is that the true value of the coefficient  $\beta$  equals the hypothesized value  $\beta_0$ .

The alternative hypothesis is the complement of  $\mathbb{H}_0$ , and is written as

$$\mathbb{H}_1 : \beta \neq \beta_0.$$

This states that the true value of  $\beta$  does not equal the hypothesized value.

We are interested in testing  $\mathbb{H}_0$  against  $\mathbb{H}_1$ . The method is to design a statistic which is informative about  $\mathbb{H}_1$ . If the observed value of the statistic is consistent with random variation under the assumption that  $\mathbb{H}_0$  is true, then we deduce that there is no evidence against  $\mathbb{H}_0$  and consequently do not reject  $\mathbb{H}_0$ . However, if the statistic takes a value which is unlikely to occur under the assumption that  $\mathbb{H}_0$  is true, then we deduce that there is evidence against  $\mathbb{H}_0$ , and consequently we reject  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$ . The steps are to design a test statistic and characterize its sampling distribution under the assumption that  $\mathbb{H}_0$  is true to control the probability of making a false rejection.

The standard statistic to test  $\mathbb{H}_0$  against  $\mathbb{H}_1$  is the absolute value of the t-statistic

$$|T| = \left| \frac{\hat{\beta} - \beta_0}{s(\hat{\beta})} \right|. \quad (5.18)$$

If  $\mathbb{H}_0$  is true, then we expect  $|T|$  to be small, but if  $\mathbb{H}_1$  is true then we would expect  $|T|$  to be large. Hence the standard rule is to reject  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  for large values of the t-statistic  $|T|$ , and otherwise fail to reject  $\mathbb{H}_0$ . Thus the hypothesis test takes the form

$$\text{Reject } \mathbb{H}_0 \text{ if } |T| > c.$$

The constant  $c$  which appears in the statement of the test is called the **critical value**. Its value is selected to control the probability of false rejections. When the null hypothesis is true,  $|T|$  has an exact student  $t$  distribution (with  $n - k$  degrees of freedom) in the normal regression model. Thus for a given value of  $c$  the probability of false rejection is

$$\begin{aligned} \Pr(\text{Reject } \mathbb{H}_0 \mid \mathbb{H}_0) &= \Pr(|T| > c \mid \mathbb{H}_0) \\ &= \Pr(T > c \mid \mathbb{H}_0) + \Pr(T < -c \mid \mathbb{H}_0) \\ &= 1 - F(c) + F(-c) \\ &= 2(1 - F(c)) \end{aligned}$$

where  $F(u)$  is the  $t_{n-k}$  distribution function. This is the probability of false rejection, and is decreasing in the critical value  $c$ . We select the value  $c$  so that this probability equals a pre-selected value called the **significance level**, which is typically written as  $\alpha$ . It is conventional to set  $\alpha = 0.05$ , though this is not a hard rule. We then select  $c$  so that  $F(c) = 1 - \alpha/2$ , which means that  $c$  is the  $1 - \alpha/2$  quantile (inverse CDF) of the  $t_{n-k}$  distribution, the same as used for confidence intervals. With this choice, the decision rule “Reject  $\mathbb{H}_0$  if  $|T| > c$ ” has a significance level (false rejection probability) of  $\alpha$ .

**Theorem 5.14.1** *In the normal regression model, if the null hypothesis (5.17) is true, then for  $|T|$  defined in (5.18),  $|T| \sim t_{n-k}$ . If  $c$  is set so that  $\Pr(|t_{n-k}| \geq c) = \alpha$ , then the test “Reject  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  if  $|T| > c$ ” has significance level  $\alpha$ .*

To report the result of a hypothesis test we need to pre-determine the significance level  $\alpha$  in order to calculate the critical value  $c$ . This can be inconvenient and arbitrary. A simplification is to report what is known as the **p-value** of the test. In general, when a test takes the form “Reject  $\mathbb{H}_0$  if  $S > c$ ” and  $S$  has null distribution  $G(u)$ , then the p-value of the test is  $p = 1 - G(S)$ . A test with significance level  $\alpha$  can be restated as “Reject  $\mathbb{H}_0$  if  $p < \alpha$ ”. It is sufficient to report the p-value  $p$ , and we can interpret the value of  $p$  as indexing the test’s strength of rejection of the null hypothesis. Thus a p-value of 0.07 might be interpreted as “nearly significant”, 0.05 as “borderline significant”, and 0.001 as “highly significant”. In the context of the normal regression model, the p-value of a t-statistic  $|T|$  is  $p = 2(1 - F_{n-k}(|T|))$  where  $F_{n-k}$  is the CDF of the student  $t$  with  $n - k$  degrees of freedom. For example, in MATLAB the calculation is `2*(1-tcdf(abs(t),n-k))`. In Stata, the default is that for any estimated regression, t-statistics for each estimated coefficient are reported along with their p-values calculated using this same formula. These t-statistics test the hypotheses that each coefficient is zero.

A p-value reports the strength of evidence against  $\mathbb{H}_0$  but is not itself a probability. A common misunderstanding is that the p-value is the “probability that the null hypothesis is true”. This is an incorrect interpretation. It is a statistic, and is random, and is a measure of the evidence against  $\mathbb{H}_0$ , nothing more.

## 5.15 Likelihood Ratio Test

In the previous section we described the t-test as the standard method to test a hypothesis on a single coefficient in a regression. In many contexts, however, we want to simultaneously assess a set of coefficients. In the normal regression model, this can be done by an  $F$  test, which can be derived from the likelihood ratio test.

Partition the regressors as  $\mathbf{x}_i = (\mathbf{x}'_{1i}, \mathbf{x}'_{2i})$  and similarly partition the coefficient vector as  $\boldsymbol{\beta} = (\boldsymbol{\beta}'_1, \boldsymbol{\beta}'_2)'$ . Then the regression model can be written as

$$y_i = \mathbf{x}'_{1i}\boldsymbol{\beta}_1 + \mathbf{x}'_{2i}\boldsymbol{\beta}_2 + e_i. \quad (5.19)$$

Let  $k = \dim(\mathbf{x}_i)$ ,  $k_1 = \dim(\mathbf{x}_{1i})$ , and  $q = \dim(\mathbf{x}_{2i})$ , so that  $k = k_1 + q$ . Partition the variables so that the hypothesis is that the second set of coefficients are zero, or

$$\mathbb{H}_0 : \boldsymbol{\beta}_2 = 0. \quad (5.20)$$

If  $\mathbb{H}_0$  is true, then the regressors  $\mathbf{x}_{2i}$  can be omitted from the regression. In this case we can write (5.19) as

$$y_i = \mathbf{x}'_{1i}\boldsymbol{\beta}_1 + e_i. \quad (5.21)$$

We call (5.21) the null model. The alternative hypothesis is that at least one element of  $\boldsymbol{\beta}_2$  is non-zero and is written as

$$\mathbb{H}_1 : \boldsymbol{\beta}_2 \neq 0.$$

When models are estimated by maximum likelihood, a well-accepted testing procedure is to reject  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  for large values of the Likelihood Ratio – the ratio of the maximized likelihood function under  $\mathbb{H}_1$  and  $\mathbb{H}_0$ , respectively. We now construct this statistic in the normal regression model. Recall from (5.10) that the maximized log-likelihood equals

$$\log L(\hat{\boldsymbol{\beta}}, \hat{\sigma}^2) = -\frac{n}{2} \log(2\pi\hat{\sigma}^2) - \frac{n}{2}.$$

We similarly need to calculate the maximized log-likelihood for the constrained model (5.21). By the same steps for derivation of the unconstrained MLE, we can find that the MLE for (5.21) is OLS of  $y_i$  on  $\mathbf{x}_{1i}$ . We can write this estimator as

$$\tilde{\boldsymbol{\beta}}_1 = (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{y}$$

with residual

$$\tilde{e}_i = y_i - \mathbf{x}'_{1i} \tilde{\boldsymbol{\beta}}_1$$

and error variance estimate

$$\tilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \tilde{e}_i^2.$$

We use the tildes “~” rather than the hats “^” above the constrained estimates to distinguish them from the unconstrained estimates. You can calculate similar to (5.10) that the maximized constrained log-likelihood is

$$\log L(\tilde{\boldsymbol{\beta}}_1, \tilde{\sigma}^2) = -\frac{n}{2} \log(2\pi\tilde{\sigma}^2) - \frac{n}{2}.$$

A classic testing procedure is to reject  $\mathbb{H}_0$  for large values of the ratio of the maximized likelihoods. Equivalently, the test rejects  $\mathbb{H}_0$  for large values of twice the difference in the log-likelihood functions. (Multiplying the likelihood difference by two turns out to be a useful scaling.) This equals

$$\begin{aligned} LR &= 2 \left( \left( -\frac{n}{2} \log(2\pi\hat{\sigma}^2) - \frac{n}{2} \right) - \left( -\frac{n}{2} \log(2\pi\tilde{\sigma}^2) - \frac{n}{2} \right) \right) \\ &= n \log \left( \frac{\tilde{\sigma}^2}{\hat{\sigma}^2} \right). \end{aligned} \quad (5.22)$$

The likelihood ratio test rejects for large values of  $LR$ , or equivalently (see Exercise 5.21), for large values of

$$F = \frac{(\tilde{\sigma}^2 - \hat{\sigma}^2)/q}{\hat{\sigma}^2/(n-k)}. \quad (5.23)$$

This is known as the  $F$  statistic for the test of hypothesis  $\mathbb{H}_0$  against  $\mathbb{H}_1$ .

To develop an appropriate critical value, we need the null distribution of  $F$ . Recall from (3.35) that  $n\hat{\sigma}^2 = \mathbf{e}'\mathbf{M}\mathbf{e}$  where  $\mathbf{M} = \mathbf{I}_n - \mathbf{P}$  with  $\mathbf{P} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ . Similarly, under  $\mathbb{H}_0$ ,  $n\tilde{\sigma}^2 = \mathbf{e}'\mathbf{M}_1\mathbf{e}$  where  $\mathbf{M} = \mathbf{I}_n - \mathbf{P}_1$  with  $\mathbf{P}_1 = \mathbf{X}_1(\mathbf{X}'_1\mathbf{X}_1)^{-1}\mathbf{X}'_1$ . You can calculate that  $\mathbf{M}_1 - \mathbf{M} = \mathbf{P} - \mathbf{P}_1$  is idempotent with rank  $q$ . Furthermore,  $(\mathbf{M}_1 - \mathbf{M})\mathbf{M} = \mathbf{0}$ . It follows that  $\mathbf{e}'(\mathbf{M}_1 - \mathbf{M})\mathbf{e} \sim \chi_q^2$  and is independent of  $\mathbf{e}'\mathbf{M}\mathbf{e}$ . Hence

$$F = \frac{\mathbf{e}'(\mathbf{M}_1 - \mathbf{M})\mathbf{e}/q}{\mathbf{e}'\mathbf{M}\mathbf{e}/(n-k)} \sim \frac{\chi_q^2/q}{\chi_{n-k}^2/(n-k)} \sim F_{q,n-k}$$

an exact  $F$  distribution with degrees of freedom  $q$  and  $n-k$ , respectively. Thus under  $\mathbb{H}_0$ , the  $F$  statistic has an exact  $F$  distribution.

The critical values are selected from the upper tail of the  $F$  distribution. For a given significance level  $\alpha$  (typically  $\alpha = 0.05$ ) we select the critical value  $c$  so that  $\Pr(F_{q,n-k} \geq c) = \alpha$ . (For example, in MATLAB the expression is `finv(1- $\alpha$ ,q,n-k)`.) The test rejects  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  if  $F > c$  and does not reject  $\mathbb{H}_0$  otherwise. The p-value of the test is  $p = 1 - G_{q,n-k}(F)$  where  $G_{q,n-k}(u)$  is the  $F_{q,n-k}$  distribution function. (In MATLAB, the p-value is computed as `1-fcdf(f,q,n-k)`.) It is equivalent to reject  $\mathbb{H}_0$  if  $F > c$  or  $p < \alpha$ .

In Stata, the command to test multiple coefficients takes the form ‘test X1 X1’ where X1 and X2 are the names of the variables whose coefficients are tested. Stata then reports the F statistic for the hypothesis that the coefficients are jointly zero along with the p-value calculated using the  $F$  distribution.

**Theorem 5.15.1** *In the normal regression model, if the null hypothesis (5.20) is true, then for  $F$  defined in (5.23),  $F \sim F_{q,n-k}$ . If  $c$  is set so that  $\Pr(F_{q,n-k} \geq c) = \alpha$ , then the test “Reject  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  if  $F > c$ ” has significance level  $\alpha$ .*



Theorem 5.15.1 justifies the  $F$  test in the normal regression model with critical values taken from the  $F$  distribution.

## 5.16 Likelihood Properties

In this section we present some general properties of the likelihood which hold broadly – not just in normal regression.

Suppose that a random vector  $\mathbf{y}$  has the conditional density  $f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})$  where the function  $f$  is known, and the parameter vector  $\boldsymbol{\theta}$  takes values in a parameter space  $\Theta$ . The log-likelihood function for a random sample  $\{\mathbf{y}_i | \mathbf{x}_i : i = 1, \dots, n\}$  takes the form

$$\log L(\boldsymbol{\theta}) = \sum_{i=1}^n \log f(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}).$$

A key property is that the expected log-likelihood is maximized at the true value of the parameter vector. At this point it is useful to make a notational distinction between a generic parameter value  $\boldsymbol{\theta}$  and its true value  $\boldsymbol{\theta}_0$ . Set  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ .

**Theorem 5.16.1**  $\boldsymbol{\theta}_0 = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \mathbb{E}(\log L(\boldsymbol{\theta}) | \mathbf{X})$

This motivates estimating  $\boldsymbol{\theta}$  by finding the value which maximizes the log-likelihood function. This is the maximum likelihood estimator (MLE):

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \log L(\boldsymbol{\theta}).$$

The **score** of the likelihood function is the vector of partial derivatives with respect to the parameters, evaluated at the true values,

$$\left. \frac{\partial}{\partial \boldsymbol{\theta}} \log L(\boldsymbol{\theta}) \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} = \sum_{i=1}^n \left. \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}) \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}.$$

The covariance matrix of the score is known as the **Fisher information**:

$$\mathcal{I} = \operatorname{var} \left( \left. \frac{\partial}{\partial \boldsymbol{\theta}} \log L(\boldsymbol{\theta}_0) \right| \mathbf{X} \right).$$

Some important properties of the score and information are now presented.

**Theorem 5.16.2** *If  $\log f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})$  is second differentiable and the support of  $\mathbf{y}$  does not depend on  $\boldsymbol{\theta}$  then*

1.  $\mathbb{E} \left( \left. \frac{\partial}{\partial \boldsymbol{\theta}} \log L(\boldsymbol{\theta}) \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} | \mathbf{X} \right) = 0$
2.  $\mathcal{I} = \sum_{i=1}^n \mathbb{E} \left( \left. \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}_0) \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}_0)' \right| \mathbf{x}_i \right)$   
 $= -\mathbb{E} \left( \left. \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \log L(\boldsymbol{\theta}_0) \right| \mathbf{X} \right)$

The first result says that the score is mean zero. The second result shows that the variance of the score equals the negative expectation of the second derivative matrix. This is known as the **Information Matrix Equality**.

We now establish the famous Cramér-Rao Lower Bound.

**Theorem 5.16.3** (*Cramér-Rao*) Under the assumptions of Theorem 5.16.2, if  $\tilde{\boldsymbol{\theta}}$  is an unbiased estimator of  $\boldsymbol{\theta}$ , then  $\text{var}(\tilde{\boldsymbol{\theta}} | \mathbf{X}) \geq \mathcal{I}^{-1}$ .

Theorem 5.16.3 shows that the inverse of the information matrix is a lower bound for the covariance matrix of unbiased estimators. This result is similar to the Gauss-Markov Theorem which established a lower bound for unbiased estimators in homoskedastic linear regression.

### Ronald Fisher

The British statistician Ronald Fisher (1890-1962) is one of the core founders of modern statistical theory. His contributions include the  $F$  distribution, p-values, the concept of Fisher information, and that of sufficient statistics.

## 5.17 Information Bound for Normal Regression

Recall the normal regression log-likelihood which has the parameters  $\boldsymbol{\beta}$  and  $\sigma^2$ . The likelihood scores for this model are

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\beta}} \log L(\boldsymbol{\beta}, \sigma^2) &= \frac{1}{\sigma^2} \sum_{i=1}^n \mathbf{x}_i (y_i - \mathbf{x}_i' \boldsymbol{\beta}) \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n \mathbf{x}_i e_i \end{aligned}$$

and

$$\begin{aligned} \frac{\partial}{\partial \sigma^2} \log L(\boldsymbol{\beta}, \sigma^2) &= -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (y_i - \mathbf{x}_i' \boldsymbol{\beta})^2 \\ &= \frac{1}{2\sigma^4} \sum_{i=1}^n (e_i^2 - \sigma^2). \end{aligned}$$

It follows that the information matrix is

$$\mathcal{I} = \text{var} \left( \begin{array}{c} \frac{\partial}{\partial \boldsymbol{\beta}} \log L(\boldsymbol{\beta}, \sigma^2) \\ \frac{\partial}{\partial \sigma^2} \log L(\boldsymbol{\beta}, \sigma^2) \end{array} \mid \mathbf{X} \right) = \left( \begin{array}{cc} \frac{1}{\sigma^2} \mathbf{X}' \mathbf{X} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{array} \right) \quad (5.24)$$

(see Exercise 5.22). The Cramér-Rao Lower Bound is

$$\mathcal{I}^{-1} = \left( \begin{array}{cc} \sigma^2 (\mathbf{X}' \mathbf{X})^{-1} & 0 \\ 0 & \frac{2\sigma^4}{n} \end{array} \right).$$

This shows that the lower bound for estimation of  $\beta$  is  $\sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$  and the lower bound for  $\sigma^2$  is  $2\sigma^4/n$ .

Since in the homoskedastic linear regression model the OLS estimator is unbiased and has variance  $\sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$ , it follows that OLS is Cramér-Rao efficient in the normal regression model, in the sense that no unbiased estimator has a lower variance matrix. This expands on the Gauss-Markov theorem, which stated that no linear unbiased estimator has a lower variance matrix in the homoskedastic regression model. Notice that the results are complementary. Gauss-Markov efficiency concerns a more narrow class of estimators (linear) but allows a broader model class (linear homoskedastic rather than normal regression). The Cramér-Rao efficiency result is more powerful in that it does not restrict the class of estimators (beyond unbiasedness) but is more restrictive in the class of models allowed (normal regression).

In contrast, the unbiased estimator  $s^2$  of  $\sigma^2$  has variance  $2\sigma^4/(n-k)$  (see Exercise 5.23) which is larger than the Cramér-Rao lower bound  $2\sigma^4/n$ . Thus in contrast to the coefficient estimator, the variance estimator is not Cramér-Rao efficient.

## 5.18 Gamma Function\*

The normal and related distributions make frequent use of the what is known as the **gamma function**. For  $\alpha > 0$  it is defined as

$$\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} \exp(-x) dx. \quad (5.25)$$

While it appears quite simple, it has some advanced properties. One is that  $\Gamma(\alpha)$  does not have a close-form solution (except for special values of  $\alpha$ ). Thus it is typically represented using the symbol  $\Gamma(\alpha)$  and implemented computationally using numerical methods.

Special values include

$$\Gamma(1) = \int_0^\infty \exp(-x) dx = 1 \quad (5.26)$$

and

$$\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}. \quad (5.27)$$

The latter holds by making the change of variables  $x = u^2$  in (5.25) and applying (5.2).

By integration by parts you can show that it satisfies the property

$$\Gamma(1 + \alpha) = \Gamma(\alpha)\alpha.$$

Combined with (5.26) we find that for positive integers  $n$ ,

$$\Gamma(n) = (n-1)!$$

This shows that the gamma function is a continuous version of the factorial.

A useful fact is

$$\int_0^\infty y^{a-1} \exp(-by) dy = b^{-a} \Gamma(a) \quad (5.28)$$

which can be found by applying change-of-variables to the definition (5.25).

Another useful fact is for  $\alpha \in \mathbb{R}$

$$\lim_{n \rightarrow \infty} \frac{\Gamma(n + \alpha)}{\Gamma(n) n^\alpha} = 1. \quad (5.29)$$

## 5.19 Technical Proofs\*

**Proof of Theorem 5.2.1.** Squaring expression (5.2)

$$\begin{aligned}
 \left( \int_0^\infty \exp(-x^2/2) dx \right)^2 &= \int_0^\infty \exp(-x^2/2) dx \int_0^\infty \exp(-u^2/2) du \\
 &= \int_0^\infty \int_0^\infty \exp(-(x^2 + u^2)/2) dx du \\
 &= \int_0^\infty \int_0^{\pi/2} r \exp(-r^2/2) d\theta dr \\
 &= \frac{\pi}{2} \int_0^\infty r \exp(-r^2/2) dr \\
 &= \frac{\pi}{2}.
 \end{aligned}$$

The third equality is the key. It makes the change-of-variables to polar coordinates  $x = r \cos \theta$  and  $u = r \sin \theta$  so that  $x^2 + u^2 = r^2$ . The Jacobian of this transformation is  $r$ . The region of integration in the  $(x, u)$  units is the positive orthant (upper-right region), which corresponds to integrating  $\theta$  from 0 to  $\pi/2$  in polar coordinates. The final two equalities are simple integration. Taking the square root we obtain (5.2). ■

**Proof of Theorem 5.2.3.** Let  $M_x(\mathbf{t}) = \exp(\mathbf{t}'\boldsymbol{\mu} + \frac{1}{2}\mathbf{t}'\boldsymbol{\Sigma}\mathbf{t})$  be the moment generating function of  $\mathbf{X}$  by Theorem 5.2.2. Then the MGF of  $\mathbf{Y}$  is

$$\begin{aligned}
 \mathbb{E}(\exp(\mathbf{s}'\mathbf{Y})) &= \mathbb{E}\exp(\mathbf{s}'(\mathbf{a} + \mathbf{B}\mathbf{X})) \\
 &= \exp(\mathbf{s}'\mathbf{a}) \mathbb{E}\exp(\mathbf{s}'\mathbf{B}\mathbf{X}) \\
 &= \exp(\mathbf{s}'\mathbf{a}) M_x(\mathbf{B}'\mathbf{s}) \\
 &= \exp(\mathbf{s}'\mathbf{a}) \exp\left(\mathbf{s}'\mathbf{B}\boldsymbol{\mu} + \frac{1}{2}\mathbf{s}'\mathbf{B}\boldsymbol{\Sigma}\mathbf{B}'\mathbf{s}\right) \\
 &= \exp\left(\mathbf{s}'(\mathbf{a} + \mathbf{B}\boldsymbol{\mu}) + \frac{1}{2}\mathbf{s}'(\mathbf{B}\boldsymbol{\Sigma}\mathbf{B}')\mathbf{s}\right)
 \end{aligned}$$

which is the MGF of  $N(\mathbf{a} + \mathbf{B}\boldsymbol{\mu}, \mathbf{B}\boldsymbol{\Sigma}\mathbf{B}')$ . Thus  $\mathbf{Y} \sim N(\mathbf{a} + \mathbf{B}\boldsymbol{\mu}, \mathbf{B}\boldsymbol{\Sigma}\mathbf{B}')$  as claimed. ■

**Proof of Theorem 5.2.4.** Let  $k_1$  and  $k_2$  denote the dimensions of  $\mathbf{X}_1$  and  $\mathbf{X}_2$  and set  $k = k_1 + k_2$ . If the components are uncorrelated then the covariance matrix for  $\mathbf{X}$  takes the form

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_1 & 0 \\ 0 & \boldsymbol{\Sigma}_2 \end{bmatrix}$$

In this case the joint density function of  $\mathbf{X}$  equals

$$\begin{aligned}
 f(\mathbf{x}_1, \mathbf{x}_2) &= \frac{1}{(2\pi)^{k/2} (\det(\boldsymbol{\Sigma}_1) \det(\boldsymbol{\Sigma}_2))^{1/2}} \\
 &\quad \cdot \exp\left(-\frac{(\mathbf{x}_1 - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}_1^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1) + (\mathbf{x}_2 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}_2^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2)}{2}\right) \\
 &= \frac{1}{(2\pi)^{k_1/2} (\det(\boldsymbol{\Sigma}_1))^{1/2}} \exp\left(-\frac{(\mathbf{x}_1 - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}_1^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1)}{2}\right) \\
 &\quad \cdot \frac{1}{(2\pi)^{k_2/2} (\det(\boldsymbol{\Sigma}_2))^{1/2}} \exp\left(-\frac{(\mathbf{x}_2 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}_2^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2)}{2}\right).
 \end{aligned}$$

This is the product of two multivariate normal densities in  $\mathbf{x}_1$  and  $\mathbf{x}_2$ . Joint densities factor if (and only if) the components are independent. This shows that uncorrelatedness implies independence.

The converse (that independence implies uncorrelatedness) holds generally. ■

**Proof of Theorem 5.3.1.** We demonstrate that  $Q = \mathbf{X}'\mathbf{X}$  has density function (5.3) by verifying that both have the same moment generating function (MGF). First, the MGF of  $\mathbf{X}'\mathbf{X}$  is

$$\begin{aligned} \mathbb{E}(\exp(t\mathbf{X}'\mathbf{X})) &= \int_{-\infty}^{\infty} \exp(t\mathbf{x}'\mathbf{x}) \frac{1}{(2\pi)^{k/2}} \exp\left(-\frac{\mathbf{x}'\mathbf{x}}{2}\right) d\mathbf{x} \\ &= \int_{-\infty}^{\infty} \frac{1}{(2\pi)^{k/2}} \exp\left(-\frac{\mathbf{x}'\mathbf{x}}{2}(1-2t)\right) d\mathbf{x} \\ &= (1-2t)^{-r/2} \int_{-\infty}^{\infty} \frac{1}{(2\pi)^{k/2}} \exp\left(-\frac{\mathbf{u}'\mathbf{u}}{2}\right) d\mathbf{u} \\ &= (1-2t)^{-r/2}. \end{aligned} \quad (5.30)$$

The fourth equality uses the change of variables  $\mathbf{u} = (1-2t)^{1/2}\mathbf{x}$  and the final equality is the normal probability integral. Second, the MGF of the density (5.3) is

$$\begin{aligned} \int_0^{\infty} \exp(tq) f(q) dq &= \int_0^{\infty} \exp(tq) \frac{1}{\Gamma\left(\frac{r}{2}\right) 2^{r/2}} q^{r/2-1} \exp(-q/2) dq \\ &= \int_0^{\infty} \frac{1}{\Gamma\left(\frac{r}{2}\right) 2^{r/2}} q^{r/2-1} \exp(-q(1/2-t)) dq \\ &= \frac{1}{\Gamma\left(\frac{r}{2}\right) 2^{r/2}} (1/2-t)^{-r/2} \Gamma\left(\frac{r}{2}\right) \\ &= (1-2t)^{-r/2}, \end{aligned} \quad (5.31)$$

the third equality using the gamma integral (5.28). The MGFs (5.30) and (5.31) are equal, verifying that (5.3) is the density of  $Q$  as claimed. ■

**Proof of Theorem 5.3.2.** As in the proof of Theorem 5.3.1 we verify that the MGF of  $Q = \mathbf{X}'\mathbf{X}$  when  $\mathbf{X} \sim N(\boldsymbol{\mu}, \mathbf{I}_r)$  is equal to the MGF of the density function (5.4).

First, we calculate the MGF of  $Q = \mathbf{X}'\mathbf{X}$  when  $\mathbf{X} \sim N(\boldsymbol{\mu}, \mathbf{I}_r)$ . Construct an orthogonal  $r \times r$  matrix  $\mathbf{H} = [\mathbf{h}_1, \mathbf{H}_2]$  whose first column equals  $\mathbf{h}_1 = \boldsymbol{\mu}(\boldsymbol{\mu}'\boldsymbol{\mu})^{-1/2}$ . Note that  $\mathbf{h}_1'\boldsymbol{\mu} = \lambda^{1/2}$  and  $\mathbf{H}_2'\boldsymbol{\mu} = \mathbf{0}$ . Define  $\mathbf{Z} = \mathbf{H}'\mathbf{X} \sim N(\boldsymbol{\mu}^*, \mathbf{I}_r)$  where

$$\boldsymbol{\mu}^* = \mathbf{H}'\boldsymbol{\mu} = \begin{pmatrix} \mathbf{h}_1'\boldsymbol{\mu} \\ \mathbf{H}_2'\boldsymbol{\mu} \end{pmatrix} = \begin{pmatrix} \lambda^{1/2} \\ \mathbf{0} \end{pmatrix} \begin{pmatrix} 1 \\ r-1 \end{pmatrix}.$$

It follows that  $Q = \mathbf{X}'\mathbf{X} = \mathbf{Z}'\mathbf{Z} = Z_1^2 + \mathbf{Z}_2'\mathbf{Z}_2$  where  $Z_1 \sim N(\lambda^{1/2}, 1)$  and  $\mathbf{Z}_2 \sim N(0, \mathbf{I}_{r-1})$  are

independent. Notice that  $\mathbf{Z}'_2 \mathbf{Z}_2 \sim \chi^2_{r-1}$  so has MGF  $(1-2t)^{-(r-1)/2}$  by (5.31). The MGF of  $Z_1^2$  is

$$\begin{aligned} \mathbb{E}(\exp(tZ_1^2)) &= \int_{-\infty}^{\infty} \exp(tx^2) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x-\sqrt{\lambda})^2\right) dx \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x^2(1-2t) - 2x\sqrt{\lambda} + \lambda)\right) dx \\ &= (1-2t)^{-1/2} \exp\left(-\frac{\lambda}{2}\right) \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(u^2 - 2u\sqrt{\frac{\lambda}{1-2t}}\right)\right) du \\ &= (1-2t)^{-1/2} \exp\left(-\frac{\lambda t}{1-2t}\right) \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(u - \sqrt{\frac{\lambda}{1-2t}}\right)^2\right) du \\ &= (1-2t)^{-1/2} \exp\left(-\frac{\lambda t}{1-2t}\right) \end{aligned}$$

where the third equality uses the change of variables  $u = (1-2t)^{1/2}x$ . Thus the MGF of  $Q = Z_1^2 + \mathbf{Z}'_2 \mathbf{Z}_2$  is

$$\begin{aligned} \mathbb{E}(\exp(tQ)) &= \mathbb{E}(\exp(t(Z_1^2 + \mathbf{Z}'_2 \mathbf{Z}_2))) \\ &= \mathbb{E}(\exp(tZ_1^2)) \mathbb{E}(\exp(t\mathbf{Z}'_2 \mathbf{Z}_2)) \\ &= (1-2t)^{-r/2} \exp\left(-\frac{\lambda t}{1-2t}\right). \end{aligned} \tag{5.32}$$

Second, we calculate the MGF of (5.4). It equals

$$\begin{aligned} &\int_0^{\infty} \exp(tx) \sum_{i=0}^{\infty} \frac{e^{-\lambda/2}}{i!} \left(\frac{\lambda}{2}\right)^i f_{r+2i}(x) dx \\ &= \sum_{i=0}^{\infty} \frac{e^{-\lambda/2}}{i!} \left(\frac{\lambda}{2}\right)^i \int_0^{\infty} \exp(tx) f_{r+2i}(x) dx \\ &= \sum_{i=0}^{\infty} \frac{e^{-\lambda/2}}{i!} \left(\frac{\lambda}{2}\right)^i (1-2t)^{-(r+2i)/2} \\ &= e^{-\lambda/2} (1-2t)^{-r/2} \sum_{i=0}^{\infty} \frac{1}{i!} \left(\frac{\lambda}{2(1-2t)}\right)^i \\ &= e^{-\lambda/2} (1-2t)^{-r/2} \exp\left(\frac{\lambda}{2(1-2t)}\right) \\ &= (1-2t)^{-r/2} \exp\left(\frac{\lambda t}{1-2t}\right) \end{aligned} \tag{5.33}$$

where the second equality uses (5.31), and the fourth uses  $\exp(x) = \sum_{i=0}^{\infty} \frac{x^i}{i!}$ . We can see that (5.32) equals (5.33), verifying that (5.4) is the density of  $Q$  as stated. ■

**Proof of Theorem 5.3.3.** The fact that  $\mathbf{A} > 0$  means that we can write  $\mathbf{A} = \mathbf{C}\mathbf{C}'$  where  $\mathbf{C}$  is non-singular (see Section A.9). Then  $\mathbf{A}^{-1} = \mathbf{C}^{-1'}\mathbf{C}^{-1}$  and by Theorem 5.2.3

$$\mathbf{C}^{-1}\mathbf{X} \sim \mathbf{N}(\mathbf{C}^{-1}\boldsymbol{\mu}, \mathbf{C}^{-1}\mathbf{A}\mathbf{C}^{-1'}) = \mathbf{N}(\mathbf{C}^{-1}\boldsymbol{\mu}, \mathbf{C}^{-1}\mathbf{C}\mathbf{C}'\mathbf{C}^{-1'}) = \mathbf{N}(\boldsymbol{\mu}^*, \mathbf{I}_r)$$

where  $\boldsymbol{\mu}^* = \mathbf{C}^{-1}\boldsymbol{\mu}$ . Thus by the definition of the non-central chi-square

$$\mathbf{X}'\mathbf{A}^{-1}\mathbf{X} = \mathbf{X}'\mathbf{C}^{-1'}\mathbf{C}^{-1}\mathbf{X} = (\mathbf{C}^{-1}\mathbf{X})'(\mathbf{C}^{-1}\mathbf{X}) \sim \chi_r^2(\boldsymbol{\mu}'^*\boldsymbol{\mu}^*).$$

Since

$$\boldsymbol{\mu}'^* \boldsymbol{\mu}^* = \boldsymbol{\mu}' \mathbf{C}^{-1'} \mathbf{C}^{-1} \boldsymbol{\mu} = \boldsymbol{\mu}' \mathbf{A}^{-1} \boldsymbol{\mu} = \lambda,$$

this equals  $\chi_r^2(\lambda)$  as claimed. ■

**Proof of Theorem 5.4.1.** Using the simple law of iterated expectations,  $T$  has density

$$\begin{aligned} f(x) &= \frac{d}{dx} \Pr \left( \frac{Z}{\sqrt{Q}/r} \leq x \right) \\ &= \frac{d}{dx} \mathbb{E} \left\{ Z \leq x \sqrt{\frac{Q}{r}} \right\} \\ &= \frac{d}{dx} \mathbb{E} \left[ \Pr \left( Z \leq x \sqrt{\frac{Q}{r}} \mid Q \right) \right] \\ &= \mathbb{E} \frac{d}{dx} \Phi \left( x \sqrt{\frac{Q}{r}} \right) \\ &= \mathbb{E} \left( \phi \left( x \sqrt{\frac{Q}{r}} \right) \sqrt{\frac{Q}{r}} \right) \\ &= \int_0^\infty \left( \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{qx^2}{2r} \right) \right) \sqrt{\frac{q}{r}} \left( \frac{1}{\Gamma(\frac{r}{2}) 2^{r/2}} q^{r/2-1} \exp(-q/2) \right) dq \\ &= \frac{\Gamma(\frac{r+1}{2})}{\sqrt{r\pi} \Gamma(\frac{r}{2})} \left( 1 + \frac{x^2}{r} \right)^{-(\frac{r+1}{2})} \end{aligned}$$

using the gamma integral (5.28). ■

**Proof of Theorem 5.4.2.** Notice that for large  $r$ , by the properties of the logarithm

$$\log \left( \left( 1 + \frac{x^2}{r} \right)^{-(\frac{r+1}{2})} \right) = - \left( \frac{r+1}{2} \right) \log \left( 1 + \frac{x^2}{r} \right) \simeq - \left( \frac{r+1}{2} \right) \frac{x^2}{r} \rightarrow -\frac{x^2}{2},$$

the limit as  $r \rightarrow \infty$ , and thus

$$\lim_{r \rightarrow \infty} \left( 1 + \frac{x^2}{r} \right)^{-(\frac{r+1}{2})} = \exp \left( -\frac{x^2}{2} \right). \quad (5.34)$$

Using a property of the gamma function (5.29)

$$\lim_{n \rightarrow \infty} \frac{\Gamma(n + \alpha)}{\Gamma(n) n^\alpha} = 1$$

with  $n = r/2$  and  $\alpha = 1/2$  we find

$$\lim_{r \rightarrow \infty} \frac{\Gamma(\frac{r+1}{2})}{\sqrt{r\pi} \Gamma(\frac{r}{2})} \left( 1 + \frac{x^2}{r} \right)^{-(\frac{r+1}{2})} = \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{x^2}{2} \right) = \phi(x).$$

■

**Proof of Theorem 5.5.1.** Let  $U \sim \chi_m^2$  and  $V \sim \chi_r^2$  be independent and set  $S = U/V$ . Let  $f_m(u)$  be the  $\chi_m^2$  density. By a similar argument as in the proof of Theorem 5.4.1,  $S$  has the density

function

$$\begin{aligned}
f_S(s) &= \mathbb{E}(f_m(sV) V) \\
&= \int_0^\infty f_m(sv) v f_r(v) dv \\
&= \frac{1}{2^{(m+r)/2} \Gamma\left(\frac{m}{2}\right) \Gamma\left(\frac{r}{2}\right)} \int_0^\infty (sv)^{m/2-1} e^{-sv/2} v^{r/2} e^{-v/2} dv \\
&= \frac{s^{m/2-1}}{2^{(m+r)/2} \Gamma\left(\frac{m}{2}\right) \Gamma\left(\frac{r}{2}\right)} \int_0^\infty v^{(m+r)/2-1} e^{-(s+1)v/2} dv \\
&= \frac{s^{m/2-1}}{\Gamma\left(\frac{m}{2}\right) \Gamma\left(\frac{r}{2}\right) (1+s)^{(m+r)/2}} \int_0^\infty t^{(m+r)/2-1} e^{-t} dt \\
&= \frac{s^{m/2-1} \Gamma\left(\frac{m+r}{2}\right)}{\Gamma\left(\frac{m}{2}\right) \Gamma\left(\frac{r}{2}\right) (1+s)^{(m+r)/2}}.
\end{aligned}$$

The fifth equality make the change-of variables  $v = 2t/(1+s)$ , and the sixth uses the definition of the Gamma function  $\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$ . Making the change-of-variables  $x = sr/m$ , we obtain the density as stated. ■

**Proof of Theorem 5.5.2.** The density of  $mF$  is

$$\frac{x^{m/2-1} \Gamma\left(\frac{m+r}{2}\right)}{r^{m/2} \Gamma\left(\frac{m}{2}\right) \Gamma\left(\frac{r}{2}\right) \left(1 + \frac{x}{r}\right)^{(m+r)/2}} \quad (5.35)$$

Using (5.29) with  $n = r/2$  and  $\alpha = m/2$  we have

$$\lim_{r \rightarrow \infty} \frac{\Gamma\left(\frac{m+r}{2}\right)}{r^{m/2} \Gamma\left(\frac{r}{2}\right)} = 2^{-m/2}$$

and similarly to (5.34) we have

$$\lim_{r \rightarrow \infty} \left(1 + \frac{x}{r}\right)^{-(\frac{m+r}{2})} = \exp\left(-\frac{x}{2}\right).$$

Together, (5.35) tends to

$$\frac{x^{m/2-1} \exp\left(-\frac{x}{2}\right)}{2^{m/2} \Gamma\left(\frac{m}{2}\right)}$$

which is the  $\chi_m^2$  density. ■

**Proof of Theorem 5.16.1.** Since  $\log(u)$  is concave we apply Jensen's inequality (B.5), take expectations are with respect to the true density  $f(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta}_0)$ , and note that the density  $f(\mathbf{y}_i \mid \mathbf{x}_i, \boldsymbol{\theta})$ , integrates to 1 for any  $\boldsymbol{\theta} \in \boldsymbol{\Theta}$ , to find that

$$\begin{aligned}
\mathbb{E}\left(\log \frac{L(\boldsymbol{\theta})}{L(\boldsymbol{\theta}_0)} \mid \mathbf{X}\right) &\leq \log \mathbb{E}\left(\frac{L(\boldsymbol{\theta})}{L(\boldsymbol{\theta}_0)} \mid \mathbf{X}\right) \\
&= \log \int \cdots \int \left( \frac{\prod_{i=1}^n f(\mathbf{y}_i \mid \mathbf{x}_i, \boldsymbol{\theta})}{\prod_{i=1}^n f(\mathbf{y}_i \mid \mathbf{x}_i, \boldsymbol{\theta}_0)} \right) \prod_{i=1}^n f(\mathbf{y}_i \mid \mathbf{x}_i, \boldsymbol{\theta}_0) d\mathbf{y}_1 \cdots d\mathbf{y}_n \\
&= \log \int \cdots \int \prod_{i=1}^n f(\mathbf{y}_i \mid \mathbf{x}_i, \boldsymbol{\theta}) d\mathbf{y}_1 \cdots d\mathbf{y}_n \\
&= \log 1 \\
&= 0.
\end{aligned}$$



This implies for any  $\boldsymbol{\theta} \in \boldsymbol{\Theta}$ ,  $\mathbb{E}(\log L(\boldsymbol{\theta})) \leq \mathbb{E}(\log L(\boldsymbol{\theta}_0))$ . Hence  $\boldsymbol{\theta}_0$  maximizes  $\mathbb{E}(\log L(\boldsymbol{\theta}))$  as claimed. ■

**Proof of Theorem 5.16.2.** For part 1, Since the support of  $\mathbf{y}$  does not depend on  $\boldsymbol{\theta}$  we can exchange integration and differentiation:

$$\mathbb{E} \left( \left. \frac{\partial}{\partial \boldsymbol{\theta}} \log L(\boldsymbol{\theta}) \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \middle| \mathbf{X} \right) = \frac{\partial}{\partial \boldsymbol{\theta}} \mathbb{E} (\log L(\boldsymbol{\theta}) |_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} | \mathbf{X}).$$

Theorem 5.16.1 showed that  $\mathbb{E}(\log L(\boldsymbol{\theta}))$  is maximized at  $\boldsymbol{\theta}_0$ , which has the first-order condition

$$\frac{\partial}{\partial \boldsymbol{\theta}} \mathbb{E} (\log L(\boldsymbol{\theta}) |_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} | \mathbf{X}) = 0$$

as needed.

For part 2, using part 1 and the fact the observations are independent

$$\begin{aligned} \mathcal{I} &= \text{var} \left( \left. \frac{\partial}{\partial \boldsymbol{\theta}} \log L(\boldsymbol{\theta}_0) \right| \mathbf{X} \right) \\ &= \mathbb{E} \left( \left( \left. \frac{\partial}{\partial \boldsymbol{\theta}} \log L(\boldsymbol{\theta}_0) \right| \mathbf{X} \right) \left( \left. \frac{\partial}{\partial \boldsymbol{\theta}} \log L(\boldsymbol{\theta}_0) \right| \mathbf{X} \right)' \middle| \mathbf{X} \right) \\ &= \sum_{i=1}^n \mathbb{E} \left( \left( \left. \frac{\partial}{\partial \boldsymbol{\theta}} f(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}_0) \right| \mathbf{x}_i \right) \left( \left. \frac{\partial}{\partial \boldsymbol{\theta}} f(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}_0) \right| \mathbf{x}_i \right)' \middle| \mathbf{x}_i \right) \end{aligned}$$

which is the first equality.

For the second, observe that

$$\frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) = \frac{\frac{\partial}{\partial \boldsymbol{\theta}} f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})}{f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})}$$

and

$$\begin{aligned} \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \log f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) &= \frac{\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})}{f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})} - \frac{\frac{\partial}{\partial \boldsymbol{\theta}} f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) \frac{\partial}{\partial \boldsymbol{\theta}'} f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})'}{f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})^2} \\ &= \frac{\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})}{f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})} - \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) \frac{\partial}{\partial \boldsymbol{\theta}'} \log f(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})'. \end{aligned}$$

It follows that

$$\begin{aligned} \mathcal{I} &= \sum_{i=1}^n \mathbb{E} \left( \left( \left. \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}_0) \right| \mathbf{x}_i \right) \left( \left. \frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}_0) \right| \mathbf{x}_i \right)' \middle| \mathbf{x}_i \right) \\ &= - \sum_{i=1}^n \mathbb{E} \left( \left. \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} f(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}_0) \right| \mathbf{x}_i \right) + \sum_{i=1}^n \mathbb{E} \left( \left. \frac{\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} f(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}_0)}{f(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}_0)} \right| \mathbf{x}_i \right). \end{aligned}$$

However, by exchanging integration and differentiation we can check that the second term is zero:

$$\begin{aligned} \mathbb{E} \left( \left. \frac{\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} f(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}_0)}{f(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\theta}_0)} \right| \mathbf{x}_i \right) &= \int \left( \left. \frac{\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} f(\mathbf{y} | \mathbf{x}_i, \boldsymbol{\theta}_0)}{f(\mathbf{y} | \mathbf{x}_i, \boldsymbol{\theta}_0)} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right) f(\mathbf{y} | \boldsymbol{\theta}_0) d\mathbf{y} \\ &= \int \left. \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} f(\mathbf{y} | \mathbf{x}_i, \boldsymbol{\theta}_0) \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} d\mathbf{y} \\ &= \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \int f(\mathbf{y} | \mathbf{x}_i, \boldsymbol{\theta}_0) d\mathbf{y} |_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \\ &= \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} 1 \\ &= 0 \end{aligned}$$

This establishes the second inequality. ■

**Proof of Theorem 5.16.3** Let  $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)$  be the sample, let  $f(\mathbf{Y}, \boldsymbol{\theta}) = \prod_{i=1}^n f(\mathbf{y}_i, \boldsymbol{\theta})$  denote the joint density of the sample, and note  $\log L(\boldsymbol{\theta}) = \log f(\mathbf{Y}, \boldsymbol{\theta})$ . Set

$$\mathbf{S} = \frac{\partial}{\partial \boldsymbol{\theta}} \log L(\boldsymbol{\theta}_0)$$

which by Theorem (5.16.2) has mean zero and variance  $\mathcal{I}$  conditional on  $\mathbf{X}$ . Write the estimator  $\tilde{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}(\mathbf{Y})$  as a function of the data. Since  $\tilde{\boldsymbol{\theta}}$  is unbiased, for any  $\boldsymbol{\theta}$ ,

$$\boldsymbol{\theta} = \mathbb{E}(\tilde{\boldsymbol{\theta}} \mid \mathbf{X}) = \int \tilde{\boldsymbol{\theta}}(\mathbf{Y}) f(\mathbf{Y}, \boldsymbol{\theta}) d\mathbf{Y}.$$

Differentiating with respect to  $\boldsymbol{\theta}$

$$\begin{aligned} \mathbf{I}_n &= \int \tilde{\boldsymbol{\theta}}(\mathbf{Y}) \frac{\partial}{\partial \boldsymbol{\theta}'} f(\mathbf{Y}, \boldsymbol{\theta}) d\mathbf{Y} \\ &= \int \tilde{\boldsymbol{\theta}}(\mathbf{Y}) \frac{\partial}{\partial \boldsymbol{\theta}'} \log f(\mathbf{Y}, \boldsymbol{\theta}) f(\mathbf{Y}, \boldsymbol{\theta}) d\mathbf{Y}. \end{aligned}$$

Evaluating at  $\boldsymbol{\theta}_0$  yields

$$\mathbf{I}_n = \mathbb{E}(\tilde{\boldsymbol{\theta}} \mathbf{S}' \mid \mathbf{X}) = \mathbb{E}((\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \mathbf{S}' \mid \mathbf{X}) \quad (5.36)$$

the second equality since  $\mathbb{E}(\mathbf{S} \mid \mathbf{X}) = 0$ .

By the matrix Cauchy-Schwarz inequality (B.11), (5.36), and  $\text{var}(\mathbf{S} \mid \mathbf{X}) = \mathbb{E}(\mathbf{S} \mathbf{S}' \mid \mathbf{X}) = \mathcal{I}$ ,

$$\begin{aligned} \text{var}(\tilde{\boldsymbol{\theta}} \mid \mathbf{X}) &= \mathbb{E}((\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)' \mid \mathbf{X}) \\ &\geq \mathbb{E}((\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \mathbf{S}' \mid \mathbf{X}) (\mathbb{E}(\mathbf{S} \mathbf{S}' \mid \mathbf{X}))^{-1} \mathbb{E}(\mathbf{S}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)' \mid \mathbf{X}) \\ &= (\mathbb{E}(\mathbf{S} \mathbf{S}' \mid \mathbf{X}))^{-1} \\ &= \mathcal{I}^{-1} \end{aligned}$$

as stated. ■

## Exercises

**Exercise 5.1** For the standard normal density  $\phi(x)$ , show that  $\phi'(x) = -x\phi(x)$ .

**Exercise 5.2** Use the result in Exercise 5.1 and integration by parts to show that for  $X \sim N(0, 1)$ ,  $\mathbb{E}X^2 = 1$ .

**Exercise 5.3** Use the results in Exercises 5.1 and 5.2, plus integration by parts, to show that for  $X \sim N(0, 1)$ ,  $\mathbb{E}X^4 = 3$ .

**Exercise 5.4** Show that the moment generating function (mgf) of  $X \sim N(0, 1)$  is  $m(t) = \mathbb{E}(\exp(tX)) = \exp(t^2/2)$ . (For the definition of the mgf see Section 2.31).

**Exercise 5.5** Use the mgf from Exercise 5.4 to verify that for  $X \sim N(0, 1)$ ,  $\mathbb{E}(X^2) = m''(0) = 1$  and  $\mathbb{E}(X^4) = m^{(4)}(0) = 3$ .

**Exercise 5.6** Write the multivariate  $N(\mathbf{0}, \mathbf{I}_k)$  density as the product of  $N(0, 1)$  density functions. That is, show that

$$\frac{1}{(2\pi)^{k/2}} \exp\left(-\frac{\mathbf{x}'\mathbf{x}}{2}\right) = \phi(x_1) \cdots \phi(x_k).$$

**Exercise 5.7** Show that the mgf of  $\mathbf{X} \sim N(\mathbf{0}, \mathbf{I}_k)$  is  $\mathbb{E}(\exp(\mathbf{t}'\mathbf{X})) = \exp(\frac{1}{2}\mathbf{t}'\mathbf{t})$ .  
Hint: Use Exercise 5.4 and the fact that the elements of  $\mathbf{X}$  are independent.

**Exercise 5.8** Show that the mgf of  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is

$$M(\mathbf{t}) = \mathbb{E}(\exp(\mathbf{t}'\mathbf{X})) = \exp\left(\mathbf{t}'\boldsymbol{\mu} + \frac{1}{2}\mathbf{t}'\boldsymbol{\Sigma}\mathbf{t}\right).$$

Hint: Write  $\mathbf{X} = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2}\mathbf{Z}$  where  $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I}_k)$ .

**Exercise 5.9** Show that the characteristic function of  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is

$$C(\mathbf{t}) = \mathbb{E}(\exp(i\mathbf{t}'\mathbf{X})) = \exp\left(i\boldsymbol{\mu}'\mathbf{t} - \frac{1}{2}\mathbf{t}'\boldsymbol{\Sigma}\mathbf{t}\right).$$

For the definition of the characteristic function see Section 2.31

Hint: For  $X \sim N(0, 1)$ , establish  $\mathbb{E}(\exp(itX)) = \exp(-\frac{1}{2}t^2)$  by integration. Then generalize to  $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  using the same steps as in Exercises 5.7 and 5.8.

**Exercise 5.10** Show that if  $Q \sim \chi_r^2$ , then  $\mathbb{E}(Q) = r$  and  $\text{var}(Q) = 2r$ .

Hint: Use the representation  $Q = \sum_{i=1}^n X_i^2$  with  $X_i$  independent  $N(0, 1)$ .

**Exercise 5.11** Show that if  $Q \sim \chi_k^2(\lambda)$ , then  $\mathbb{E}(Q) = k + \lambda$ .

**Exercise 5.12** Suppose  $X_i$  are independent  $N(\mu_i, \sigma_i^2)$ . Find the distribution of the weighted sum  $\sum_{i=1}^n w_i X_i$ .

**Exercise 5.13** Show that if  $\mathbf{e} \sim N(\mathbf{0}, \mathbf{I}_n \sigma^2)$  and  $\mathbf{H}'\mathbf{H} = \mathbf{I}_n$  then  $\mathbf{u} = \mathbf{H}'\mathbf{e} \sim N(\mathbf{0}, \mathbf{I}_n \sigma^2)$ .

**Exercise 5.14** Show that if  $\mathbf{e} \sim N(\mathbf{0}, \boldsymbol{\Sigma})$  and  $\boldsymbol{\Sigma} = \mathbf{A}\mathbf{A}'$  then  $\mathbf{u} = \mathbf{A}^{-1}\mathbf{e} \sim N(\mathbf{0}, \mathbf{I}_n)$ .

**Exercise 5.15** Show that  $\hat{\boldsymbol{\theta}} = \text{argmax}_{\boldsymbol{\theta} \in \Theta} \log L(\boldsymbol{\theta}) = \text{argmax}_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta})$ .

**Exercise 5.16** For the regression in-sample predicted values  $\hat{y}_i$  show that  $\hat{y}_i | \mathbf{X} \sim N(\mathbf{x}'_i \boldsymbol{\beta}, \sigma^2 h_{ii})$  where  $h_{ii}$  are the leverage values (3.25).

**Exercise 5.17** In the normal regression model, show that the leave-one out prediction errors  $\tilde{e}_i$  and the standardized residuals  $\bar{e}_i$  are independent of  $\hat{\boldsymbol{\beta}}$ , conditional on  $\mathbf{X}$ .

Hint: Use (3.46) and (4.26).

**Exercise 5.18** In the normal regression model, show that the robust covariance matrices  $\hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}^W$ ,  $\hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}$ , and  $\bar{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}$  are independent of the OLS estimate  $\hat{\boldsymbol{\beta}}$ , conditional on  $\mathbf{X}$ .

**Exercise 5.19** Let  $F(u)$  be the distribution function of a random variable  $X$  whose density is symmetric about zero. (This includes the standard normal and the student  $t$ .) Show that  $F(-u) = 1 - F(u)$ .

**Exercise 5.20** Let  $C_\beta = [L, U]$  be a  $1 - \alpha$  confidence interval for  $\beta$ , and consider the transformation  $\theta = g(\beta)$  where  $g(\cdot)$  is monotonically increasing. Consider the confidence interval  $C_\theta = [g(L), g(U)]$  for  $\theta$ . Show that  $\Pr(\theta \in C_\theta) = \Pr(\beta \in C_\beta)$ . Use this result to develop a confidence interval for  $\sigma$ .

**Exercise 5.21** Show that the test “Reject  $\mathbb{H}_0$  if  $LR \geq c_1$ ” for  $LR$  defined in (5.22), and the test “Reject  $\mathbb{H}_0$  if  $F \geq c_2$ ” for  $F$  defined in (5.23), yield the same decisions if  $c_2 = (\exp(c_1/n) - 1)(n - k)/q$ . Why does this mean that the two tests are *equivalent*?

**Exercise 5.22** Show (5.24).

**Exercise 5.23** In the normal regression model, let  $s^2$  be the unbiased estimator of the error variance  $\sigma^2$  from (4.30).

(a) Show that  $\text{var}(s^2) = 2\sigma^4/(n - k)$ .

(b) Show that  $\text{var}(s^2)$  is strictly larger than the Cramér-Rao Lower Bound for  $\sigma^2$ .

## Chapter 6

# An Introduction to Large Sample Asymptotics

### 6.1 Introduction

For inference (confidence intervals and hypothesis testing) on unknown parameters we need sampling distributions, either exact or approximate, of estimates and other statistics.

In Chapter 4 we derived the mean and variance of the least-squares estimator in the context of the linear regression model, but this is not a complete description of the sampling distribution and is thus not sufficient for inference. Furthermore, the theory does not apply in the context of the linear projection model, which is more relevant for empirical applications.

In Chapter 5 we derived the exact sampling distribution of the OLS estimator, t-statistics, and F-statistics for the normal regression model, allowing for inference. But these results are narrowly confined to the normal regression model, which requires the unrealistic assumption that the regression error is normally distributed and independent of the regressors. Perhaps we can view these results as some sort of approximation to the sampling distributions without requiring the assumption of normality, but how can we be precise about this?

To illustrate the situation with an example, let  $y_i$  and  $x_i$  be drawn from the joint density

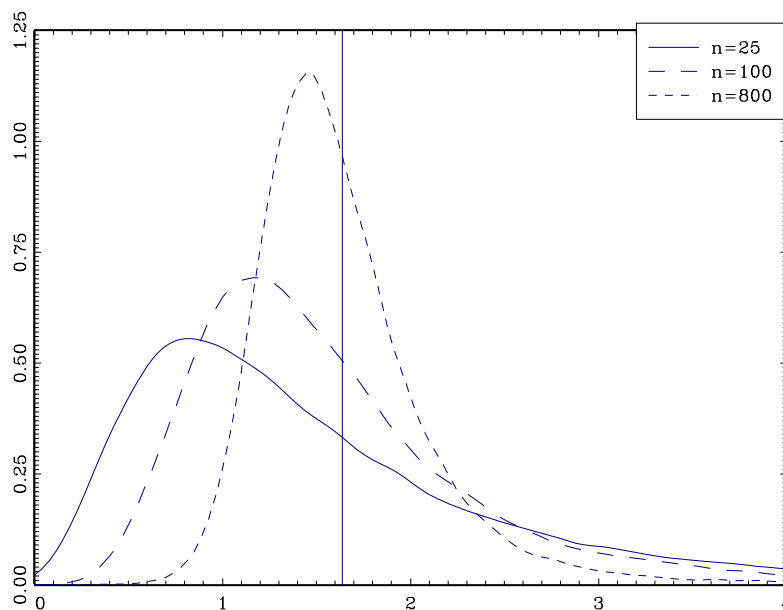
$$f(x, y) = \frac{1}{2\pi xy} \exp\left(-\frac{1}{2}(\log y - \log x)^2\right) \exp\left(-\frac{1}{2}(\log x)^2\right)$$

and let  $\hat{\beta}$  be the slope coefficient estimate from a least-squares regression of  $y_i$  on  $x_i$  and a constant. Using simulation methods, the density function of  $\hat{\beta}$  was computed and plotted in Figure 6.1 for sample sizes of  $n = 25$ ,  $n = 100$  and  $n = 800$ . The vertical line marks the true projection coefficient.

From the figure we can see that the density functions are dispersed and highly non-normal. As the sample size increases the density becomes more concentrated about the population coefficient. Is there a simple way to characterize the sampling distribution of  $\hat{\beta}$ ?

In principle the sampling distribution of  $\hat{\beta}$  is a function of the joint distribution of  $(y_i, x_i)$  and the sample size  $n$ , but in practice this function is extremely complicated so it is not feasible to analytically calculate the exact distribution of  $\hat{\beta}$  except in very special cases. Therefore we typically rely on approximation methods.

In this chapter we introduce asymptotic theory, which approximates by taking the limit of the finite sample distribution as the sample size  $n$  tends to infinity. It is important to understand that this is an approximation technique, as the asymptotic distributions are used to assess the finite sample distributions of our estimators in actual practical samples. The primary tools of asymptotic theory are the weak law of large numbers (WLLN), central limit theorem (CLT), and continuous mapping theorem (CMT). With these tools we can approximate the sampling distributions of most econometric estimators.

Figure 6.1: Sampling Density of  $\hat{\beta}$ 

In this chapter we provide a concise summary. It will be useful for most students to review this material, even if most is familiar.

## 6.2 Asymptotic Limits

“Asymptotic analysis” is a method of approximation obtained by taking a suitable limit. There is more than one method to take limits, but the most common is to take the limit of the sequence of sampling distributions as the sample size tends to positive infinity, written “as  $n \rightarrow \infty$ .” It is not meant to be interpreted literally, but rather as an approximating device.

The first building block for asymptotic analysis is the concept of a limit of a sequence.

**Definition 6.2.1** A sequence  $a_n$  has the **limit**  $a$ , written  $a_n \rightarrow a$  as  $n \rightarrow \infty$ , or alternatively as  $\lim_{n \rightarrow \infty} a_n = a$ , if for all  $\delta > 0$  there is some  $n_\delta < \infty$  such that for all  $n \geq n_\delta$ ,  $|a_n - a| \leq \delta$ .

In words,  $a_n$  has the limit  $a$  if the sequence gets closer and closer to  $a$  as  $n$  gets larger. If a sequence has a limit, that limit is unique (a sequence cannot have two distinct limits). If  $a_n$  has the limit  $a$ , we also say that  $a_n$  **converges** to  $a$  as  $n \rightarrow \infty$ .

Not all sequences have limits. For example, the sequence  $\{1, 2, 1, 2, 1, 2, \dots\}$  does not have a limit. It is therefore sometimes useful to have a more general definition of limits which always exist, and these are the limit superior and limit inferior of a sequence.

**Definition 6.2.2**  $\liminf_{n \rightarrow \infty} a_n \stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} \inf_{m \geq n} a_m$

**Definition 6.2.3**  $\limsup_{n \rightarrow \infty} a_n \stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} \sup_{m \geq n} a_m$

The limit inferior and limit superior always exist (including  $\pm\infty$  as possibilities), and equal when the limit exists. In the example given earlier, the limit inferior of  $\{1, 2, 1, 2, 1, 2, \dots\}$  is 1, and the limit superior is 2.

### 6.3 Convergence in Probability

A sequence of numbers may converge to a limit, but what about a sequence of random variables? For example, consider a sample mean  $\bar{y} = n^{-1} \sum_{i=1}^n y_i$  based on an random sample of  $n$  observations. As  $n$  increases, the distribution of  $\bar{y}$  changes. In what sense can we describe the “limit” of  $\bar{y}$ ? In what sense does it converge?

Since  $\bar{y}$  is a random variable, we cannot directly apply the deterministic concept of a sequence of numbers. Instead, we require a definition of convergence which is appropriate for random variables. There are more than one such definition, but the most commonly used is called convergence in probability.

**Definition 6.3.1** A random variable  $z_n \in \mathbb{R}$  **converges in probability** to  $z$  as  $n \rightarrow \infty$ , denoted  $z_n \xrightarrow{p} z$ , or alternatively  $\text{plim}_{n \rightarrow \infty} z_n = z$ , if for all  $\delta > 0$ ,

$$\lim_{n \rightarrow \infty} \Pr(|z_n - z| \leq \delta) = 1. \quad (6.1)$$

We call  $z$  the **probability limit** (or **plim**) of  $z_n$ .

The definition looks quite abstract, but it formalizes the concept of a sequence of random variables concentrating about a point. The event  $\{|z_n - z| \leq \delta\}$  occurs when  $z_n$  is within  $\delta$  of the point  $z$ .  $\Pr(|z_n - z| \leq \delta)$  is the probability of this event – that  $z_n$  is within  $\delta$  of the point  $z$ . Equation (6.1) states that this probability approaches 1 as the sample size  $n$  increases. The definition of convergence in probability requires that this holds for any  $\delta$ . So for any small interval about  $z$  the distribution of  $z_n$  concentrates within this interval for large  $n$ .

You may notice that the definition concerns the distribution of the random variables  $z_n$ , not their realizations. Furthermore, notice that the definition uses the concept of a conventional (deterministic) limit, but the latter is applied to a sequence of probabilities, not directly to the random variables  $z_n$  or their realizations.

Two comments about the notation are worth mentioning. First, it is conventional to write the convergence symbol as  $\xrightarrow{p}$  where the “ $p$ ” above the arrow indicates that the convergence is “in probability”. You should try and adhere to this notation, and not simply write  $z_n \rightarrow z$ . Second, it is important to include the phrase “as  $n \rightarrow \infty$ ” to be specific about how the limit is obtained.

A common mistake is to confuse convergence in probability with convergence in expectation:

$$\mathbb{E}(z_n) \rightarrow \mathbb{E}(z). \quad (6.2)$$

They are related but distinct concepts. Neither (6.1) nor (6.2) implies the other.

To see the distinction it might be helpful to think through a stylized example. Consider a discrete random variable  $z_n$  which takes the value 0 with probability  $1 - n^{-1}$  and the value  $a_n \neq 0$  with probability  $n^{-1}$ , or

$$\begin{aligned} \Pr(z_n = 0) &= 1 - \frac{1}{n} \\ \Pr(z_n = a_n) &= \frac{1}{n}. \end{aligned} \quad (6.3)$$

In this example the probability distribution of  $z_n$  concentrates at zero as  $n$  increases, regardless of the sequence  $a_n$ . You can check that  $z_n \xrightarrow{p} 0$  as  $n \rightarrow \infty$ .

In this example we can also calculate that the expectation of  $z_n$  is

$$\mathbb{E}(z_n) = \frac{a_n}{n}.$$

Despite the fact that  $z_n$  converges in probability to zero, its expectation will not decrease to zero unless  $a_n/n \rightarrow 0$ . If  $a_n$  diverges to infinity at a rate equal to  $n$  (or faster) then  $\mathbb{E}(z_n)$  will not converge to zero. For example, if  $a_n = n$ , then  $\mathbb{E}(z_n) = 1$  for all  $n$ , even though  $z_n \xrightarrow{p} 0$ . This example might seem a bit artificial, but the point is that the concepts of convergence in probability and convergence in expectation are distinct, so it is important not to confuse one with the other.

Another common source of confusion with the notation surrounding probability limits is that the expression to the right of the arrow “ $\xrightarrow{p}$ ” must be free of dependence on the sample size  $n$ . Thus expressions of the form “ $z_n \xrightarrow{p} c_n$ ” are notationally meaningless and should not be used.

## 6.4 Weak Law of Large Numbers

In large samples we expect parameter estimates to be close to the population values. For example, in Section 4.3 we saw that the sample mean  $\bar{y}$  is unbiased for  $\mu = \mathbb{E}(y)$  and has variance  $\sigma^2/n$ . As  $n$  gets large its variance decreases and thus the distribution of  $\bar{y}$  concentrates about the population mean  $\mu$ . It turns out that this implies that the sample mean converges in probability to the population mean.

When  $y$  has a finite variance there is a fairly straightforward proof by applying Chebyshev's inequality.

**Theorem 6.4.1 Chebyshev's Inequality.** *For any random variable  $z_n$  and constant  $\delta > 0$*

$$\Pr(|z_n - \mathbb{E}z_n| \geq \delta) \leq \frac{\text{var}(z_n)}{\delta^2}.$$

Chebyshev's inequality is terrifically important in asymptotic theory. While its proof is a technical exercise in probability theory, it is quite simple so we discuss it forthwith. Let  $F_n(u)$  denote the distribution of  $z_n - \mathbb{E}z_n$ . Then

$$\Pr(|z_n - \mathbb{E}z_n| \geq \delta) = \Pr((z_n - \mathbb{E}z_n)^2 \geq \delta^2) = \int_{\{u^2 \geq \delta^2\}} dF_n(u).$$

The integral is over the event  $\{u^2 \geq \delta^2\}$ , so that the inequality  $1 \leq \frac{u^2}{\delta^2}$  holds throughout. Thus

$$\int_{\{u^2 \geq \delta^2\}} dF_n(u) \leq \int_{\{u^2 \geq \delta^2\}} \frac{u^2}{\delta^2} dF_n(u) \leq \int \frac{u^2}{\delta^2} dF_n(u) = \frac{\mathbb{E}(z_n - \mathbb{E}z_n)^2}{\delta^2} = \frac{\text{var}(z_n)}{\delta^2},$$

which establishes the desired inequality.

Applied to the sample mean  $\bar{y}$  which has variance  $\sigma^2/n$ , Chebyshev's inequality shows that for any  $\delta > 0$

$$\Pr(|\bar{y} - \mathbb{E}(\bar{y})| \geq \delta) \leq \frac{\sigma^2/n}{\delta^2}.$$



For fixed  $\sigma^2$  and  $\delta$ , the bound on the right-hand-side shrinks to zero as  $n \rightarrow \infty$ . (Specifically, for any  $\varepsilon > 0$  set  $n \geq \sigma^2 / (\delta^2 \varepsilon)$ . Then the right-hand-side is less than or equal to  $\varepsilon$ .) Thus the probability that  $\bar{y}$  is within  $\delta$  of  $\mathbb{E}(\bar{y}) = \mu$  approaches 1 as  $n$  gets large, or

$$\lim_{n \rightarrow \infty} \Pr(|\bar{y} - \mu| < \delta) = 1.$$

This means that  $\bar{y}$  converges in probability to  $\mu$  as  $n \rightarrow \infty$ .

This result is called the **weak law of large numbers**. Our derivation assumed that  $y$  has a finite variance, but with a more careful derivation all that is necessary is a finite mean.

**Theorem 6.4.2 Weak Law of Large Numbers (WLLN)**

If  $y_i$  are independent and identically distributed and  $\mathbb{E}|y| < \infty$ , then as  $n \rightarrow \infty$ ,

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \xrightarrow{p} \mathbb{E}(y).$$

The proof of Theorem 6.4.2 is presented in Section 6.16.

The WLLN shows that the estimator  $\bar{y}$  converges in probability to the true population mean  $\mu$ . In general, an estimator which converges in probability to the population value is called **consistent**.

**Definition 6.4.1** An estimator  $\hat{\theta}$  of a parameter  $\theta$  is **consistent** if  $\hat{\theta} \xrightarrow{p} \theta$  as  $n \rightarrow \infty$ .

**Theorem 6.4.3** If  $y_i$  are independent and identically distributed and  $\mathbb{E}|y| < \infty$ , then  $\hat{\mu} = \bar{y}$  is consistent for the population mean  $\mu$ .

Consistency is a good property for an estimator to possess. It means that for any given data distribution, there is a sample size  $n$  sufficiently large such that the estimator  $\hat{\theta}$  will be arbitrarily close to the true value  $\theta$  with high probability. The theorem does not tell us, however, how large this  $n$  has to be. Thus the theorem does not give practical guidance for empirical practice. Still, it is a minimal property for an estimator to be considered a “good” estimator, and provides a foundation for more useful approximations.

## 6.5 Almost Sure Convergence and the Strong Law\*

Convergence in probability is sometimes called **weak convergence**. A related concept is **almost sure convergence**, also known as **strong convergence**. (In probability theory the term “almost sure” means “with probability equal to one”. An event which is random but occurs with probability equal to one is said to be **almost sure**.)

**Definition 6.5.1** A random variable  $z_n \in \mathbb{R}$  *converges almost surely* to  $z$  as  $n \rightarrow \infty$ , denoted  $z_n \xrightarrow{a.s.} z$ , if for every  $\delta > 0$

$$\Pr \left( \lim_{n \rightarrow \infty} |z_n - z| \leq \delta \right) = 1. \quad (6.4)$$

The convergence (6.4) is stronger than (6.1) because it computes the probability of a limit rather than the limit of a probability. Almost sure convergence is stronger than convergence in probability in the sense that  $z_n \xrightarrow{a.s.} z$  implies  $z_n \xrightarrow{p} z$ .

In the example (6.3) of Section 6.3, the sequence  $z_n$  converges in probability to zero for any sequence  $a_n$ , but this is not sufficient for  $z_n$  to converge almost surely. In order for  $z_n$  to converge to zero almost surely, it is necessary that  $a_n \rightarrow 0$ .

In the random sampling context the sample mean can be shown to converge almost surely to the population mean. This is called the **strong law of large numbers**.

**Theorem 6.5.1 Strong Law of Large Numbers (SLLN)**

If  $y_i$  are independent and identically distributed and  $\mathbb{E}|y| < \infty$ , then as  $n \rightarrow \infty$ ,

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \xrightarrow{a.s.} \mathbb{E}(y).$$

The proof of the SLLN is technically quite advanced so is not presented here. For a proof see Billingsley (1995, Theorem 22.1) or Ash (1972, Theorem 7.2.5).

The WLLN is sufficient for most purposes in econometrics, so we will not use the SLLN in this text.

## 6.6 Vector-Valued Moments

Our preceding discussion focused on the case where  $y$  is real-valued (a scalar), but nothing important changes if we generalize to the case where  $\mathbf{y} \in \mathbb{R}^m$  is a vector. To fix notation, the elements of  $\mathbf{y}$  are

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}.$$

The population mean of  $\mathbf{y}$  is just the vector of marginal means

$$\boldsymbol{\mu} = \mathbb{E}(\mathbf{y}) = \begin{pmatrix} \mathbb{E}(y_1) \\ \mathbb{E}(y_2) \\ \vdots \\ \mathbb{E}(y_m) \end{pmatrix}.$$

When working with random vectors  $\mathbf{y}$  it is convenient to measure their magnitude by their Euclidean length or Euclidean norm

$$\|\mathbf{y}\| = (y_1^2 + \cdots + y_m^2)^{1/2}.$$

In vector notation we have

$$\|\mathbf{y}\|^2 = \mathbf{y}'\mathbf{y}.$$

It turns out that it is equivalent to describe finiteness of moments in terms of the Euclidean norm of a vector or all individual components.

**Theorem 6.6.1** For  $\mathbf{y} \in \mathbb{R}^m$ ,  $\mathbb{E}\|\mathbf{y}\| < \infty$  if and only if  $\mathbb{E}|y_j| < \infty$  for  $j = 1, \dots, m$ .

The  $m \times m$  variance matrix of  $\mathbf{y}$  is

$$\mathbf{V} = \text{var}(\mathbf{y}) = \mathbb{E}((\mathbf{y} - \boldsymbol{\mu})(\mathbf{y} - \boldsymbol{\mu})').$$

$\mathbf{V}$  is often called a variance-covariance matrix. You can show that the elements of  $\mathbf{V}$  are finite if  $\mathbb{E}\|\mathbf{y}\|^2 < \infty$ .

A random sample  $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$  consists of  $n$  observations of independent and identically distributed draws from the distribution of  $\mathbf{y}$ . (Each draw is an  $m$ -vector.) The vector sample mean

$$\bar{\mathbf{y}} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i = \begin{pmatrix} \bar{y}_1 \\ \bar{y}_2 \\ \vdots \\ \bar{y}_m \end{pmatrix}$$

is the vector of sample means of the individual variables.

Convergence in probability of a vector can be defined as convergence in probability of all elements in the vector. Thus  $\bar{\mathbf{y}} \xrightarrow{p} \boldsymbol{\mu}$  if and only if  $\bar{y}_j \xrightarrow{p} \mu_j$  for  $j = 1, \dots, m$ . Since the latter holds if  $\mathbb{E}|y_j| < \infty$  for  $j = 1, \dots, m$ , or equivalently  $\mathbb{E}\|\mathbf{y}\| < \infty$ , we can state this formally as follows.

**Theorem 6.6.2 WLLN for random vectors**

If  $\mathbf{y}_i$  are independent and identically distributed and  $\mathbb{E}\|\mathbf{y}\| < \infty$ , then as  $n \rightarrow \infty$ ,

$$\bar{\mathbf{y}} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i \xrightarrow{p} \mathbb{E}(\mathbf{y}).$$

## 6.7 Convergence in Distribution

The WLLN is a useful first step, but does not give an approximation to the distribution of an estimator. A **large-sample** or **asymptotic** approximation can be obtained using the concept of **convergence in distribution**.

We say that a sequence of random vectors  $\mathbf{z}_n$  converges in distribution if the sequence of distribution functions  $F_n(\mathbf{u}) = \Pr(\mathbf{z}_n \leq \mathbf{u})$  converges to a limit distribution function.

**Definition 6.7.1** Let  $\mathbf{z}_n$  be a random vector with distribution  $F_n(\mathbf{u}) = \Pr(\mathbf{z}_n \leq \mathbf{u})$ . We say that  $\mathbf{z}_n$  **converges in distribution** to  $\mathbf{z}$  as  $n \rightarrow \infty$ , denoted  $\mathbf{z}_n \xrightarrow{d} \mathbf{z}$ , if for all  $\mathbf{u}$  at which  $F(\mathbf{u}) = \Pr(\mathbf{z} \leq \mathbf{u})$  is continuous,  $F_n(\mathbf{u}) \rightarrow F(\mathbf{u})$  as  $n \rightarrow \infty$ .

Under these conditions, it is also said that  $F_n$  **converges weakly** to  $F$ . It is common to refer to  $\mathbf{z}$  and its distribution  $F(\mathbf{u})$  as the **asymptotic distribution**, **large sample distribution**, or **limit distribution** of  $\mathbf{z}_n$ .

When the limit distribution  $\mathbf{z}$  is degenerate (that is,  $\Pr(\mathbf{z} = \mathbf{c}) = 1$  for some  $\mathbf{c}$ ) we can write the convergence as  $\mathbf{z}_n \xrightarrow{d} \mathbf{c}$ , which is equivalent to convergence in probability,  $\mathbf{z}_n \xrightarrow{p} \mathbf{c}$ .

Technically, in most cases of interest it is difficult to establish the limit distributions of sample statistics  $\mathbf{z}_n$  by working directly with their distribution function. It turns out that in most cases it is easier to work with their characteristic function  $C_n(\boldsymbol{\lambda}) = \mathbb{E}(\exp(i\boldsymbol{\lambda}'\mathbf{z}_n))$ , which is a transformation of the distribution. (See Section 2.31 for the definition.) While this is more technical than needed for most applied economists, we introduce this material to give a complete reference for large sample approximations.

The characteristic function  $C_n(\mathbf{t})$  completely describes the distribution of  $\mathbf{z}_n$ . It therefore seems reasonable to expect that if  $C_n(\mathbf{t})$  converges to a limit function  $C(\mathbf{t})$ , then the the distribution of  $\mathbf{z}_n$  converges as well. This turns out to be true, and is known as Lévy's continuity theorem.

**Theorem 6.7.1 Lévy's Continuity Theorem.**  $\mathbf{z}_n \xrightarrow{d} \mathbf{z}$  if and only if  $\mathbb{E}(\exp(i\mathbf{t}'\mathbf{z}_n)) \rightarrow \mathbb{E}(\exp(i\mathbf{t}'\mathbf{z}))$  for every  $\mathbf{t} \in \mathbb{R}^k$ .

While this result seems quite intuitive, a rigorous proof is quite advanced and so is not presented here. See Van der Vaart (2008) Theorem 2.13.

Finally, we mention a standard trick which is commonly used to establish multivariate convergence results.

**Theorem 6.7.2 Cramér-Wold Device.**  $\mathbf{z}_n \xrightarrow{d} \mathbf{z}$  if and only if  $\boldsymbol{\lambda}'\mathbf{z}_n \xrightarrow{d} \boldsymbol{\lambda}'\mathbf{z}$  for every  $\boldsymbol{\lambda} \in \mathbb{R}^k$  with  $\boldsymbol{\lambda}'\boldsymbol{\lambda} = 1$ .

We present a proof in Section 6.16 which is a simple application of Lévy's continuity theorem.

## 6.8 Central Limit Theorem

We would like to obtain a distributional approximation to the sample mean  $\bar{y}$ . We start under the random sampling assumption so that the observations are independent and identically distributed, and have a finite mean  $\mu = \mathbb{E}(y)$  and variance  $\sigma^2 = \text{var}(y)$ .

Let's start by finding the asymptotic distribution of  $\bar{y}$ , in the sense that  $\bar{y} \xrightarrow{d} \mathbf{z}$  for some random variable  $\mathbf{z}$ . From the WLLN we know that  $\bar{y} \xrightarrow{p} \mu$ . Since convergence in probability to a constant is the same as convergence in distribution, this means that  $\bar{y} \xrightarrow{d} \mu$  as well. This is not a useful distributional result as the limit distribution is a constant. To obtain a non-degenerate distribution

we need to rescale  $\bar{y}$ . Recall that  $\text{var}(\bar{y} - \mu) = \sigma^2/n$ , which means that  $\text{var}(\sqrt{n}(\bar{y} - \mu)) = \sigma^2$ . This suggests renormalizing the statistic as

$$z_n = \sqrt{n}(\bar{y} - \mu).$$

Notice that  $\mathbb{E}(z_n) = 0$  and  $\text{var}(z_n) = \sigma^2$ . This shows that the mean and variance have been stabilized. We now seek to determine the asymptotic distribution of  $z_n$ .

The answer is provided by the central limit theorem (CLT) which states that standardized sample averages converge in distribution to normal random vectors. There are several versions of the CLT. The most basic is the case where the observations are independent and identically distributed.

**Theorem 6.8.1 Lindeberg–Lévy Central Limit Theorem.** *If  $y_i$  are independent and identically distributed and  $\mathbb{E}(y_i^2) < \infty$ , then as  $n \rightarrow \infty$*

$$\sqrt{n}(\bar{y} - \mu) \xrightarrow{d} N(0, \sigma^2)$$

*where  $\mu = \mathbb{E}(y)$  and  $\sigma^2 = \mathbb{E}(y_i - \mu)^2$ .*

The proof of the CLT is rather technical (so is presented in Section 6.16) but at the core is a quadratic approximation of the log of the characteristic function.

As we discussed above, in finite samples the standardized sum  $z_n = \sqrt{n}(\bar{y}_n - \mu)$  has mean zero and variance  $\sigma^2$ . What the CLT adds is that  $z_n$  is also approximately normally distributed, and that the normal approximation improves as  $n$  increases.

The CLT is one of the most powerful and mysterious results in statistical theory. It shows that the simple process of averaging induces normality. The first version of the CLT (for the number of heads resulting from many tosses of a fair coin) was established by the French mathematician Abraham de Moivre in an article published in 1733. This was extended to cover an approximation to the binomial distribution in 1812 by Pierre-Simon Laplace in his book *Théorie Analytique des Probabilités*, and the most general statements are credited to articles by the Russian mathematician Aleksandr Lyapunov (1901) and the Finnish mathematician Jarl Waldemar Lindeberg (1920, 1922). The above statement is known as the classic (or Lindeberg–Lévy) CLT due to contributions by Lindeberg (1920) and the French mathematician Paul Pierre Lévy.

A more general version which allows heterogeneous distributions was provided by Lindeberg (1922). The following is the most general statement.

**Theorem 6.8.2 Lindeberg–Feller Central Limit Theorem.** *Suppose  $y_{ni}$  are independent but not necessarily identically distributed with finite means  $\mu_{ni} = \mathbb{E}(y_{ni})$  and variances  $\sigma_{ni}^2 = \mathbb{E}(y_{ni} - \mu_{ni})^2$ . Set  $\bar{\sigma}_n^2 = n^{-1} \sum_{i=1}^n \sigma_{ni}^2$ . If  $\bar{\sigma}_n^2 > 0$  and for all  $\varepsilon > 0$*

$$\lim_{n \rightarrow \infty} \frac{1}{n\bar{\sigma}_n^2} \sum_{i=1}^n \mathbb{E} \left( (y_{ni} - \mu_{ni})^2 1 \left( (y_{ni} - \mu_{ni})^2 \geq \varepsilon n \bar{\sigma}_n^2 \right) \right) = 0 \quad (6.5)$$

*then as  $n \rightarrow \infty$*

$$\frac{\sqrt{n}(\bar{y} - \mathbb{E}(\bar{y}))}{\bar{\sigma}_n^{1/2}} \xrightarrow{d} N(0, 1).$$

The proof of the Lindeberg-Feller CLT is substantially more technical, so we do not present it here. See Billingsley (1995, Theorem 27.2).

The Lindeberg-Feller CLT is quite general as it puts minimal conditions on the sequence of means and variances. The key assumption is equation (6.5) which is known as **Lindeberg's Condition**. In its raw form it is difficult to interpret. The intuition for (6.5) is that it excludes any single observation from dominating the asymptotic distribution. Since (6.5) is quite abstract, in most contexts we use more elementary conditions which are simpler to interpret.

One such alternative is called **Lyapunov's condition**: For some  $\delta > 0$

$$\lim_{n \rightarrow \infty} \frac{1}{n^{1+\delta/2} \bar{\sigma}_n^{2+\delta}} \sum_{i=1}^n \mathbb{E} \left( |y_{ni} - \mu_{ni}|^{2+\delta} \right) = 0. \quad (6.6)$$

Lyapunov's condition implies Lindeberg's condition, and hence the CLT. Indeed, the left-side of (6.5) is bounded by

$$\begin{aligned} & \lim_{n \rightarrow \infty} \frac{1}{n \bar{\sigma}_n^2} \sum_{i=1}^n \mathbb{E} \left( \frac{|y_{ni} - \mu_{ni}|^{2+\delta}}{|y_{ni} - \mu_{ni}|^\delta} 1 \left( |y_{ni} - \mu_{ni}|^2 \geq \varepsilon n \bar{\sigma}_n^2 \right) \right) \\ & \leq \lim_{n \rightarrow \infty} \frac{1}{\varepsilon^{\delta/2} n^{1+\delta/2} \bar{\sigma}_n^{2+\delta}} \sum_{i=1}^n \mathbb{E} \left( |y_{ni} - \mu_{ni}|^{2+\delta} \right) \\ & = 0 \end{aligned}$$

by (6.6).

Lyapunov's condition is still awkward to interpret. A still simpler condition is a uniform moment bound: For some  $\delta > 0$

$$\sup_{n,i} \mathbb{E} |y_{ni}|^{2+\delta} < \infty. \quad (6.7)$$

This is typically combined with the lower variance bound

$$\liminf_{n \rightarrow \infty} \bar{\sigma}_n^2 > 0. \quad (6.8)$$

These bounds together imply Lyapunov's condition. To see this, (6.7) and (6.8) imply there is some  $C < \infty$  such that  $\sup_{n,i} \mathbb{E} |y_{ni}|^{2+\delta} \leq C$  and  $\liminf_{n \rightarrow \infty} \bar{\sigma}_n^2 \geq C^{-1}$ . Without loss of generality assume  $\mu_{ni} = 0$ . Then the left side of (6.6) is bounded by

$$\lim_{n \rightarrow \infty} \frac{C^{2+\delta/2}}{n^{\delta/2}} = 0,$$

so Lyapunov's condition holds and hence the CLT.

An alternative to (6.8) is to assume that the average variance  $\bar{\sigma}_n^2$  converges to a constant, that is,

$$\bar{\sigma}_n^2 = n^{-1} \sum_{i=1}^n \sigma_{ni}^2 \rightarrow \sigma^2 < \infty. \quad (6.9)$$

This assumption is reasonable in many applications.

We now state the simplest and most commonly used version of a heterogeneous CLT based on the Lindeberg-Feller Theorem.

**Theorem 6.8.3** *Suppose  $y_{ni}$  are independent but not necessarily identically distributed. If (6.7) and (6.9) hold, then as  $n \rightarrow \infty$*

$$\sqrt{n} (\bar{y} - \mathbb{E}(\bar{y})) \xrightarrow{d} N(0, \sigma^2). \quad (6.10)$$

One advantage of Theorem 6.8.3 is that it allows  $\sigma^2 = 0$  (unlike Theorem 6.8.2).

## 6.9 Multivariate Central Limit Theorem

Multivariate central limit theory applies when we consider vector-valued observations  $\mathbf{y}_i$  and sample averages  $\bar{\mathbf{y}}$ . In the i.i.d. case we know that the mean of  $\bar{\mathbf{y}}$  is the mean vector  $\boldsymbol{\mu} = \mathbb{E}(\mathbf{y})$  and its variance is  $n^{-1}\mathbf{V}$  where  $\mathbf{V} = \mathbb{E}((\mathbf{y} - \boldsymbol{\mu})(\mathbf{y} - \boldsymbol{\mu})')$ . Again we wish to transform  $\bar{\mathbf{y}}$  so that its mean and variance do not depend on  $n$ . We do this again by centering and scaling, by setting  $\mathbf{z}_n = \sqrt{n}(\bar{\mathbf{y}}_n - \boldsymbol{\mu})$ . This has mean  $\mathbf{0}$  and variance  $\mathbf{V}$ , which are independent of  $n$  as desired.

To develop a distributional approximation for  $\mathbf{z}_n$  we use a multivariate central limit theorem. We present three such results, corresponding to the three univariate results from the previous section. Each is derived from the univariate theory by the Cramér-Wold device (Theorem 6.7.2).

We first present the multivariate version of Theorem 6.8.1.

**Theorem 6.9.1 Multivariate Lindeberg–Lévy Central Limit Theorem.** *If  $\mathbf{y}_i \in \mathbb{R}^k$  are independent and identically distributed and  $\mathbb{E}\|\mathbf{y}_i\|^2 < \infty$ , then as  $n \rightarrow \infty$*

$$\sqrt{n}(\bar{\mathbf{y}} - \boldsymbol{\mu}) \xrightarrow{d} N(\mathbf{0}, \mathbf{V})$$

*where  $\boldsymbol{\mu} = \mathbb{E}(\mathbf{y})$  and  $\mathbf{V} = \mathbb{E}((\mathbf{y} - \boldsymbol{\mu})(\mathbf{y} - \boldsymbol{\mu})')$ .*

We next present a multivariate version of Theorem 6.8.2.

**Theorem 6.9.2 Multivariate Lindeberg-Feller CLT.** *Suppose  $\mathbf{y}_{ni} \in \mathbb{R}^k$  are independent but not necessarily identically distributed with finite means  $\boldsymbol{\mu}_{ni} = \mathbb{E}(\mathbf{y}_{ni})$  and variance matrices  $\mathbf{V}_{ni} = \mathbb{E}((\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni})(\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni})')$ . Set  $\bar{\mathbf{V}}_n = n^{-1} \sum_{i=1}^n \mathbf{V}_{ni}$  and  $\nu_n^2 = \lambda_{\min}(\bar{\mathbf{V}}_n)$ . If  $\nu_n^2 > 0$  and for all  $\varepsilon > 0$*

$$\lim_{n \rightarrow \infty} \frac{1}{n\nu_n^2} \sum_{i=1}^n \mathbb{E} \left( \|\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni}\|^2 1 \left( \|\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni}\|^2 \geq \varepsilon n \nu_n^2 \right) \right) = 0 \quad (6.11)$$

*then as  $n \rightarrow \infty$*

$$\bar{\mathbf{V}}_n^{-1/2} \sqrt{n}(\bar{\mathbf{y}} - \mathbb{E}(\bar{\mathbf{y}})) \xrightarrow{d} N(\mathbf{0}, \mathbf{I}_k).$$

We finally present a multivariate version of Theorem 6.8.3.

**Theorem 6.9.3** *Suppose  $\mathbf{y}_{ni} \in \mathbb{R}^k$  are independent but not necessarily identically distributed with finite means  $\boldsymbol{\mu}_{ni} = \mathbb{E}(\mathbf{y}_{ni})$  and variance matrices  $\mathbf{V}_{ni} = \mathbb{E}((\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni})(\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni})')$ . Set  $\bar{\mathbf{V}}_n = n^{-1} \sum_{i=1}^n \mathbf{V}_{ni}$ . If*

$$\bar{\mathbf{V}}_n \rightarrow \mathbf{V} > 0 \quad (6.12)$$

*and for some  $\delta > 0$*

$$\sup_{n,i} \mathbb{E} \|\mathbf{y}_{ni}\|^{2+\delta} < \infty \quad (6.13)$$

*then as  $n \rightarrow \infty$*

$$\sqrt{n}(\bar{\mathbf{y}} - \mathbb{E}(\bar{\mathbf{y}})) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}).$$

Similarly to Theorem 6.8.3, an advantage of Theorem 6.9.3 is that it allows the variance matrix  $\mathbf{V}$  to be singular.

## 6.10 Higher Moments

Often we want to estimate a parameter  $\boldsymbol{\mu}$  which is the expected value of a transformation of a random vector  $\mathbf{y}$ . That is,  $\boldsymbol{\mu}$  can be written as

$$\boldsymbol{\mu} = \mathbb{E}(\mathbf{h}(\mathbf{y}))$$

for some function  $\mathbf{h} : \mathbb{R}^m \rightarrow \mathbb{R}^k$ . For example, the second moment of  $y$  is  $\mathbb{E}(y^2)$ , the  $r^{th}$  is  $\mathbb{E}(y^r)$ , the moment generating function is  $\mathbb{E}(\exp(ty))$ , and the distribution function is  $\mathbb{E}(1\{y \leq x\})$ .

Estimating parameters of this form fits into our previous analysis by defining the random variable  $\mathbf{z} = \mathbf{h}(\mathbf{y})$  for then  $\boldsymbol{\mu} = \mathbb{E}(\mathbf{z})$  is just a simple moment of  $\mathbf{z}$ . This suggests the moment estimator

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i = \frac{1}{n} \sum_{i=1}^n \mathbf{h}(\mathbf{y}_i).$$

For example, the moment estimator of  $\mathbb{E}(y^r)$  is  $n^{-1} \sum_{i=1}^n y_i^r$ , that of the moment generating function is  $n^{-1} \sum_{i=1}^n \exp(ty_i)$ , and for the distribution function the estimator is  $n^{-1} \sum_{i=1}^n 1\{y_i \leq x\}$ .

Since  $\hat{\boldsymbol{\mu}}$  is a sample average, and transformations of iid variables are also i.i.d., the asymptotic results of the previous sections immediately apply.

**Theorem 6.10.1** *If  $\mathbf{y}_i$  are independent and identically distributed,  $\boldsymbol{\mu} = \mathbb{E}(\mathbf{h}(\mathbf{y}))$ , and  $\mathbb{E}\|\mathbf{h}(\mathbf{y})\| < \infty$ , then for  $\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{h}(\mathbf{y}_i)$ , as  $n \rightarrow \infty$ ,  $\hat{\boldsymbol{\mu}} \xrightarrow{p} \boldsymbol{\mu}$ .*

**Theorem 6.10.2** *If  $\mathbf{y}_i$  are independent and identically distributed,  $\boldsymbol{\mu} = \mathbb{E}(\mathbf{h}(\mathbf{y}))$ , and  $\mathbb{E}\|\mathbf{h}(\mathbf{y})\|^2 < \infty$ , then for  $\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{h}(\mathbf{y}_i)$ , as  $n \rightarrow \infty$ ,*

$$\sqrt{n}(\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}) \xrightarrow{d} N(\mathbf{0}, \mathbf{V})$$

where  $\mathbf{V} = \mathbb{E}((\mathbf{h}(\mathbf{y}) - \boldsymbol{\mu})(\mathbf{h}(\mathbf{y}) - \boldsymbol{\mu})')$ .

Theorems 6.10.1 and 6.10.2 show that the estimate  $\hat{\boldsymbol{\mu}}$  is consistent for  $\boldsymbol{\mu}$  and asymptotically normally distributed, so long as the stated moment conditions hold.

A word of caution. Theorems 6.10.1 and 6.10.2 give the impression that it is possible to estimate any moment of  $y$ . Technically this is the case so long as that moment is finite. What is hidden by the notation, however, is that estimates of high order moments can be quite imprecise. For example, consider the sample 8<sup>th</sup> moment  $\hat{\mu}_8 = \frac{1}{n} \sum_{i=1}^n y_i^8$ , and suppose for simplicity that  $y$  is  $N(0, 1)$ . Then we can calculate<sup>1</sup> that  $\text{var}(\hat{\mu}_8) = n^{-1} 2,016,000$ , which is immense, even for large  $n$ ! In general, higher-order moments are challenging to estimate because their variance depends upon even higher moments which can be quite large in some cases.

<sup>1</sup>By the formula for the variance of a mean  $\text{var}(\hat{\mu}_8) = n^{-1} (\mathbb{E}(y^{16}) - (\mathbb{E}(y^8))^2)$ . Since  $y$  is  $N(0, 1)$ ,  $\mathbb{E}(y^{16}) = 15!! = 2,027,025$  and  $\mathbb{E}(y^8) = 7!! = 105$  where  $k!!$  is the double factorial.



## 6.11 Functions of Moments

We now expand our investigation and consider estimation of parameters which can be written as a continuous function of  $\boldsymbol{\mu} = \mathbb{E}(\mathbf{h}(\mathbf{y}))$ . That is, the parameter of interest can be written as

$$\boldsymbol{\beta} = \mathbf{g}(\boldsymbol{\mu}) = \mathbf{g}(\mathbb{E}(\mathbf{h}(\mathbf{y}))) \quad (6.14)$$

for some functions  $\mathbf{g} : \mathbb{R}^k \rightarrow \mathbb{R}^\ell$  and  $\mathbf{h} : \mathbb{R}^m \rightarrow \mathbb{R}^k$ .

As one example, the geometric mean of wages  $w$  is

$$\gamma = \exp(\mathbb{E}(\log(w))). \quad (6.15)$$

This is (6.14) with  $g(u) = \exp(u)$  and  $h(w) = \log(w)$ .

A simple yet common example is the variance

$$\begin{aligned} \sigma^2 &= \mathbb{E}(w - \mathbb{E}(w))^2 \\ &= \mathbb{E}(w^2) - (\mathbb{E}(w))^2. \end{aligned}$$

This is (6.14) with

$$\mathbf{h}(w) = \begin{pmatrix} w \\ w^2 \end{pmatrix}$$

and

$$g(\mu_1, \mu_2) = \mu_2 - \mu_1^2.$$

Similarly, the skewness of the wage distribution is

$$sk = \frac{\mathbb{E}((w - \mathbb{E}(w))^3)}{(\mathbb{E}((w - \mathbb{E}(w))^2))^{3/2}}.$$

This is (6.14) with

$$\mathbf{h}(w) = \begin{pmatrix} w \\ w^2 \\ w^3 \end{pmatrix}$$

and

$$g(\mu_1, \mu_2, \mu_3) = \frac{\mu_3 - 3\mu_2\mu_1 + 2\mu_1^3}{(\mu_2 - \mu_1^2)^{3/2}}. \quad (6.16)$$

The parameter  $\boldsymbol{\beta} = \mathbf{g}(\boldsymbol{\mu})$  is not a population moment, so it does not have a direct moment estimator. Instead, it is common to use a **plug-in estimate** formed by replacing the unknown  $\boldsymbol{\mu}$  with its point estimate  $\hat{\boldsymbol{\mu}}$  and then “plugging” this into the expression for  $\boldsymbol{\beta}$ . The first step is

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{h}(\mathbf{y}_i)$$

and the second step is

$$\hat{\boldsymbol{\beta}} = \mathbf{g}(\hat{\boldsymbol{\mu}}).$$

Again, the hat “ $\hat{\cdot}$ ” indicates that  $\hat{\boldsymbol{\beta}}$  is a sample estimate of  $\boldsymbol{\beta}$ .

For example, the plug-in estimate of the geometric mean  $\gamma$  of the wage distribution from (6.15) is

$$\hat{\gamma} = \exp(\hat{\mu})$$

with

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \log(\text{wage}_i).$$

The plug-in estimate of the variance is

$$\begin{aligned}\hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n w_i^2 - \left( \frac{1}{n} \sum_{i=1}^n w_i \right)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (w_i - \bar{w})^2.\end{aligned}$$

The estimator for the skewness is

$$\begin{aligned}\widehat{sk} &= \frac{\hat{\mu}_3 - 3\hat{\mu}_2\hat{\mu}_1 + 2\hat{\mu}_1^3}{(\hat{\mu}_2 - \hat{\mu}_1^2)^{3/2}} \\ &= \frac{\frac{1}{n} \sum_{i=1}^n (w_i - \bar{w})^3}{\left( \frac{1}{n} \sum_{i=1}^n (w_i - \bar{w})^2 \right)^{3/2}}\end{aligned}$$

where

$$\hat{\mu}_j = \frac{1}{n} \sum_{i=1}^n w_i^j.$$

A useful property is that continuous functions are limit-preserving.

**Theorem 6.11.1 Continuous Mapping Theorem (CMT).** If  $z_n \xrightarrow{p} c$  as  $n \rightarrow \infty$  and  $g(\cdot)$  is continuous at  $c$ , then  $g(z_n) \xrightarrow{p} g(c)$  as  $n \rightarrow \infty$ .

The proof of Theorem 6.11.1 is given in Section 6.16.

For example, if  $z_n \xrightarrow{p} c$  as  $n \rightarrow \infty$  then

$$\begin{aligned}z_n + a &\xrightarrow{p} c + a \\ az_n &\xrightarrow{p} ac \\ z_n^2 &\xrightarrow{p} c^2\end{aligned}$$

as the functions  $g(u) = u + a$ ,  $g(u) = au$ , and  $g(u) = u^2$  are continuous. Also

$$\frac{a}{z_n} \xrightarrow{p} \frac{a}{c}$$

if  $c \neq 0$ . The condition  $c \neq 0$  is important as the function  $g(u) = a/u$  is not continuous at  $u = 0$ .

If  $y_i$  are independent and identically distributed,  $\mu = \mathbb{E}(\mathbf{h}(\mathbf{y}))$ , and  $\mathbb{E}\|\mathbf{h}(\mathbf{y})\| < \infty$ , then for  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{h}(\mathbf{y}_i)$ , as  $n \rightarrow \infty$ ,  $\hat{\mu} \xrightarrow{p} \mu$ . Applying the CMT,  $\hat{\beta} = g(\hat{\mu}) \xrightarrow{p} g(\mu) = \beta$ .

**Theorem 6.11.2** If  $y_i$  are independent and identically distributed,  $\beta = g(\mathbb{E}(\mathbf{h}(\mathbf{y})))$ ,  $\mathbb{E}\|\mathbf{h}(\mathbf{y})\| < \infty$ , and  $g(\mathbf{u})$  is continuous at  $\mathbf{u} = \mu$ , then for  $\hat{\beta} = g(\frac{1}{n} \sum_{i=1}^n \mathbf{h}(\mathbf{y}_i))$ , as  $n \rightarrow \infty$ ,  $\hat{\beta} \xrightarrow{p} \beta$ .

To apply Theorem 6.11.2 it is necessary to check if the function  $g$  is continuous at  $\mu$ . In our first example  $g(u) = \exp(u)$  is continuous everywhere. It therefore follows from Theorem 6.6.2 and Theorem 6.11.2 that if  $\mathbb{E}|\log(wage)| < \infty$  then as  $n \rightarrow \infty$ ,  $\hat{\gamma} \xrightarrow{p} \gamma$ .

In the example of the variance,  $g$  is continuous for all  $\mu$ . Thus if  $\mathbb{E}(w^2) < \infty$  then as  $n \rightarrow \infty$ ,  $\hat{\sigma}^2 \xrightarrow{p} \sigma^2$ .

In our third example  $g$  defined in (6.16) is continuous for all  $\mu$  such that  $\text{var}(w) = \mu_2 - \mu_1^2 > 0$ , which holds unless  $w$  has a degenerate distribution. Thus if  $\mathbb{E}|w|^3 < \infty$  and  $\text{var}(w) > 0$  then as  $n \rightarrow \infty$ ,  $\widehat{sk} \xrightarrow{p} sk$ .

## 6.12 Delta Method

In this section we introduce two tools – an extended version of the CMT and the Delta Method – which allow us to calculate the asymptotic distribution of the parameter estimate  $\hat{\beta}$ .

We first present an extended version of the continuous mapping theorem which allows convergence in distribution.

### Theorem 6.12.1 Continuous Mapping Theorem

If  $\mathbf{z}_n \xrightarrow{d} \mathbf{z}$  as  $n \rightarrow \infty$  and  $\mathbf{g} : \mathbb{R}^m \rightarrow \mathbb{R}^k$  has the set of discontinuity points  $D_g$  such that  $\Pr(\mathbf{z} \in D_g) = 0$ , then  $\mathbf{g}(\mathbf{z}_n) \xrightarrow{d} \mathbf{g}(\mathbf{z})$  as  $n \rightarrow \infty$ .

For a proof of Theorem 6.12.1 see Theorem 2.3 of van der Vaart (1998). It was first proved by Mann and Wald (1943) and is therefore sometimes referred to as the Mann-Wald Theorem.

Theorem 6.12.1 allows the function  $\mathbf{g}$  to be discontinuous only if the probability at being at a discontinuity point is zero. For example, the function  $g(u) = u^{-1}$  is discontinuous at  $u = 0$ , but if  $\mathbf{z}_n \xrightarrow{d} \mathbf{z} \sim N(0, 1)$  then  $\Pr(\mathbf{z} = 0) = 0$  so  $\mathbf{z}_n^{-1} \xrightarrow{d} \mathbf{z}^{-1}$ .

A special case of the Continuous Mapping Theorem is known as Slutsky's Theorem.

### Theorem 6.12.2 Slutsky's Theorem

If  $\mathbf{z}_n \xrightarrow{d} \mathbf{z}$  and  $c_n \xrightarrow{p} c$  as  $n \rightarrow \infty$ , then

1.  $\mathbf{z}_n + c_n \xrightarrow{d} \mathbf{z} + c$
2.  $\mathbf{z}_n c_n \xrightarrow{d} \mathbf{z} c$
3.  $\frac{\mathbf{z}_n}{c_n} \xrightarrow{d} \frac{\mathbf{z}}{c}$  if  $c \neq 0$

Even though Slutsky's Theorem is a special case of the CMT, it is a useful statement as it focuses on the most common applications – addition, multiplication, and division.

Despite the fact that the plug-in estimator  $\hat{\beta}$  is a function of  $\hat{\mu}$  for which we have an asymptotic distribution, Theorem 6.12.1 does not directly give us an asymptotic distribution for  $\hat{\beta}$ . This is because  $\hat{\beta} = \mathbf{g}(\hat{\mu})$  is written as a function of  $\hat{\mu}$ , not of the standardized sequence  $\sqrt{n}(\hat{\mu} - \mu)$ . We need an intermediate step – a first order Taylor series expansion. This step is so critical to statistical theory that it has its own name – **The Delta Method**.

### Theorem 6.12.3 Delta Method:

If  $\sqrt{n}(\hat{\mu} - \mu) \xrightarrow{d} \xi$ , where  $\mathbf{g}(\mathbf{u})$  is continuously differentiable in a neighborhood of  $\mu$  then as  $n \rightarrow \infty$

$$\sqrt{n}(\mathbf{g}(\hat{\mu}) - \mathbf{g}(\mu)) \xrightarrow{d} \mathbf{G}'\xi \quad (6.17)$$

where  $\mathbf{G}(\mathbf{u}) = \frac{\partial}{\partial \mathbf{u}} \mathbf{g}(\mathbf{u})'$  and  $\mathbf{G} = \mathbf{G}(\mu)$ . In particular, if  $\xi \sim N(0, \mathbf{V})$  then as  $n \rightarrow \infty$

$$\sqrt{n}(\mathbf{g}(\hat{\mu}) - \mathbf{g}(\mu)) \xrightarrow{d} N(0, \mathbf{G}'\mathbf{V}\mathbf{G}). \quad (6.18)$$

The Delta Method allows us to complete our derivation of the asymptotic distribution of the estimator  $\hat{\beta}$  of  $\beta$ . By combining Theorems 6.10.2 and 6.12.3 we can find the asymptotic distribution of the plug-in estimator  $\hat{\beta}$ .

**Theorem 6.12.4** *If  $y_i$  are independent and identically distributed,  $\mu = \mathbb{E}(\mathbf{h}(\mathbf{y}))$ ,  $\beta = \mathbf{g}(\mu)$ ,  $\mathbb{E}\|\mathbf{h}(\mathbf{y})\|^2 < \infty$ , and  $\mathbf{G}(\mathbf{u}) = \frac{\partial}{\partial \mathbf{u}} \mathbf{g}(\mathbf{u})'$  is continuous in a neighborhood of  $\mu$ , then for  $\hat{\beta} = \mathbf{g}(\frac{1}{n} \sum_{i=1}^n \mathbf{h}(\mathbf{y}_i))$ , as  $n \rightarrow \infty$*

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(\mathbf{0}, \mathbf{G}' \mathbf{V} \mathbf{G})$$

*where  $\mathbf{V} = \mathbb{E}((\mathbf{h}(\mathbf{y}) - \mu)(\mathbf{h}(\mathbf{y}) - \mu)')$  and  $\mathbf{G} = \mathbf{G}(\mu)$ .*

Theorem 6.11.2 established the consistency of  $\hat{\beta}$  for  $\beta$ , and Theorem 6.12.4 established its asymptotic normality. It is instructive to compare the conditions required for these results. Consistency required that  $\mathbf{h}(\mathbf{y})$  have a finite mean, while asymptotic normality requires that this variable have a finite variance. Consistency required that  $\mathbf{g}(\mathbf{u})$  be continuous, while our proof of asymptotic normality used the assumption that  $\mathbf{g}(\mathbf{u})$  is continuously differentiable.

## 6.13 Stochastic Order Symbols

It is convenient to have simple symbols for random variables and vectors which converge in probability to zero or are stochastically bounded. In this section we introduce some of the most commonly found notation.

It might be useful to review the common notation for non-random convergence and boundedness. Let  $x_n$  and  $a_n$ ,  $n = 1, 2, \dots$ , be non-random sequences. The notation

$$x_n = o(1)$$

(pronounced “small oh-one”) is equivalent to  $x_n \rightarrow 0$  as  $n \rightarrow \infty$ . The notation

$$x_n = o(a_n)$$

is equivalent to  $a_n^{-1}x_n \rightarrow 0$  as  $n \rightarrow \infty$ . The notation

$$x_n = O(1)$$

(pronounced “big oh-one”) means that  $x_n$  is bounded uniformly in  $n$  – there exists an  $M < \infty$  such that  $|x_n| \leq M$  for all  $n$ . The notation

$$x_n = O(a_n)$$

is equivalent to  $a_n^{-1}x_n = O(1)$ .

We now introduce similar concepts for sequences of random variables. Let  $z_n$  and  $a_n$ ,  $n = 1, 2, \dots$  be sequences of random variables. (In most applications,  $a_n$  is non-random.) The notation

$$z_n = o_p(1)$$

(“small oh-P-one”) means that  $z_n \xrightarrow{p} 0$  as  $n \rightarrow \infty$ . For example, for any consistent estimator  $\hat{\beta}$  for  $\beta$  we can write

$$\hat{\beta} = \beta + o_p(1).$$

We also write

$$z_n = o_p(a_n)$$

if  $a_n^{-1}z_n = o_p(1)$ .

Similarly, the notation  $z_n = O_p(1)$  (“big oh-P-one”) means that  $z_n$  is bounded in probability. Precisely, for any  $\varepsilon > 0$  there is a constant  $M_\varepsilon < \infty$  such that

$$\limsup_{n \rightarrow \infty} \Pr(|z_n| > M_\varepsilon) \leq \varepsilon.$$

Furthermore, we write

$$z_n = O_p(a_n)$$

if  $a_n^{-1}z_n = O_p(1)$ .

$O_p(1)$  is weaker than  $o_p(1)$  in the sense that  $z_n = o_p(1)$  implies  $z_n = O_p(1)$  but not the reverse. However, if  $z_n = O_p(a_n)$  then  $z_n = O_p(b_n)$  for any  $b_n$  such that  $a_n/b_n \rightarrow 0$ .

If a random vector converges in distribution  $\mathbf{z}_n \xrightarrow{d} \mathbf{z}$  (for example, if  $\mathbf{z} \sim N(\mathbf{0}, \mathbf{V})$ ) then  $\mathbf{z}_n = O_p(1)$ . It follows that for estimators  $\hat{\boldsymbol{\beta}}$  which satisfy the convergence of Theorem 6.12.4 then we can write

$$\hat{\boldsymbol{\beta}} = \boldsymbol{\beta} + O_p(n^{-1/2}).$$

In words, this statement says that the estimator  $\hat{\boldsymbol{\beta}}$  equals the true coefficient  $\boldsymbol{\beta}$  plus a random component which is bounded when scaled by  $n^{1/2}$ . Equivalently, we can write

$$n^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = O_p(1).$$

Another useful observation is that a random sequence with a bounded moment is stochastically bounded.

**Theorem 6.13.1** *If  $\mathbf{z}_n$  is a random vector which satisfies*

$$\mathbb{E} \|\mathbf{z}_n\|^\delta = O(a_n)$$

*for some sequence  $a_n$  and  $\delta > 0$ , then*

$$\mathbf{z}_n = O_p(a_n^{1/\delta}).$$

*Similarly,  $\mathbb{E} \|\mathbf{z}_n\|^\delta = o(a_n)$  implies  $\mathbf{z}_n = o_p(a_n^{1/\delta})$ .*

This can be shown using Markov’s inequality (B.14). The assumptions imply that there is some  $M < \infty$  such that  $\mathbb{E} \|\mathbf{z}_n\|^\delta \leq Ma_n$  for all  $n$ . For any  $\varepsilon$  set  $B = \left(\frac{M}{\varepsilon}\right)^{1/\delta}$ . Then

$$\Pr\left(a_n^{-1/\delta} \|\mathbf{z}_n\| > B\right) = \Pr\left(\|\mathbf{z}_n\|^\delta > \frac{Ma_n}{\varepsilon}\right) \leq \frac{\varepsilon}{Ma_n} \mathbb{E} \|\mathbf{z}_n\|^\delta \leq \varepsilon$$

as required.

There are many simple rules for manipulating  $o_p(1)$  and  $O_p(1)$  sequences which can be deduced from the continuous mapping theorem or Slutsky’s Theorem. For example,

$$\begin{aligned} o_p(1) + o_p(1) &= o_p(1) \\ o_p(1) + O_p(1) &= O_p(1) \\ O_p(1) + O_p(1) &= O_p(1) \\ o_p(1)o_p(1) &= o_p(1) \\ o_p(1)O_p(1) &= o_p(1) \\ O_p(1)O_p(1) &= O_p(1). \end{aligned}$$

## 6.14 Uniform Stochastic Bounds\*

For some applications it can be useful to obtain the stochastic order of the random variable

$$\max_{1 \leq i \leq n} |y_i|.$$

This is the magnitude of the largest observation in the sample  $\{y_1, \dots, y_n\}$ . If the support of the distribution of  $y_i$  is unbounded, then as the sample size  $n$  increases, the largest observation will also tend to increase. It turns out that there is a simple characterization.

**Theorem 6.14.1** *If  $y_i$  are identically distributed and  $\mathbb{E}|y|^r < \infty$ , then as  $n \rightarrow \infty$*

$$n^{-1/r} \max_{1 \leq i \leq n} |y_i| \xrightarrow{p} 0. \quad (6.19)$$

*Furthermore, if  $\mathbb{E}(\exp(ty)) < \infty$  for some  $t > 0$  then for any  $\eta > 0$*

$$(\log n)^{-(1+\eta)} \max_{1 \leq i \leq n} |y_i| \xrightarrow{p} 0. \quad (6.20)$$

The proof of Theorem 6.14.1 is presented in Section 6.16.

Equivalently, (6.19) can be written as

$$\max_{1 \leq i \leq n} |y_i| = o_p(n^{1/r}) \quad (6.21)$$

and (6.22) as

$$\max_{1 \leq i \leq n} |y_i| = o_p(\log n). \quad (6.22)$$

Equation (6.21) says that if  $y$  has  $r$  finite moments, then the largest observation will diverge at a rate slower than  $n^{1/r}$ . As  $r$  increases this rate decreases. Equation (6.22) shows that if we strengthen this to  $y$  having all finite moments and a finite moment generating function (for example, if  $y$  is normally distributed) then the largest observation will diverge slower than  $\log n$ . Thus the higher the moments, the slower the rate of divergence.

To simplify the notation, we write (6.21) as  $y_i = o_p(n^{1/r})$  uniformly in  $1 \leq i \leq n$ . It is important to understand when the  $O_p$  or  $o_p$  symbols are applied to subscript  $i$  random variables whether the convergence is pointwise in  $i$ , or is uniform in  $i$  in the sense of (6.21)-(6.22).

Theorem 6.14.1 applies to random vectors. For example, if  $\mathbb{E}\|\mathbf{y}\|^r < \infty$  then

$$\max_{1 \leq i \leq n} \|\mathbf{y}_i\| = o_p(n^{1/r}).$$

## 6.15 Semiparametric Efficiency

In this section we argue that the sample mean  $\hat{\boldsymbol{\mu}}$  and plug-in estimator  $\hat{\boldsymbol{\beta}} = \mathbf{g}(\hat{\boldsymbol{\mu}})$  are efficient estimators of the parameters  $\boldsymbol{\mu}$  and  $\boldsymbol{\beta}$ . Our demonstration is based on the rich but technically challenging theory of semiparametric efficiency bounds. An excellent accessible review has been provided by Newey (1990). We will also appeal to the asymptotic theory of maximum likelihood estimation (see Chapter 5).

We start by examining the sample mean  $\hat{\boldsymbol{\mu}}$ , for the asymptotic efficiency of  $\hat{\boldsymbol{\beta}}$  will follow from that of  $\hat{\boldsymbol{\mu}}$ .

Recall, we know that if  $\mathbb{E}(\|\mathbf{y}\|^2) < \infty$  then the sample mean has the asymptotic distribution  $\sqrt{n}(\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}) \xrightarrow{d} N(0, \mathbf{V})$ . We want to know if  $\hat{\boldsymbol{\mu}}$  is the best feasible estimator, or if there is another

estimator with a smaller asymptotic variance. While it seems intuitively unlikely that another estimator could have a smaller asymptotic variance, how do we know that this is not the case?

When we ask if  $\hat{\mu}$  is the best estimator, we need to be clear about the class of models – the class of permissible distributions. For estimation of the mean  $\mu$  of the distribution of  $\mathbf{y}$  the broadest conceivable class is  $\mathcal{L}_1 = \{F : \mathbb{E} \|\mathbf{y}\| < \infty\}$ . This class is too broad for our current purposes, as  $\hat{\mu}$  is not asymptotically  $N(0, \mathbf{V})$  for all  $F \in \mathcal{L}_1$ . A more realistic choice is  $\mathcal{L}_2 = \left\{F : \mathbb{E} \left( \|\mathbf{y}\|^2 \right) < \infty \right\}$  – the class of finite-variance distributions. When we seek an efficient estimator of the mean  $\mu$  in the class of models  $\mathcal{L}_2$  what we are seeking is the best estimator, given that all we know is that  $F \in \mathcal{L}_2$ .

To show that the answer is not immediately obvious, it might be helpful to review a setting where the sample mean is inefficient. Suppose that  $y \in \mathbb{R}$  has the double exponential density  $f(y | \mu) = 2^{-1/2} \exp(-|y - \mu| \sqrt{2})$ . Since  $\text{var}(y) = 1$  we see that the sample mean satisfies  $\sqrt{n}(\bar{y} - \mu) \xrightarrow{d} N(0, 1)$ . In this model the maximum likelihood estimator (MLE)  $\tilde{\mu}$  for  $\mu$  is the sample median. Recall from the theory of maximum likelihood that the MLE satisfies  $\sqrt{n}(\tilde{\mu} - \mu) \xrightarrow{d} N\left(0, (\mathbb{E}(S^2))^{-1}\right)$  where  $S = \frac{\partial}{\partial \mu} \log f(y | \mu) = -\sqrt{2} \text{sgn}(y - \mu)$  is the score. We can calculate that  $\mathbb{E}(S^2) = 2$  and thus conclude that  $\sqrt{n}(\tilde{\mu} - \mu) \xrightarrow{d} N(0, 1/2)$ . The asymptotic variance of the MLE is one-half that of the sample mean. Thus when the true density is known to be double exponential the sample mean is inefficient.

But the estimator which achieves this improved efficiency – the sample median – is not generically consistent for the population mean. It is inconsistent if the density is asymmetric or skewed. So the improvement comes at a great cost. Another way of looking at this is that the sample median is efficient in the class of densities  $\{f(y | \mu) = 2^{-1/2} \exp(-|y - \mu| \sqrt{2})\}$  but unless it is known that this is the correct distribution class this knowledge is not very useful.

The relevant question is whether or not the sample mean is efficient when the form of the distribution is unknown. We call this setting **semiparametric** as the parameter of interest (the mean) is finite dimensional while the remaining features of the distribution are unspecified. In the semiparametric context an estimator is called **semiparametrically efficient** if it has the smallest asymptotic variance among all semiparametric estimators.

The mathematical trick is to reduce the semiparametric model to a set of parametric “submodels”. The Cramer-Rao variance bound can be found for each parametric submodel. The variance bound for the semiparametric model (the union of the submodels) is then defined as the supremum of the individual variance bounds.

Formally, suppose that the true density of  $\mathbf{y}$  is the unknown function  $f(\mathbf{y})$  with mean  $\mu = \mathbb{E}(\mathbf{y}) = \int \mathbf{y} f(\mathbf{y}) d\mathbf{y}$ . A parametric submodel  $\eta$  for  $f(\mathbf{y})$  is a density  $f_\eta(\mathbf{y} | \theta)$  which is a smooth function of a parameter  $\theta$ , and there is a true value  $\theta_0$  such that  $f_\eta(\mathbf{y} | \theta_0) = f(\mathbf{y})$ . The index  $\eta$  indicates the submodels. The equality  $f_\eta(\mathbf{y} | \theta_0) = f(\mathbf{y})$  means that the submodel class passes through the true density, so the submodel is a true model. The class of submodels  $\eta$  and parameter  $\theta_0$  depend on the true density  $f$ . In the submodel  $f_\eta(\mathbf{y} | \theta)$ , the mean is  $\mu_\eta(\theta) = \int \mathbf{y} f_\eta(\mathbf{y} | \theta) d\mathbf{y}$  which varies with the parameter  $\theta$ . Let  $\eta \in \mathbb{N}$  be the class of all submodels for  $f$ .

Since each submodel  $\eta$  is parametric we can calculate the efficiency bound for estimation of  $\mu$  within this submodel. Specifically, given the density  $f_\eta(\mathbf{y} | \theta)$  its likelihood score is

$$\mathbf{S}_\eta = \frac{\partial}{\partial \theta} \log f_\eta(\mathbf{y} | \theta),$$

so the Cramer-Rao lower bound for estimation of  $\theta$  is  $\left(\mathbb{E}(\mathbf{S}_\eta \mathbf{S}_\eta')\right)^{-1}$ . Defining  $\mathbf{M}_\eta = \frac{\partial}{\partial \theta} \mu_\eta(\theta_0)'$ , by Theorem 5.16.3 the Cramer-Rao lower bound for estimation of  $\mu$  within the submodel  $\eta$  is  $\mathbf{V}_\eta = \mathbf{M}_\eta' \left(\mathbb{E}(\mathbf{S}_\eta \mathbf{S}_\eta')\right)^{-1} \mathbf{M}_\eta$ .

As  $\mathbf{V}_\eta$  is the efficiency bound for the submodel class  $f_\eta(\mathbf{y} | \theta)$ , no estimator can have an asymptotic variance smaller than  $\mathbf{V}_\eta$  for any density  $f_\eta(\mathbf{y} | \theta)$  in the submodel class, including the

true density  $f$ . This is true for all submodels  $\eta$ . Thus the asymptotic variance of any semiparametric estimator cannot be smaller than  $\mathbf{V}_\eta$  for any conceivable submodel. Taking the supremum of the Cramer-Rao bounds from all conceivable submodels we define<sup>2</sup>

$$\bar{\mathbf{V}} = \sup_{\eta \in \mathbb{N}} \mathbf{V}_\eta.$$

The asymptotic variance of any semiparametric estimator cannot be smaller than  $\bar{\mathbf{V}}$ , since it cannot be smaller than any individual  $\mathbf{V}_\eta$ . We call  $\bar{\mathbf{V}}$  the **semiparametric asymptotic variance bound** or **semiparametric efficiency bound** for estimation of  $\boldsymbol{\mu}$ , as it is a lower bound on the asymptotic variance for any semiparametric estimator. If the asymptotic variance of a specific semiparametric estimator equals the bound  $\bar{\mathbf{V}}$  we say that the estimator is **semiparametrically efficient**.

For many statistical problems it is quite challenging to calculate the semiparametric variance bound. However, in some cases there is a simple method to find the solution. Suppose that we can find a submodel  $\eta_0$  whose Cramer-Rao lower bound satisfies  $\mathbf{V}_{\eta_0} = \mathbf{V}_\mu$  where  $\mathbf{V}_\mu$  is the asymptotic variance of a known semiparametric estimator. In this case, we can deduce that  $\bar{\mathbf{V}} = \mathbf{V}_{\eta_0} = \mathbf{V}_\mu$ . Otherwise (that is, if  $\mathbf{V}_{\eta_0}$  is not the efficiency bound) there would exist another submodel  $\eta_1$  whose Cramer-Rao lower bound satisfies  $\mathbf{V}_{\eta_0} < \mathbf{V}_{\eta_1}$  (because  $\mathbf{V}_{\eta_0}$  is not the supremum). This would imply  $\mathbf{V}_\mu < \mathbf{V}_{\eta_1}$  which contradicts the Cramer-Rao Theorem (since when submodel  $\eta_1$  is true then no estimator can have a lower variance than  $\mathbf{V}_{\eta_1}$ ).

We now find this submodel for the sample mean  $\hat{\boldsymbol{\mu}}$ . Our goal is to find a parametric submodel whose Cramer-Rao bound for  $\boldsymbol{\mu}$  is  $\mathbf{V}$ . This can be done by creating a tilted version of the true density. Consider the parametric submodel

$$f_\eta(\mathbf{y} \mid \boldsymbol{\theta}) = f(\mathbf{y}) (1 + \boldsymbol{\theta}' \mathbf{V}^{-1} (\mathbf{y} - \boldsymbol{\mu})) \quad (6.23)$$

where  $f(\mathbf{y})$  is the true density and  $\boldsymbol{\mu} = \mathbb{E}\mathbf{y}$ . Note that

$$\int f_\eta(\mathbf{y} \mid \boldsymbol{\theta}) d\mathbf{y} = \int f(\mathbf{y}) d\mathbf{y} + \boldsymbol{\theta}' \mathbf{V}^{-1} \int f(\mathbf{y}) (\mathbf{y} - \boldsymbol{\mu}) d\mathbf{y} = 1$$

and for all  $\boldsymbol{\theta}$  close to zero  $f_\eta(\mathbf{y} \mid \boldsymbol{\theta}) \geq 0$ . Thus  $f_\eta(\mathbf{y} \mid \boldsymbol{\theta})$  is a valid density function. It is a parametric submodel since  $f_\eta(\mathbf{y} \mid \boldsymbol{\theta}_0) = f(\mathbf{y})$  when  $\boldsymbol{\theta}_0 = 0$ . This parametric submodel has the mean

$$\begin{aligned} \boldsymbol{\mu}(\boldsymbol{\theta}) &= \int \mathbf{y} f_\eta(\mathbf{y} \mid \boldsymbol{\theta}) d\mathbf{y} \\ &= \int \mathbf{y} f(\mathbf{y}) d\mathbf{y} + \int f(\mathbf{y}) \mathbf{y} (\mathbf{y} - \boldsymbol{\mu})' \mathbf{V}^{-1} \boldsymbol{\theta} d\mathbf{y} \\ &= \boldsymbol{\mu} + \boldsymbol{\theta} \end{aligned}$$

which is a smooth function of  $\boldsymbol{\theta}$ .

Since

$$\frac{\partial}{\partial \boldsymbol{\theta}} \log f_\eta(\mathbf{y} \mid \boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \log (1 + \boldsymbol{\theta}' \mathbf{V}^{-1} (\mathbf{y} - \boldsymbol{\mu})) = \frac{\mathbf{V}^{-1} (\mathbf{y} - \boldsymbol{\mu})}{1 + \boldsymbol{\theta}' \mathbf{V}^{-1} (\mathbf{y} - \boldsymbol{\mu})}$$

it follows that the score function for  $\boldsymbol{\theta}$  is

$$\mathbf{S}_\eta = \frac{\partial}{\partial \boldsymbol{\theta}} \log f_\eta(\mathbf{y} \mid \boldsymbol{\theta}_0) = \mathbf{V}^{-1} (\mathbf{y} - \boldsymbol{\mu}). \quad (6.24)$$

By Theorem 5.16.3 the Cramer-Rao lower bound for  $\boldsymbol{\theta}$  is

$$(\mathbb{E}(\mathbf{S}_\eta \mathbf{S}_\eta'))^{-1} = (\mathbf{V}^{-1} \mathbb{E}((\mathbf{y} - \boldsymbol{\mu})(\mathbf{y} - \boldsymbol{\mu})') \mathbf{V}^{-1})^{-1} = \mathbf{V}. \quad (6.25)$$

---

<sup>2</sup>It is not obvious that this supremum exists, as  $\mathbf{V}_\eta$  is a matrix so there is not a unique ordering of matrices. However, in many cases (including the ones we study) the supremum exists and is unique.



The Cramer-Rao lower bound for  $\boldsymbol{\mu}(\boldsymbol{\theta}) = \boldsymbol{\mu} + \boldsymbol{\theta}$  is also  $\mathbf{V}$ , and this equals the asymptotic variance of the moment estimator  $\hat{\boldsymbol{\mu}}$ . This was what we set out to show.

In summary, we have shown that in the submodel (6.23) the Cramer-Rao lower bound for estimation of  $\boldsymbol{\mu}$  is  $\mathbf{V}$  which equals the asymptotic variance of the sample mean. This establishes the following result.

**Proposition 6.15.1** *In the class of distributions  $F \in \mathcal{L}_2$ , the semiparametric variance bound for estimation of  $\boldsymbol{\mu}$  is  $\mathbf{V} = \text{var}(y_i)$ , and the sample mean  $\hat{\boldsymbol{\mu}}$  is a semiparametrically efficient estimator of the population mean  $\boldsymbol{\mu}$ .*

We call this result a proposition rather than a theorem as we have not attended to the regularity conditions.

It is a simple matter to extend this result to the plug-in estimator  $\hat{\boldsymbol{\beta}} = \mathbf{g}(\hat{\boldsymbol{\mu}})$ . We know from Theorem 6.12.4 that if  $\mathbb{E} \|\mathbf{y}\|^2 < \infty$  and  $\mathbf{g}(\mathbf{u})$  is continuously differentiable at  $\mathbf{u} = \boldsymbol{\mu}$  then the plug-in estimator has the asymptotic distribution  $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} N(0, \mathbf{G}'\mathbf{V}\mathbf{G})$ . We therefore consider the class of distributions

$$\mathcal{L}_2(\mathbf{g}) = \left\{ F : \mathbb{E} \|\mathbf{y}\|^2 < \infty, \mathbf{g}(\mathbf{u}) \text{ is continuously differentiable at } \mathbf{u} = \mathbb{E}(\mathbf{y}) \right\}.$$

For example, if  $\boldsymbol{\beta} = \mu_1/\mu_2$  where  $\mu_1 = \mathbb{E}(y_1)$  and  $\mu_2 = \mathbb{E}(y_2)$  then

$$\mathcal{L}_2(g) = \left\{ F : \mathbb{E}(y_1^2) < \infty, \mathbb{E}(y_2^2) < \infty, \text{ and } \mathbb{E}(y_2) \neq 0 \right\}.$$

For any submodel  $\eta$  the Cramer-Rao lower bound for estimation of  $\boldsymbol{\beta} = \mathbf{g}(\boldsymbol{\mu})$  is  $\mathbf{G}'\mathbf{V}_\eta\mathbf{G}$ . For the submodel (6.23) this bound is  $\mathbf{G}'\mathbf{V}\mathbf{G}$  which equals the asymptotic variance of  $\hat{\boldsymbol{\beta}}$  from Theorem 6.12.4. Thus  $\hat{\boldsymbol{\beta}}$  is semiparametrically efficient.

**Proposition 6.15.2** *In the class of distributions  $F \in \mathcal{L}_2(\mathbf{g})$  the semiparametric variance bound for estimation of  $\boldsymbol{\beta} = \mathbf{g}(\boldsymbol{\mu})$  is  $\mathbf{G}'\mathbf{V}\mathbf{G}$ , and the plug-in estimator  $\hat{\boldsymbol{\beta}} = \mathbf{g}(\hat{\boldsymbol{\mu}})$  is a semiparametrically efficient estimator of  $\boldsymbol{\beta}$ .*

The result in Proposition 6.15.2 is quite general. Smooth functions of sample moments are efficient estimators for their population counterparts. This is a very powerful result, as most econometric estimators can be written (or approximated) as smooth functions of sample means.

## 6.16 Technical Proofs\*

In this section we provide proofs of some of the more technical points in the chapter. These proofs may only be of interest to more mathematically inclined students.

**Proof of Theorem 6.4.2:** Without loss of generality, we can assume  $\mathbb{E}(y_i) = 0$  by recentering  $y_i$  on its expectation.

We need to show that for all  $\delta > 0$  and  $\eta > 0$  there is some  $N < \infty$  so that for all  $n \geq N$ ,  $\Pr(|\bar{y}| > \delta) \leq \eta$ . Fix  $\delta$  and  $\eta$ . Set  $\varepsilon = \delta\eta/3$ . Pick  $C < \infty$  large enough so that

$$\mathbb{E}(|y_i| \mathbf{1}(|y_i| > C)) \leq \varepsilon \tag{6.26}$$

(where  $1(\cdot)$  is the indicator function) which is possible since  $\mathbb{E}|y_i| < \infty$ . Define the random variables

$$\begin{aligned} w_i &= y_i 1(|y_i| \leq C) - \mathbb{E}(y_i 1(|y_i| \leq C)) \\ z_i &= y_i 1(|y_i| > C) - \mathbb{E}(y_i 1(|y_i| > C)) \end{aligned}$$

so that

$$\bar{y} = \bar{w} + \bar{z}$$

and

$$\mathbb{E}|\bar{y}| \leq \mathbb{E}|\bar{w}| + \mathbb{E}|\bar{z}|. \quad (6.27)$$

We now show that sum of the expectations on the right-hand-side can be bounded below  $3\varepsilon$ .

First, by the Triangle Inequality (A.26) and the Expectation Inequality (B.8),

$$\begin{aligned} \mathbb{E}|z_i| &= \mathbb{E}|y_i 1(|y_i| > C) - \mathbb{E}(y_i 1(|y_i| > C))| \\ &\leq \mathbb{E}|y_i 1(|y_i| > C)| + |\mathbb{E}(y_i 1(|y_i| > C))| \\ &\leq 2\mathbb{E}|y_i 1(|y_i| > C)| \\ &\leq 2\varepsilon, \end{aligned} \quad (6.28)$$

and thus by the Triangle Inequality (A.26) and (6.28)

$$\mathbb{E}|\bar{z}| = \mathbb{E}\left|\frac{1}{n} \sum_{i=1}^n z_i\right| \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}|z_i| \leq 2\varepsilon. \quad (6.29)$$

Second, by a similar argument

$$\begin{aligned} |w_i| &= |y_i 1(|y_i| \leq C) - \mathbb{E}(y_i 1(|y_i| \leq C))| \\ &\leq |y_i 1(|y_i| \leq C)| + |\mathbb{E}(y_i 1(|y_i| \leq C))| \\ &\leq 2|y_i 1(|y_i| \leq C)| \\ &\leq 2C \end{aligned} \quad (6.30)$$

where the final inequality is (6.26). Then by Jensen's Inequality (B.5), the fact that the  $w_i$  are iid and mean zero, and (6.30),

$$(\mathbb{E}|\bar{w}|)^2 \leq \mathbb{E}(|\bar{w}|^2) = \frac{\mathbb{E}(w_i^2)}{n} \leq \frac{4C^2}{n} \leq \varepsilon^2 \quad (6.31)$$

the final inequality holding for  $n \geq 4C^2/\varepsilon^2 = 36C^2/\delta^2\eta^2$ . Equations (6.27), (6.29) and (6.31) together show that

$$\mathbb{E}|\bar{y}| \leq 3\varepsilon \quad (6.32)$$

as desired.

Finally, by Markov's Inequality (B.14) and (6.32),

$$\Pr(|\bar{y}| > \delta) \leq \frac{\mathbb{E}|\bar{y}|}{\delta} \leq \frac{3\varepsilon}{\delta} = \eta,$$

the final equality by the definition of  $\varepsilon$ . We have shown that for any  $\delta > 0$  and  $\eta > 0$  then for all  $n \geq 36C^2/\delta^2\eta^2$ ,  $\Pr(|\bar{y}| > \delta) \leq \eta$ , as needed. ■

**Proof of Theorem 6.6.1:** By Loève's  $c_r$  Inequality (A.16)

$$\|\mathbf{y}\| = \left( \sum_{j=1}^m y_j^2 \right)^{1/2} \leq \sum_{j=1}^m |y_j|.$$

Thus if  $\mathbb{E}|y_j| < \infty$  for  $j = 1, \dots, m$ , then

$$\mathbb{E}\|\mathbf{y}\| \leq \sum_{j=1}^m \mathbb{E}|y_j| < \infty.$$

For the reverse inequality, the Euclidean norm of a vector is larger than the length of any individual component, so for any  $j$ ,  $|y_j| \leq \|\mathbf{y}\|$ . Thus, if  $\mathbb{E}\|\mathbf{y}\| < \infty$ , then  $\mathbb{E}|y_j| < \infty$  for  $j = 1, \dots, m$ . ■

**Proof of Theorem 6.7.2:** By Lévy's Continuity Theorem (Theorem 6.7.1),  $\mathbf{z}_n \xrightarrow{d} \mathbf{z}$  if and only if  $\mathbb{E}(\exp(\mathbf{i}s'\mathbf{z}_n)) \rightarrow \mathbb{E}(\exp(\mathbf{i}s'\mathbf{z}))$  for every  $\mathbf{s} \in \mathbb{R}^k$ . We can write  $\mathbf{s} = t\boldsymbol{\lambda}$  where  $t \in \mathbb{R}$  and  $\boldsymbol{\lambda} \in \mathbb{R}^k$  with  $\boldsymbol{\lambda}'\boldsymbol{\lambda} = 1$ . Thus the convergence holds if and only if  $\mathbb{E}(\exp(it\boldsymbol{\lambda}'\mathbf{z}_n)) \rightarrow \mathbb{E}(\exp(it\boldsymbol{\lambda}'\mathbf{z}))$  for every  $t \in \mathbb{R}$  and  $\boldsymbol{\lambda} \in \mathbb{R}^k$  with  $\boldsymbol{\lambda}'\boldsymbol{\lambda} = 1$ . Again by Lévy's Continuity Theorem, this holds if and only if  $\boldsymbol{\lambda}'\mathbf{z}_n \xrightarrow{d} \boldsymbol{\lambda}'\mathbf{z}$  for every  $\boldsymbol{\lambda} \in \mathbb{R}^k$  and with  $\boldsymbol{\lambda}'\boldsymbol{\lambda} = 1$ . ■

**Proof of Theorem 6.8.1:** The moment bound  $\mathbb{E}(y_i^2) < \infty$  is sufficient to guarantee that  $\mu$  and  $\sigma^2$  are well defined and finite. Without loss of generality, it is sufficient to consider the case  $\mu = 0$ .

Our proof method is to calculate the characteristic function of  $\sqrt{n}\bar{y}_n$  and show that it converges pointwise to  $\exp(-\lambda^2\sigma^2/2)$ , the characteristic function of  $N(0, \sigma^2)$ . By Lévy's Continuity Theorem (Theorem 6.7.1) this implies  $\sqrt{n}\bar{y}_n \xrightarrow{d} N(0, \sigma^2)$ .

Let  $C(t) = \mathbb{E}\exp(it y_i)$  denote the characteristic function of  $y_i$  and set  $c(t) = \log C(t)$ , which is sometimes called the cumulant generating function. We start by calculating a second order Taylor series expansion of  $c(t)$  about  $t = 0$ , which requires computing the first two derivatives of  $c(t)$  at  $t = 0$ . These derivatives are

$$\begin{aligned} c'(t) &= \frac{C'(t)}{C(t)} \\ c''(t) &= \frac{C''(t)}{C(t)} - \left( \frac{C'(t)}{C(t)} \right)^2. \end{aligned}$$

Using (2.61) and  $\mu = 0$  we find

$$\begin{aligned} c(0) &= 0 \\ c'(0) &= 0 \\ c''(0) &= -\sigma^2. \end{aligned}$$

Then the second-order Taylor series expansion of  $c(t)$  about  $t = 0$  equals

$$\begin{aligned} c(t) &= c(0) + c'(0)t + \frac{1}{2}c''(t^*)t^2 \\ &= \frac{1}{2}c''(t^*)t^2 \end{aligned} \tag{6.33}$$

where  $t^*$  lies on the line segment joining 0 and  $t$ .

We now compute  $C_n(t) = \mathbb{E}\exp(it\sqrt{n}\bar{y}_n)$ , the characteristic function of  $\sqrt{n}\bar{y}_n$ . By the properties

of the exponential function, the independence of the  $y_i$ , and the definition of  $c(t)$

$$\begin{aligned}
 \log C_n(t) &= \log \mathbb{E} \left( \exp \left( i \frac{1}{\sqrt{n}} \sum_{i=1}^n t y_i \right) \right) \\
 &= \log \mathbb{E} \left( \prod_{i=1}^n \exp \left( i \frac{1}{\sqrt{n}} t y_i \right) \right) \\
 &= \log \prod_{i=1}^n \mathbb{E} \left( \exp \left( i \frac{1}{\sqrt{n}} t y_i \right) \right) \\
 &= \sum_{i=1}^n \log \mathbb{E} \left( \exp \left( i \frac{1}{\sqrt{n}} t y_i \right) \right) \\
 &= n c \left( \frac{t}{\sqrt{n}} \right) \\
 &= \frac{1}{2} c''(t_n) t^2
 \end{aligned}$$

For  $n$  large the argument  $t/\sqrt{n}$  is in a neighborhood of 0. Since the second moment of  $y_i$  is finite,  $c''(t)$  is continuous at  $t = 0$ . Thus we can apply a second order Taylor series expansion about 0, and apply  $c(0) = c'(0) = 0$  to find that

$$\begin{aligned}
 \log C_n(t) &= n c \left( \frac{t}{\sqrt{n}} \right) \\
 &= n \left( c(0) + c'(0) \frac{t}{\sqrt{n}} + \frac{1}{2} c'' \left( \frac{t_n}{\sqrt{n}} \right) \left( \frac{t}{\sqrt{n}} \right)^2 \right) \\
 &= \frac{1}{2} c'' \left( \frac{t_n}{\sqrt{n}} \right) t^2
 \end{aligned}$$

where  $t_n$  lies on the line segment joining 0 and  $t$ . Since  $t_n$  is bounded we deduce that  $c''(t_n/\sqrt{n}) \rightarrow c''(0) = -\sigma^2$ . Hence, as  $n \rightarrow \infty$ ,

$$\log C_n(t) \rightarrow -\frac{1}{2} \sigma^2 t^2$$

and

$$C_n(t) \rightarrow \exp \left( -\frac{1}{2} \sigma^2 t^2 \right)$$

which is the characteristic function of the  $N(0, \sigma^2)$  distribution, as shown in Exercise 5.9. This completes the proof. ■

**Proof of Theorem 6.8.3:** Suppose that  $\sigma^2 = 0$ . Then  $\text{var}(\sqrt{n}(\bar{y} - \mathbb{E}(\bar{y}))) = \bar{\sigma}_n^2 \rightarrow \sigma^2 = 0$  so  $\sqrt{n}(\bar{y} - \mathbb{E}(\bar{y})) \xrightarrow{p} 0$  and hence  $\sqrt{n}(\bar{y} - \mathbb{E}(\bar{y})) \xrightarrow{d} 0$ . The random variable  $N(0, \sigma^2) = N(0, 0)$  is 0 with probability 1, so this is  $\sqrt{n}(\bar{y} - \mathbb{E}(\bar{y})) \xrightarrow{d} N(0, \sigma^2)$  as stated.

Now suppose that  $\sigma^2 > 0$ . This implies (6.8). Together with (6.7) this implies Lyapunov's condition, and hence Lindeberg's condition, and hence Theorem 6.8.2, which states

$$\frac{\sqrt{n}(\bar{y} - \mathbb{E}(\bar{y}))}{\bar{\sigma}_n^{1/2}} \xrightarrow{d} N(0, 1).$$

Combined with (6.9) we deduce  $\sqrt{n}(\bar{y} - \mathbb{E}(\bar{y})) \xrightarrow{d} N(0, \sigma^2)$  as stated. ■

**Proof of Theorem 6.9.1:** Set  $\lambda \in \mathbb{R}^k$  with  $\lambda' \lambda = 1$  and define  $u_i = \lambda' (\mathbf{y}_i - \boldsymbol{\mu})$ . The  $u_i$  are i.i.d with  $\mathbb{E}(u_i^2) = \lambda' \mathbf{V} \lambda < \infty$ . By Theorem 6.8.1,

$$\lambda' \sqrt{n} (\bar{\mathbf{y}} - \boldsymbol{\mu}) = \frac{1}{\sqrt{n}} \sum_{i=1}^n u_i \xrightarrow{d} N(0, \lambda' \mathbf{V} \lambda)$$

Notice that if  $\mathbf{z} \sim N(\mathbf{0}, \mathbf{V})$  then  $\lambda' \mathbf{z} \sim N(0, \lambda' \mathbf{V} \lambda)$ . Thus

$$\lambda' \sqrt{n} (\bar{\mathbf{y}} - \boldsymbol{\mu}) \xrightarrow{d} \lambda' \mathbf{z}.$$

Since this holds for all  $\lambda$ , the conditions of Theorem 6.7.2 are satisfied and we deduce that

$$\sqrt{n} (\bar{\mathbf{y}} - \boldsymbol{\mu}) \xrightarrow{d} \mathbf{z} \sim N(\mathbf{0}, \mathbf{V})$$

as stated. ■

**Proof of Theorem 6.9.2:** Set  $\lambda \in \mathbb{R}^k$  with  $\lambda' \lambda = 1$  and define  $u_{ni} = \lambda' \bar{\mathbf{V}}_n^{-1/2} (\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni})$ . Notice that  $u_{ni}$  are independent and has variance  $\sigma_{ni}^2 = \lambda' \bar{\mathbf{V}}_n^{-1/2} \mathbf{V}_{ni} \bar{\mathbf{V}}_n^{-1/2} \lambda$  and  $\bar{\sigma}_n^2 = n^{-1} \sum_{i=1}^n \sigma_{ni}^2 = 1$ . It is sufficient to verify (6.5). By the Cauchy-Schwarz inequality,

$$\begin{aligned} u_{ni}^2 &= \left( \lambda' \bar{\mathbf{V}}_n^{-1/2} (\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni}) \right)^2 \\ &\leq \lambda' \bar{\mathbf{V}}_n^{-1} \lambda \|\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni}\|^2 \\ &\leq \frac{\|\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni}\|^2}{\lambda_{\min}(\bar{\mathbf{V}}_n)} \\ &= \frac{\|\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni}\|^2}{\nu_n^2}. \end{aligned}$$

Then

$$\begin{aligned} \frac{1}{n \bar{\sigma}_n^2} \sum_{i=1}^n \mathbb{E}(u_{ni}^2 1(u_{ni}^2 \geq \varepsilon n \bar{\sigma}_n^2)) &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}(u_{ni}^2 1(u_{ni}^2 \geq \varepsilon n)) \\ &\leq \frac{1}{n \nu_n^2} \sum_{i=1}^n \mathbb{E}(\|\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni}\|^2 1(\|\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni}\|^2 \geq \varepsilon n \nu_n^2)) \\ &\rightarrow 0 \end{aligned}$$

by (6.11). This establishes (6.5). We deduce from Theorem 6.8.2 that

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n u_{ni} = \lambda' \sqrt{n} \bar{\mathbf{V}}_n^{-1/2} (\bar{\mathbf{y}} - \mathbb{E}(\bar{\mathbf{y}})) \xrightarrow{d} N(0, 1) = \lambda' \mathbf{z}$$

where  $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I}_k)$ . Since this holds for all  $\lambda$ , the conditions of Theorem 6.7.2 are satisfied and we deduce that

$$\sqrt{n} \bar{\mathbf{V}}_n^{-1/2} (\bar{\mathbf{y}} - \mathbb{E}(\bar{\mathbf{y}})) \xrightarrow{d} N(0, \mathbf{I}_k)$$

as stated. ■

**Proof of Theorem 6.9.3:** Set  $\lambda \in \mathbb{R}^k$  with  $\lambda' \lambda = 1$  and define  $u_{ni} = \lambda' (\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni})$ . Using the triangle inequality and (6.13) we obtain

$$\sup_{n,i} \mathbb{E}(|u_{ni}|^{2+\delta}) \leq \sup_{n,i} \mathbb{E}(\|\mathbf{y}_{ni} - \boldsymbol{\mu}_{ni}\|^{2+\delta}) < \infty$$

which is (6.7). Notice that

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}(u_{ni}^2) = \boldsymbol{\lambda}' \frac{1}{n} \sum_{i=1}^n \mathbf{V}_{ni} \boldsymbol{\lambda} = \boldsymbol{\lambda}' \overline{\mathbf{V}}_n \boldsymbol{\lambda} \rightarrow \boldsymbol{\lambda}' \mathbf{V} \boldsymbol{\lambda}$$

which is (6.9). Since the  $u_{ni}$  are independent, by Theorem 6.9.1,

$$\boldsymbol{\lambda}' \sqrt{n} (\overline{\mathbf{y}} - \mathbb{E}(\overline{\mathbf{y}})) = \frac{1}{\sqrt{n}} \sum_{i=1}^n u_{ni} \xrightarrow{d} N(0, \boldsymbol{\lambda}' \mathbf{V} \boldsymbol{\lambda}) = \boldsymbol{\lambda}' \mathbf{z}$$

where  $\mathbf{z} \sim N(\mathbf{0}, \mathbf{V})$ . Since this holds for all  $\boldsymbol{\lambda}$ , the conditions of Theorem 6.7.2 are satisfied and we deduce that

$$\sqrt{n} (\overline{\mathbf{y}} - \mathbb{E}(\overline{\mathbf{y}})) \xrightarrow{d} N(\mathbf{0}, \mathbf{V})$$

as stated. ■

**Proof of Theorem 6.12.3:** By a vector Taylor series expansion, for each element of  $\mathbf{g}$ ,

$$g_j(\boldsymbol{\theta}_n) = g_j(\boldsymbol{\theta}) + g_{j\boldsymbol{\theta}}(\boldsymbol{\theta}_{jn}^*) (\boldsymbol{\theta}_n - \boldsymbol{\theta})$$

where  $\boldsymbol{\theta}_{jn}^*$  lies on the line segment between  $\boldsymbol{\theta}_n$  and  $\boldsymbol{\theta}$  and therefore converges in probability to  $\boldsymbol{\theta}$ . It follows that  $a_{jn} = g_{j\boldsymbol{\theta}}(\boldsymbol{\theta}_{jn}^*) - g_{j\boldsymbol{\theta}} \xrightarrow{p} 0$ . Stacking across elements of  $\mathbf{g}$ , we find

$$\sqrt{n} (\mathbf{g}(\boldsymbol{\theta}_n) - \mathbf{g}(\boldsymbol{\theta})) = (\mathbf{G} + \mathbf{a}_n)' \sqrt{n} (\boldsymbol{\theta}_n - \boldsymbol{\theta}) \xrightarrow{d} \mathbf{G}' \boldsymbol{\xi}. \quad (6.34)$$

The convergence is by Theorem 6.12.1, as  $\mathbf{G} + \mathbf{a}_n \xrightarrow{d} \mathbf{G}$ ,  $\sqrt{n} (\boldsymbol{\theta}_n - \boldsymbol{\theta}) \xrightarrow{d} \boldsymbol{\xi}$ , and their product is continuous. This establishes (6.17)

When  $\boldsymbol{\xi} \sim N(\mathbf{0}, \mathbf{V})$ , the right-hand-side of (6.34) equals

$$\mathbf{G}' \boldsymbol{\xi} = \mathbf{G}' N(\mathbf{0}, \mathbf{V}) = N(\mathbf{0}, \mathbf{G}' \mathbf{V} \mathbf{G})$$

establishing (6.18). ■

**Proof of Theorem 6.14.1:** First consider (6.19). Take any  $\delta > 0$ . The event  $\{\max_{1 \leq i \leq n} |y_i| > \delta n^{1/r}\}$  means that at least one of the  $|y_i|$  exceeds  $\delta n^{1/r}$ , which is the same as the event  $\bigcup_{i=1}^n \{|y_i| > \delta n^{1/r}\}$  or equivalently  $\bigcup_{i=1}^n \{|y_i|^r > \delta^r n\}$ . Since the probability of the union of events is smaller than the sum of the probabilities,

$$\begin{aligned} \Pr \left( n^{-1/r} \max_{1 \leq i \leq n} |y_i| > \delta \right) &= \Pr \left( \bigcup_{i=1}^n \{|y_i|^r > \delta^r n\} \right) \\ &\leq \sum_{i=1}^n \Pr(|y_i|^r > n \delta^r) \\ &\leq \frac{1}{n \delta^r} \sum_{i=1}^n \mathbb{E}(|y_i|^r \mathbf{1}(|y_i|^r > n \delta^r)) \\ &= \frac{1}{\delta^r} \mathbb{E}(|y_1|^r \mathbf{1}(|y_1|^r > n \delta^r)) \end{aligned}$$

where the second inequality is the strong form of Markov's inequality (Theorem B.15) and the final equality is since the  $y_i$  are iid. Since  $\mathbb{E}(|y|^r) < \infty$  this final expectation converges to zero as  $n \rightarrow \infty$ . This is because

$$\mathbb{E}(|y|^r) = \int |y|^r dF(y) < \infty$$

implies

$$\mathbb{E}(|y_i|^r \mathbf{1}(|y_i|^r > c)) = \int_{|y|^r > c} |y|^r dF(y) \rightarrow 0 \quad (6.35)$$

as  $c \rightarrow \infty$ . This establishes (6.19).

Now consider (6.20). Take any  $\delta > 0$  and pick  $n$  large enough so that  $(\log n)^\eta t\delta \geq 1$ . By a similar calculation

$$\begin{aligned} \Pr\left((\log n)^{-(1+\eta)} \max_{1 \leq i \leq n} |y_i| > \delta\right) &= \Pr\left(\bigcup_{i=1}^n \left\{\exp |ty_i| > \exp\left((\log n)^{1+\eta} t\delta\right)\right\}\right) \\ &\leq \sum_{i=1}^n \Pr(\exp |ty_i| > n) \\ &\leq \mathbb{E}(\exp |ty| \mathbf{1}(\exp |ty| > n)) \end{aligned}$$

where the second line uses  $\exp\left((\log n)^{1+\eta} t\delta\right) \geq \exp(\log n) = n$ . The assumption  $\mathbb{E}(\exp(ty)) < \infty$  means  $\mathbb{E}(\exp |ty| \mathbf{1}(\exp |ty| > n)) \rightarrow 0$  as  $n \rightarrow \infty$  by the same argument as in (6.35). This establishes (6.20). ■

## Exercises

**Exercise 6.1** For the following sequences, show  $a_n \rightarrow 0$  as  $n \rightarrow \infty$

- (a)  $a_n = 1/n$
- (b)  $a_n = \frac{1}{n} \sin\left(\frac{\pi}{2}n\right)$

**Exercise 6.2** Does the sequence  $a_n = \sin\left(\frac{\pi}{2}n\right)$  converge? Find the liminf and limsup as  $n \rightarrow \infty$ .

**Exercise 6.3** A weighted sample mean takes the form  $\bar{y}^* = \frac{1}{n} \sum_{i=1}^n w_i y_i$  for some non-negative constants  $w_i$  satisfying  $\frac{1}{n} \sum_{i=1}^n w_i = 1$ . Assume  $y_i$  is iid.

- (a) Show that  $\bar{y}^*$  is unbiased for  $\mu = \mathbb{E}(y_i)$ .
- (b) Calculate  $\text{var}(\bar{y}^*)$ .
- (c) Show that a sufficient condition for  $\bar{y}^* \xrightarrow{p} \mu$  is that  $\frac{1}{n^2} \sum_{i=1}^n w_i^2 \rightarrow 0$ .
- (d) Show that a sufficient condition for the condition in part 3 is  $\max_{i \leq n} w_i = o(n)$ .

**Exercise 6.4** Consider a random variable  $X_n$  with the probability distribution

$$X_n = \begin{cases} -n & \text{with probability } 1/n \\ 0 & \text{with probability } 1 - 2/n \\ n & \text{with probability } 1/n \end{cases}$$

- (a) Does  $X_n \rightarrow_p 0$  as  $n \rightarrow \infty$ ?
- (b) Calculate  $\mathbb{E}(X_n)$
- (c) Calculate  $\text{var}(X_n)$
- (d) Now suppose the distribution is

$$X_n = \begin{cases} 0 & \text{with probability } 1 - n \\ n & \text{with probability } 1/n \end{cases}$$

Calculate  $\mathbb{E}(X_n)$

- (e) Conclude that  $X_n \rightarrow_p 0$  as  $n \rightarrow \infty$  and  $\mathbb{E}(X_n) \rightarrow 0$  are unrelated.

**Exercise 6.5** A weighted sample mean takes the form  $\bar{y}^* = \frac{1}{n} \sum_{i=1}^n w_i y_i$  for some non-negative constants  $w_i$  satisfying  $\frac{1}{n} \sum_{i=1}^n w_i = 1$ . Assume  $y_i$  is iid.

- (a) Show that  $\bar{y}^*$  is unbiased for  $\mu = \mathbb{E}(y_i)$ .
- (b) Calculate  $\text{var}(\bar{y}^*)$ .
- (c) Show that a sufficient condition for  $\bar{y}^* \xrightarrow{p} \mu$  is that  $\frac{1}{n^2} \sum_{i=1}^n w_i^2 \rightarrow 0$ .
- (d) Show that a sufficient condition for the condition in part c is  $\max_{i \leq n} w_i/n \rightarrow 0$ .

**Exercise 6.6** Take a random sample  $\{y_1, \dots, y_n\}$ . Which statistics converge in probability by the weak law of large numbers and continuous mapping theorem, assuming the moment exists?

- (a)  $\frac{1}{n} \sum_{i=1}^n y_i^2$



- (b)  $\frac{1}{n} \sum_{i=1}^n y_i^3$
- (c)  $\max_{i \leq n} y_i$
- (d)  $\frac{1}{n} \sum_{i=1}^n y_i^2 - \left(\frac{1}{n} \sum_{i=1}^n y_i\right)^2$
- (e)  $\frac{\sum_{i=1}^n y_i^2}{\sum_{i=1}^n y_i}$  assuming  $\mathbb{E}(y_i) > 0$
- (f)  $1\left(\frac{1}{n} \sum_{i=1}^n y_i > 0\right)$  where

$$1(a) = \begin{cases} 1 & \text{if } a \text{ is true} \\ 0 & \text{if } a \text{ is not true} \end{cases}$$

**Exercise 6.7** Take a random sample  $\{X_1, \dots, X_n\}$  where  $X > 0$ . Consider the sample geometric mean

$$\hat{\mu} = \left( \prod_{i=1}^n X_i \right)^{1/n}$$

and population geometric mean

$$\mu = \exp(\mathbb{E}(\log X))$$

Assuming  $\mu$  is finite, show that  $\hat{\mu} \rightarrow_p \mu$  as  $n \rightarrow \infty$ .

**Exercise 6.8** Take a random variable  $Z$  such that  $\mathbb{E}(Z) = 0$  and  $\text{var}(Z) = 1$ . Use Chebyshev's inequality to find a  $\delta$  such that  $\Pr(|Z| > \delta) \leq 0.05$ . Contrast this with the exact  $\delta$  which solves  $\Pr(|Z| > \delta) = 0.05$  when  $Z \sim N(0, 1)$ . Comment on the difference.

**Exercise 6.9** Find the moment estimator  $\hat{\mu}_3$  of  $\mu_3 = \mathbb{E}(y_i^3)$  and show that  $\sqrt{n}(\hat{\mu}_3 - \mu_3) \xrightarrow{d} N(0, v^2)$  for some  $v^2$ . Write  $v^2$  as a function of the moments of  $y_i$ .

**Exercise 6.10** Suppose  $z_n \xrightarrow{p} c$  as  $n \rightarrow \infty$ . Show that  $z_n^2 \xrightarrow{p} c^2$  as  $n \rightarrow \infty$  using the definition of convergence in probability, but not appealing to the CMT.

**Exercise 6.11** Let  $\mu_k = \mathbb{E}(y^k)$  for some integer  $k \geq 1$ .

- (a) Write down the natural moment estimator  $\hat{\mu}_k$  of  $\mu_k$ .
- (b) Find the asymptotic distribution of  $\sqrt{n}(\hat{\mu}_k - \mu_k)$  as  $n \rightarrow \infty$ . (Assume  $\mathbb{E}(X^{2k}) < \infty$ .)

**Exercise 6.12** Let  $m_k = (\mathbb{E}(y^k))^{1/k}$  for some integer  $k \geq 1$ .

- (a) Write down an estimator  $\hat{m}_k$  of  $m_k$ .
- (b) Find the asymptotic distribution of  $\sqrt{n}(\hat{m}_k - m_k)$  as  $n \rightarrow \infty$ .

**Exercise 6.13** Suppose  $\sqrt{n}(\hat{\mu} - \mu) \xrightarrow{d} N(0, v^2)$  and set  $\beta = \mu^2$  and  $\hat{\beta} = \hat{\mu}^2$ .

- (a) Use the Delta Method to obtain an asymptotic distribution for  $\sqrt{n}(\hat{\beta} - \beta)$ .
- (b) Now suppose  $\mu = 0$ . Describe what happens to the asymptotic distribution from the previous part.
- (c) Improve on the previous answer. Under the assumption  $\mu = 0$ , find the asymptotic distribution for  $n\hat{\beta} = n\hat{\mu}^2$ .
- (d) Comment on the differences between the answers in parts 1 and 3.

**Exercise 6.14** Let  $y$  be distributed Bernoulli  $P(y = 1) = p$  and  $P(y = 0) = 1 - p$  for some unknown  $0 < p < 1$ .

- (a) Show that  $p = \mathbb{E}(y)$
- (b) Write down the natural moment estimator  $\hat{p}$  of  $p$ .
- (c) Find  $\text{var}(\hat{p})$
- (d) Find the asymptotic distribution of  $\sqrt{n}(\hat{p} - p)$  as  $n \rightarrow \infty$ .

## Chapter 7

# Asymptotic Theory for Least Squares

### 7.1 Introduction

It turns out that the asymptotic theory of least-squares estimation applies equally to the projection model and the linear CEF model, and therefore the results in this chapter will be stated for the broader projection model described in Section 2.18. Recall that the model is

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

for  $i = 1, \dots, n$ , where the linear projection  $\boldsymbol{\beta}$  is

$$\boldsymbol{\beta} = (\mathbb{E}(\mathbf{x}_i \mathbf{x}_i'))^{-1} \mathbb{E}(\mathbf{x}_i y_i).$$

Some of the results of this section hold under random sampling (Assumption 1.5.2) and finite second moments (Assumption 2.18.1). We restate this condition here for clarity.

**Assumption 7.1.1**

1. *The observations  $(y_i, \mathbf{x}_i)$ ,  $i = 1, \dots, n$ , are independent and identically distributed.*
2.  $\mathbb{E}(y^2) < \infty$ .
3.  $\mathbb{E}\|\mathbf{x}\|^2 < \infty$ .
4.  $\mathbf{Q}_{\mathbf{x}\mathbf{x}} = \mathbb{E}(\mathbf{x}\mathbf{x}')$  is positive definite.

Some of the results will require a strengthening to finite fourth moments.

**Assumption 7.1.2** *In addition to Assumption 7.1.1,  $\mathbb{E}(y_i^4) < \infty$  and  $\mathbb{E}\|\mathbf{x}_i\|^4 < \infty$ .*

## 7.2 Consistency of Least-Squares Estimator

In this section we use the weak law of large numbers (WLLN, Theorem 6.4.2 and Theorem 6.6.2) and continuous mapping theorem (CMT, Theorem 6.11.1) to show that the least-squares estimator  $\hat{\beta}$  is consistent for the projection coefficient  $\beta$ .

This derivation is based on three key components. First, the OLS estimator can be written as a continuous function of a set of sample moments. Second, the WLLN shows that sample moments converge in probability to population moments. And third, the CMT states that continuous functions preserve convergence in probability. We now explain each step in brief and then in greater detail.

First, observe that the OLS estimator

$$\hat{\beta} = \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i y_i \right) = \hat{\mathbf{Q}}_{\mathbf{x}\mathbf{x}}^{-1} \hat{\mathbf{Q}}_{\mathbf{x}y}$$

is a function of the sample moments  $\hat{\mathbf{Q}}_{\mathbf{x}\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'$  and  $\hat{\mathbf{Q}}_{\mathbf{x}y} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i y_i$ .

Second, by an application of the WLLN these sample moments converge in probability to the population moments. Specifically, the fact that  $(y_i, \mathbf{x}_i)$  are mutually independent and identically distributed implies that any function of  $(y_i, \mathbf{x}_i)$  is iid, including  $\mathbf{x}_i \mathbf{x}_i'$  and  $\mathbf{x}_i y_i$ . These variables also have finite expectations under Assumption 7.1.1. Under these conditions, the WLLN (Theorem 6.6.2) implies that as  $n \rightarrow \infty$ ,

$$\hat{\mathbf{Q}}_{\mathbf{x}\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \xrightarrow{p} \mathbb{E}(\mathbf{x}_i \mathbf{x}_i') = \mathbf{Q}_{\mathbf{x}\mathbf{x}} \quad (7.1)$$

and

$$\hat{\mathbf{Q}}_{\mathbf{x}y} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i y_i \xrightarrow{p} \mathbb{E}(\mathbf{x}_i y_i) = \mathbf{Q}_{\mathbf{x}y}. \quad (7.2)$$

Third, the CMT (Theorem 6.11.1) allows us to combine these equations to show that  $\hat{\beta}$  converges in probability to  $\beta$ . Specifically, as  $n \rightarrow \infty$ ,

$$\begin{aligned} \hat{\beta} &= \hat{\mathbf{Q}}_{\mathbf{x}\mathbf{x}}^{-1} \hat{\mathbf{Q}}_{\mathbf{x}y} \\ &\xrightarrow{p} \mathbf{Q}_{\mathbf{x}\mathbf{x}}^{-1} \mathbf{Q}_{\mathbf{x}y} \\ &= \beta. \end{aligned} \quad (7.3)$$

We have shown that  $\hat{\beta} \xrightarrow{p} \beta$ , as  $n \rightarrow \infty$ . In words, the OLS estimator converges in probability to the projection coefficient vector  $\beta$  as the sample size  $n$  gets large.

To fully understand the application of the CMT we walk through it in detail. We can write

$$\hat{\beta} = g(\hat{\mathbf{Q}}_{\mathbf{x}\mathbf{x}}, \hat{\mathbf{Q}}_{\mathbf{x}y})$$

where  $g(\mathbf{A}, \mathbf{b}) = \mathbf{A}^{-1} \mathbf{b}$  is a function of  $\mathbf{A}$  and  $\mathbf{b}$ . The function  $g(\mathbf{A}, \mathbf{b})$  is a continuous function of  $\mathbf{A}$  and  $\mathbf{b}$  at all values of the arguments such that  $\mathbf{A}^{-1}$  exists. Assumption 7.1.1 specifies that  $\mathbf{Q}_{\mathbf{x}\mathbf{x}}^{-1}$  exists and thus  $g(\mathbf{A}, \mathbf{b})$  is continuous at  $\mathbf{A} = \mathbf{Q}_{\mathbf{x}\mathbf{x}}$ . This justifies the application of the CMT in (7.3).

For a slightly different demonstration of (7.3), recall that (4.7) implies that

$$\hat{\beta} - \beta = \hat{\mathbf{Q}}_{\mathbf{x}\mathbf{x}}^{-1} \hat{\mathbf{Q}}_{\mathbf{x}e} \quad (7.4)$$

where

$$\hat{\mathbf{Q}}_{\mathbf{x}e} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i e_i.$$

The WLLN and (2.27) imply

$$\widehat{\mathbf{Q}}_{xe} \xrightarrow{p} \mathbb{E}(\mathbf{x}_i e_i) = \mathbf{0}. \quad (7.5)$$

Therefore

$$\begin{aligned} \widehat{\boldsymbol{\beta}} - \boldsymbol{\beta} &= \widehat{\mathbf{Q}}_{xx}^{-1} \widehat{\mathbf{Q}}_{xe} \\ &\xrightarrow{p} \mathbf{Q}_{xx}^{-1} \mathbf{0} \\ &= \mathbf{0} \end{aligned}$$

which is the same as  $\widehat{\boldsymbol{\beta}} \xrightarrow{p} \boldsymbol{\beta}$ .

**Theorem 7.2.1 Consistency of Least-Squares**

Under Assumption 7.1.1,  $\widehat{\mathbf{Q}}_{xx} \xrightarrow{p} \mathbf{Q}_{xx}$ ,  $\widehat{\mathbf{Q}}_{xy} \xrightarrow{p} \mathbf{Q}_{xy}$ ,  $\widehat{\mathbf{Q}}_{xx}^{-1} \xrightarrow{p} \mathbf{Q}_{xx}^{-1}$ ,  $\widehat{\mathbf{Q}}_{xe} \xrightarrow{p} \mathbf{0}$ , and  $\widehat{\boldsymbol{\beta}} \xrightarrow{p} \boldsymbol{\beta}$  as  $n \rightarrow \infty$ .

Theorem 7.2.1 states that the OLS estimator  $\widehat{\boldsymbol{\beta}}$  converges in probability to  $\boldsymbol{\beta}$  as  $n$  increases, and thus  $\widehat{\boldsymbol{\beta}}$  is consistent for  $\boldsymbol{\beta}$ . In the stochastic order notation, Theorem 7.2.1 can be equivalently written as

$$\widehat{\boldsymbol{\beta}} = \boldsymbol{\beta} + o_p(1). \quad (7.6)$$

To illustrate the effect of sample size on the least-squares estimator consider the least-squares regression

$$\ln(\text{Wage}_i) = \beta_1 \text{Education}_i + \beta_2 \text{Experience}_i + \beta_3 \text{Experience}_i^2 + \beta_4 + e_i.$$

We use the sample of 24,344 white men from the March 2009 CPS. Randomly sorting the observations, and sequentially estimating the model by least-squares, starting with the first 5 observations, and continuing until the full sample is used, the sequence of estimates are displayed in Figure 7.1. You can see how the least-squares estimate changes with the sample size, but as the number of observations increases it settles down to the full-sample estimate  $\widehat{\beta}_1 = 0.114$ .

### 7.3 Asymptotic Normality

We started this chapter discussing the need for an approximation to the distribution of the OLS estimator  $\widehat{\boldsymbol{\beta}}$ . In Section 7.2 we showed that  $\widehat{\boldsymbol{\beta}}$  converges in probability to  $\boldsymbol{\beta}$ . Consistency is a good first step, but in itself does not describe the distribution of the estimator. In this section we derive an approximation typically called the **asymptotic distribution**.

The derivation starts by writing the estimator as a function of sample moments. One of the moments must be written as a sum of zero-mean random vectors and normalized so that the central limit theorem can be applied. The steps are as follows.

Take equation (7.4) and multiply it by  $\sqrt{n}$ . This yields the expression

$$\sqrt{n}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i e_i \right). \quad (7.7)$$

This shows that the normalized and centered estimator  $\sqrt{n}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})$  is a function of the sample average  $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'$  and the normalized sample average  $\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i e_i$ . Furthermore, the latter has mean zero so the central limit theorem (CLT, Theorem 6.8.1) applies.

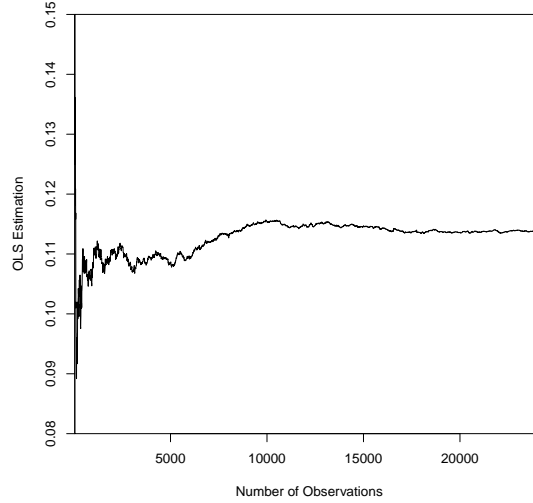


Figure 7.1: The least-squares estimator  $\hat{\beta}_1$  as a function of sample size  $n$

The product  $\mathbf{x}_i e_i$  is iid (since the observations are iid) and mean zero (since  $\mathbb{E}(\mathbf{x}_i e_i) = \mathbf{0}$ ). Define the  $k \times k$  covariance matrix

$$\mathbf{\Omega} = \mathbb{E}(\mathbf{x}_i \mathbf{x}_i' e_i^2). \quad (7.8)$$

We require the elements of  $\mathbf{\Omega}$  to be finite, written  $\mathbf{\Omega} < \infty$ . It will be useful to recall that Theorem 2.18.1.6 shows that Assumption 7.1.2 implies that  $\mathbb{E}(e_i^4) < \infty$ .

The  $j\ell^{th}$  element of  $\mathbf{\Omega}$  is  $\mathbb{E}(x_{ji} x_{\ell i} e_i^2)$ . By the Expectation Inequality (B.8), the  $j\ell^{th}$  element of  $\mathbf{\Omega}$  is

$$|\mathbb{E}(x_{ji} x_{\ell i} e_i^2)| \leq \mathbb{E}|x_{ji} x_{\ell i} e_i^2| = \mathbb{E}(|x_{ji}| |x_{\ell i}| e_i^2).$$

By two applications of the Cauchy-Schwarz Inequality (B.10), this is smaller than

$$(\mathbb{E}(x_{ji}^2 x_{\ell i}^2))^{1/2} (\mathbb{E}(e_i^4))^{1/2} \leq (\mathbb{E}(x_{ji}^4))^{1/4} (\mathbb{E}(x_{\ell i}^4))^{1/4} (\mathbb{E}(e_i^4))^{1/2} < \infty$$

where the finiteness holds under Assumption 7.1.2.

An alternative way to show that the elements of  $\mathbf{\Omega}$  are finite is by using a matrix norm  $\|\cdot\|$  (See Appendix A.18). Then by the Expectation Inequality, the Cauchy-Schwarz Inequality, and Assumption 7.1.2

$$\|\mathbf{\Omega}\| \leq \mathbb{E}\|\mathbf{x}_i \mathbf{x}_i' e_i^2\| = \mathbb{E}(\|\mathbf{x}_i\|^2 e_i^2) \leq (\mathbb{E}\|\mathbf{x}_i\|^4)^{1/2} (\mathbb{E}(e_i^4))^{1/2} < \infty.$$

This is a more compact argument (often described as more *elegant*) but such manipulations should not be done without understanding the notation and the applicability of each step of the argument.

Regardless, the finiteness of the covariance matrix means that we can then apply the CLT (Theorem 6.8.1).

**Theorem 7.3.1** Under Assumption 7.1.2,

$$\mathbf{\Omega} < \infty \quad (7.9)$$

and

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i e_i \xrightarrow{d} N(\mathbf{0}, \mathbf{\Omega}) \quad (7.10)$$

as  $n \rightarrow \infty$ .

Putting together (7.1), (7.7), and (7.10),

$$\begin{aligned}\sqrt{n}(\hat{\beta} - \beta) &\xrightarrow{d} \mathbf{Q}_{xx}^{-1} \mathbf{N}(\mathbf{0}, \mathbf{\Omega}) \\ &= \mathbf{N}(\mathbf{0}, \mathbf{Q}_{xx}^{-1} \mathbf{\Omega} \mathbf{Q}_{xx}^{-1})\end{aligned}$$

as  $n \rightarrow \infty$ , where the final equality follows from the property that linear combinations of normal vectors are also normal (Theorem 5.2.3).

We have derived the asymptotic normal approximation to the distribution of the least-squares estimator.

**Theorem 7.3.2 Asymptotic Normality of Least-Squares Estimator**

*Under Assumption 7.1.2, as  $n \rightarrow \infty$*

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{V}_{\beta})$$

where

$$\mathbf{V}_{\beta} = \mathbf{Q}_{xx}^{-1} \mathbf{\Omega} \mathbf{Q}_{xx}^{-1}, \quad (7.11)$$

$$\mathbf{Q}_{xx} = \mathbb{E}(\mathbf{x}_i \mathbf{x}_i'), \text{ and } \mathbf{\Omega} = \mathbb{E}(\mathbf{x}_i \mathbf{x}_i' e_i^2).$$

In the stochastic order notation, Theorem 7.3.2 implies that

$$\hat{\beta} = \beta + O_p(n^{-1/2}) \quad (7.12)$$

which is stronger than (7.6).

The matrix  $\mathbf{V}_{\beta} = \mathbf{Q}_{xx}^{-1} \mathbf{\Omega} \mathbf{Q}_{xx}^{-1}$  is the variance of the asymptotic distribution of  $\sqrt{n}(\hat{\beta} - \beta)$ . Consequently,  $\mathbf{V}_{\beta}$  is often referred to as the **asymptotic covariance matrix** of  $\hat{\beta}$ . The expression  $\mathbf{V}_{\beta} = \mathbf{Q}_{xx}^{-1} \mathbf{\Omega} \mathbf{Q}_{xx}^{-1}$  is called a **sandwich** form, as the matrix  $\mathbf{\Omega}$  is sandwiched between two copies of  $\mathbf{Q}_{xx}^{-1}$ .

It is useful to compare the variance of the asymptotic distribution given in (7.11) and the finite-sample conditional variance in the CEF model as given in (4.12):

$$\mathbf{V}_{\hat{\beta}} = \text{var}(\hat{\beta} | \mathbf{X}) = (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{D}\mathbf{X}) (\mathbf{X}'\mathbf{X})^{-1}. \quad (7.13)$$

Notice that  $\mathbf{V}_{\hat{\beta}}$  is the exact conditional variance of  $\hat{\beta}$  and  $\mathbf{V}_{\beta}$  is the asymptotic variance of  $\sqrt{n}(\hat{\beta} - \beta)$ . Thus  $\mathbf{V}_{\beta}$  should be (roughly)  $n$  times as large as  $\mathbf{V}_{\hat{\beta}}$ , or  $\mathbf{V}_{\beta} \approx n\mathbf{V}_{\hat{\beta}}$ . Indeed, multiplying (7.13) by  $n$  and distributing, we find

$$n\mathbf{V}_{\hat{\beta}} = \left(\frac{1}{n}\mathbf{X}'\mathbf{X}\right)^{-1} \left(\frac{1}{n}\mathbf{X}'\mathbf{D}\mathbf{X}\right) \left(\frac{1}{n}\mathbf{X}'\mathbf{X}\right)^{-1}$$

which looks like an estimator of  $\mathbf{V}_{\beta}$ . Indeed, as  $n \rightarrow \infty$

$$n\mathbf{V}_{\hat{\beta}} \xrightarrow{p} \mathbf{V}_{\beta}.$$

The expression  $\mathbf{V}_{\hat{\beta}}$  is useful for practical inference (such as computation of standard errors and tests) since it is the variance of the estimator  $\hat{\beta}$ , while  $\mathbf{V}_{\beta}$  is useful for asymptotic theory as it

is well defined in the limit as  $n$  goes to infinity. We will make use of both symbols and it will be advisable to adhere to this convention.

There is a special case where  $\mathbf{\Omega}$  and  $\mathbf{V}_\beta$  simplify. Suppose that

$$\text{cov}(\mathbf{x}_i \mathbf{x}_i', e_i^2) = \mathbf{0}. \quad (7.14)$$

Condition (7.14) holds in the homoskedastic linear regression model, but is somewhat broader. Under (7.14) the asymptotic variance formulae simplify as

$$\mathbf{\Omega} = \mathbb{E}(\mathbf{x}_i \mathbf{x}_i') \mathbb{E}(e_i^2) = \mathbf{Q}_{xx} \sigma^2 \quad (7.15)$$

$$\mathbf{V}_\beta = \mathbf{Q}_{xx}^{-1} \mathbf{\Omega} \mathbf{Q}_{xx}^{-1} = \mathbf{Q}_{xx}^{-1} \sigma^2 \equiv \mathbf{V}_\beta^0 \quad (7.16)$$

In (7.16) we define  $\mathbf{V}_\beta^0 = \mathbf{Q}_{xx}^{-1} \sigma^2$  whether (7.14) is true or false. When (7.14) is true then  $\mathbf{V}_\beta = \mathbf{V}_\beta^0$ , otherwise  $\mathbf{V}_\beta \neq \mathbf{V}_\beta^0$ . We call  $\mathbf{V}_\beta^0$  the **homoskedastic asymptotic covariance matrix**.

Theorem 7.3.2 states that the sampling distribution of the least-squares estimator, after rescaling, is approximately normal when the sample size  $n$  is sufficiently large. This holds true for all joint distributions of  $(y_i, \mathbf{x}_i)$  which satisfy the conditions of Assumption 7.1.2, and is therefore broadly applicable. Consequently, asymptotic normality is routinely used to approximate the finite sample distribution of  $\sqrt{n}(\hat{\beta} - \beta)$ .

A difficulty is that for any fixed  $n$  the sampling distribution of  $\hat{\beta}$  can be arbitrarily far from the normal distribution. In Figure 6.1 we have already seen a simple example where the least-squares estimate is quite asymmetric and non-normal even for reasonably large sample sizes. The normal approximation improves as  $n$  increases, but how large should  $n$  be in order for the approximation to be useful? Unfortunately, there is no simple answer to this reasonable question. The trouble is that no matter how large is the sample size, the normal approximation is arbitrarily poor for some data distribution satisfying the assumptions. We illustrate this problem using a simulation. Let  $y_i = \beta_1 x_i + \beta_2 + e_i$  where  $x_i$  is  $N(0, 1)$ , and  $e_i$  is independent of  $x_i$  with the Double Pareto density  $f(e) = \frac{\alpha}{2} |e|^{-\alpha-1}$ ,  $|e| \geq 1$ . If  $\alpha > 2$  the error  $e_i$  has zero mean and variance  $\alpha/(\alpha - 2)$ . As  $\alpha$  approaches 2, however, its variance diverges to infinity. In this context the normalized least-squares slope estimator  $\sqrt{n \frac{\alpha-2}{\alpha}} (\hat{\beta}_1 - \beta_1)$  has the  $N(0, 1)$  asymptotic distribution for any  $\alpha > 2$ .

In Figure 7.2 we display the finite sample densities of the normalized estimator  $\sqrt{n \frac{\alpha-2}{\alpha}} (\hat{\beta}_1 - \beta_1)$ , setting  $n = 100$  and varying the parameter  $\alpha$ . For  $\alpha = 3.0$  the density is very close to the  $N(0, 1)$  density. As  $\alpha$  diminishes the density changes significantly, concentrating most of the probability mass around zero.

Another example is shown in Figure 7.3. Here the model is  $y_i = \beta + e_i$  where

$$e_i = \frac{u_i^r - \mathbb{E}(u_i^r)}{\left(\mathbb{E}(u_i^{2r}) - (\mathbb{E}(u_i^r))^2\right)^{1/2}} \quad (7.17)$$

and  $u_i \sim N(0, 1)$  and some integer  $r \geq 1$ . We show the sampling distribution of  $\sqrt{n}(\hat{\beta} - \beta)$  setting  $n = 100$ , for  $r = 1, 4, 6$  and  $8$ . As  $r$  increases, the sampling distribution becomes highly skewed and non-normal. The lesson from Figures 7.2 and 7.3 is that the  $N(0, 1)$  asymptotic approximation is never guaranteed to be accurate.

## 7.4 Joint Distribution

Theorem 7.3.2 gives the joint asymptotic distribution of the coefficient estimates. We can use the result to study the covariance between the coefficient estimates. For simplicity, suppose  $k = 2$  with no intercept, both regressors are mean zero and the error is homoskedastic. Let  $\sigma_1^2$  and  $\sigma_2^2$  be



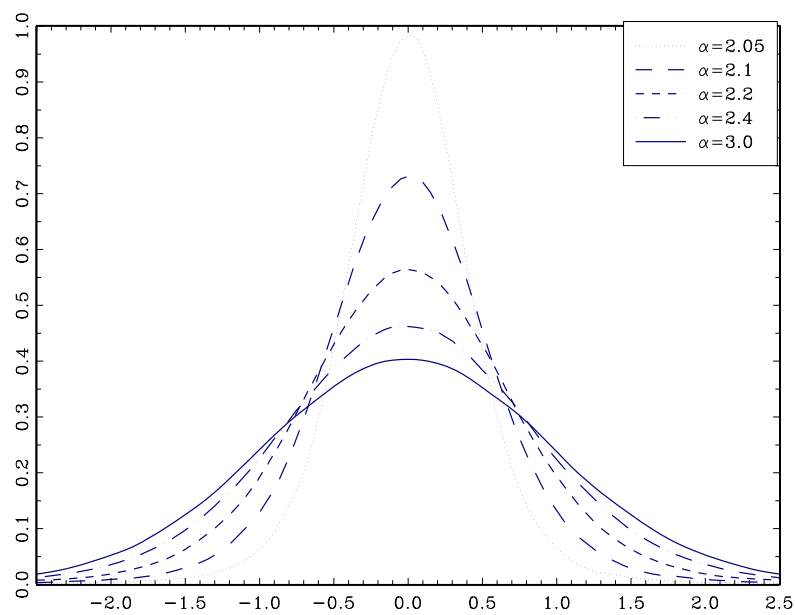


Figure 7.2: Density of Normalized OLS estimator with Double Pareto Error

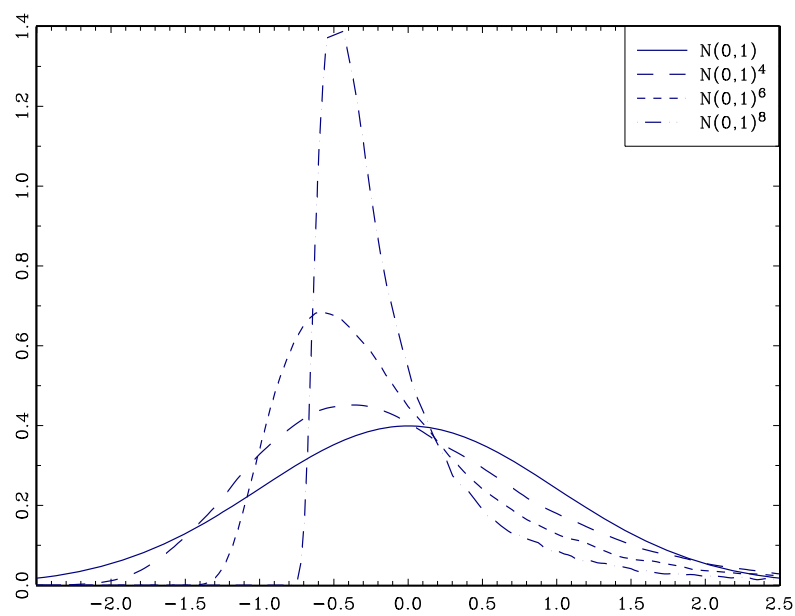


Figure 7.3: Density of Normalized OLS estimator with error process (7.17)

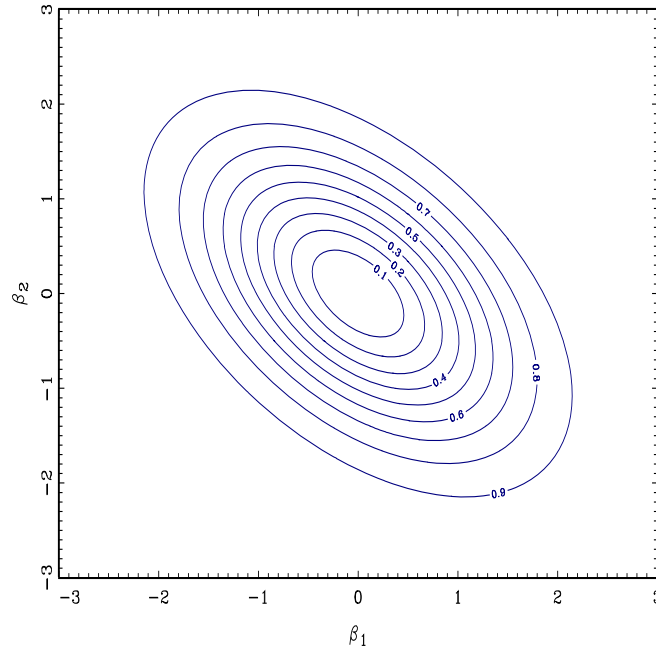


Figure 7.4: Contours of Joint Distribution of  $(\hat{\beta}_1, \hat{\beta}_2)$ , homoskedastic case

the variances of  $x_{1i}$  and  $x_{2i}$ , and  $\rho$  be their correlation. Then using the formula for inversion of a  $2 \times 2$  matrix,

$$\mathbf{V}_{\beta}^0 = \sigma^2 \mathbf{Q}_{xx}^{-1} = \frac{\sigma^2}{\sigma_1^2 \sigma_2^2 (1 - \rho^2)} \begin{bmatrix} \sigma_2^2 & -\rho \sigma_1 \sigma_2 \\ -\rho \sigma_1 \sigma_2 & \sigma_1^2 \end{bmatrix}.$$

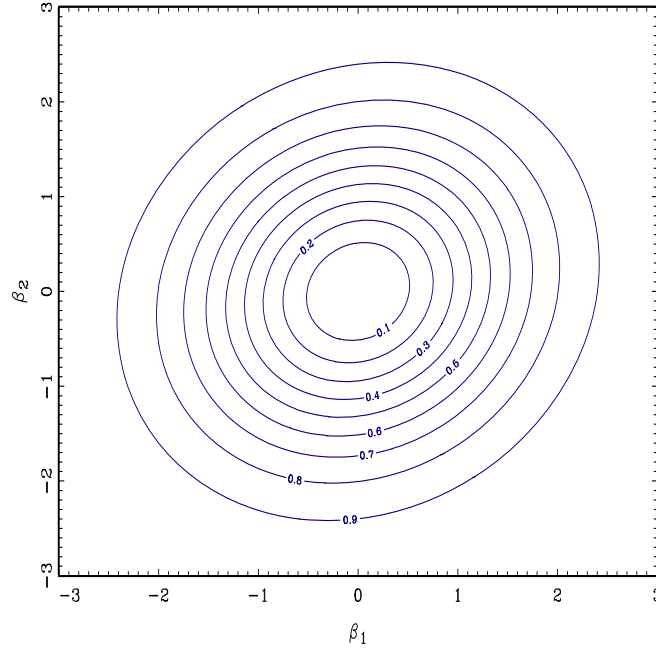
Thus if  $x_{1i}$  and  $x_{2i}$  are positively correlated ( $\rho > 0$ ) then  $\hat{\beta}_1$  and  $\hat{\beta}_2$  are negatively correlated (and vice-versa).

For illustration, Figure 7.4 displays the probability contours of the joint asymptotic distribution of  $\hat{\beta}_1 - \beta_1$  and  $\hat{\beta}_2 - \beta_2$  when  $\beta_1 = \beta_2 = 0$ ,  $\sigma_1^2 = \sigma_2^2 = \sigma^2 = 1$ , and  $\rho = 0.5$ . The coefficient estimates are negatively correlated since the regressors are positively correlated. This means that if  $\hat{\beta}_1$  is unusually negative, it is likely that  $\hat{\beta}_2$  is unusually positive, or conversely. It is also unlikely that we will observe both  $\hat{\beta}_1$  and  $\hat{\beta}_2$  unusually large and of the same sign.

This finding that the correlation of the regressors is of opposite sign of the correlation of the coefficient estimates is sensitive to the assumption of homoskedasticity. If the errors are heteroskedastic then this relationship is not guaranteed.

This can be seen through a simple constructed example. Suppose that  $x_{1i}$  and  $x_{2i}$  only take the values  $\{-1, +1\}$ , symmetrically, with  $\Pr(x_{1i} = x_{2i} = 1) = \Pr(x_{1i} = x_{2i} = -1) = 3/8$ , and  $\Pr(x_{1i} = 1, x_{2i} = -1) = \Pr(x_{1i} = -1, x_{2i} = 1) = 1/8$ . You can check that the regressors are mean zero, unit variance and correlation 0.5, which is identical with the setting displayed in Figure 7.4.

Now suppose that the error is heteroskedastic. Specifically, suppose that  $\mathbb{E}(e_i^2 | x_{1i} = x_{2i}) = \frac{5}{4}$  and  $\mathbb{E}(e_i^2 | x_{1i} \neq x_{2i}) = \frac{1}{4}$ . You can check that  $\mathbb{E}(e_i^2) = 1$ ,  $\mathbb{E}(x_{1i}^2 e_i^2) = \mathbb{E}(x_{2i}^2 e_i^2) = 1$  and

Figure 7.5: Contours of Joint Distribution of  $\hat{\beta}_1$  and  $\hat{\beta}_2$ , heteroskedastic case

$\mathbb{E}(x_{1i}x_{2i}e_i^2) = \frac{7}{8}$ . Therefore

$$\begin{aligned} \mathbf{V}_\beta &= \mathbf{Q}_{xx}^{-1} \mathbf{\Omega} \mathbf{Q}_{xx}^{-1} \\ &= \frac{9}{16} \begin{bmatrix} 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{bmatrix} \begin{bmatrix} 1 & \frac{7}{8} \\ \frac{7}{8} & 1 \end{bmatrix} \begin{bmatrix} 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{bmatrix} \\ &= \frac{4}{3} \begin{bmatrix} 1 & \frac{1}{4} \\ \frac{1}{4} & 1 \end{bmatrix}. \end{aligned}$$

Thus the coefficient estimates  $\hat{\beta}_1$  and  $\hat{\beta}_2$  are positively correlated (their correlation is  $1/4$ .) The joint probability contours of their asymptotic distribution is displayed in Figure 7.5. We can see how the two estimates are positively associated.

What we found through this example is that in the presence of heteroskedasticity there is no simple relationship between the correlation of the regressors and the correlation of the parameter estimates.

We can extend the above analysis to study the covariance between coefficient sub-vectors. For example, partitioning  $\mathbf{x}'_i = (\mathbf{x}'_{1i}, \mathbf{x}'_{2i})$  and  $\boldsymbol{\beta}' = (\boldsymbol{\beta}'_1, \boldsymbol{\beta}'_2)$ , we can write the general model as

$$y_i = \mathbf{x}'_{1i}\boldsymbol{\beta}_1 + \mathbf{x}'_{2i}\boldsymbol{\beta}_2 + e_i$$

and the coefficient estimates as  $\hat{\boldsymbol{\beta}}' = (\hat{\boldsymbol{\beta}}'_1, \hat{\boldsymbol{\beta}}'_2)$ . Make the partitions

$$\mathbf{Q}_{xx} = \begin{bmatrix} \mathbf{Q}_{11} & \mathbf{Q}_{12} \\ \mathbf{Q}_{21} & \mathbf{Q}_{22} \end{bmatrix}, \quad \mathbf{\Omega} = \begin{bmatrix} \mathbf{\Omega}_{11} & \mathbf{\Omega}_{12} \\ \mathbf{\Omega}_{21} & \mathbf{\Omega}_{22} \end{bmatrix}. \quad (7.18)$$

From (2.41)

$$\mathbf{Q}_{xx}^{-1} = \begin{bmatrix} \mathbf{Q}_{11}^{-1} & -\mathbf{Q}_{11}^{-1}\mathbf{Q}_{12}\mathbf{Q}_{22}^{-1} \\ -\mathbf{Q}_{22}^{-1}\mathbf{Q}_{21}\mathbf{Q}_{11}^{-1} & \mathbf{Q}_{22}^{-1} \end{bmatrix}$$

where  $\mathbf{Q}_{11.2} = \mathbf{Q}_{11} - \mathbf{Q}_{12}\mathbf{Q}_{22}^{-1}\mathbf{Q}_{21}$  and  $\mathbf{Q}_{22.1} = \mathbf{Q}_{22} - \mathbf{Q}_{21}\mathbf{Q}_{11}^{-1}\mathbf{Q}_{12}$ . Thus when the error is homoskedastic,

$$\text{cov}(\hat{\beta}_1, \hat{\beta}_2) = -\sigma^2 \mathbf{Q}_{11.2}^{-1} \mathbf{Q}_{12} \mathbf{Q}_{22}^{-1}$$

which is a matrix generalization of the two-regressor case.

In the general case, you can show that (Exercise 7.5)

$$\mathbf{V}_\beta = \begin{bmatrix} \mathbf{V}_{11} & \mathbf{V}_{12} \\ \mathbf{V}_{21} & \mathbf{V}_{22} \end{bmatrix} \quad (7.19)$$

where

$$\mathbf{V}_{11} = \mathbf{Q}_{11.2}^{-1} (\mathbf{\Omega}_{11} - \mathbf{Q}_{12}\mathbf{Q}_{22}^{-1}\mathbf{\Omega}_{21} - \mathbf{\Omega}_{12}\mathbf{Q}_{22}^{-1}\mathbf{Q}_{21} + \mathbf{Q}_{12}\mathbf{Q}_{22}^{-1}\mathbf{\Omega}_{22}\mathbf{Q}_{22}^{-1}\mathbf{Q}_{21}) \mathbf{Q}_{11.2}^{-1} \quad (7.20)$$

$$\mathbf{V}_{21} = \mathbf{Q}_{22.1}^{-1} (\mathbf{\Omega}_{21} - \mathbf{Q}_{21}\mathbf{Q}_{11}^{-1}\mathbf{\Omega}_{11} - \mathbf{\Omega}_{22}\mathbf{Q}_{22}^{-1}\mathbf{Q}_{21} + \mathbf{Q}_{21}\mathbf{Q}_{11}^{-1}\mathbf{\Omega}_{12}\mathbf{Q}_{22}^{-1}\mathbf{Q}_{21}) \mathbf{Q}_{11.2}^{-1} \quad (7.21)$$

$$\mathbf{V}_{22} = \mathbf{Q}_{22.1}^{-1} (\mathbf{\Omega}_{22} - \mathbf{Q}_{21}\mathbf{Q}_{11}^{-1}\mathbf{\Omega}_{12} - \mathbf{\Omega}_{21}\mathbf{Q}_{11}^{-1}\mathbf{Q}_{12} + \mathbf{Q}_{21}\mathbf{Q}_{11}^{-1}\mathbf{\Omega}_{11}\mathbf{Q}_{11}^{-1}\mathbf{Q}_{12}) \mathbf{Q}_{22.1}^{-1} \quad (7.22)$$

Unfortunately, these expressions are not easily interpretable.

## 7.5 Consistency of Error Variance Estimators

Using the methods of Section 7.2 we can show that the estimators  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{e}_i^2$  and  $s^2 = \frac{1}{n-k} \sum_{i=1}^n \hat{e}_i^2$  are consistent for  $\sigma^2$ .

The trick is to write the residual  $\hat{e}_i$  as equal to the error  $e_i$  plus a deviation term

$$\begin{aligned} \hat{e}_i &= y_i - \mathbf{x}_i' \hat{\beta} \\ &= e_i + \mathbf{x}_i' \beta - \mathbf{x}_i' \hat{\beta} \\ &= e_i - \mathbf{x}_i' (\hat{\beta} - \beta). \end{aligned}$$

Thus the squared residual equals the squared error plus a deviation

$$\hat{e}_i^2 = e_i^2 - 2e_i \mathbf{x}_i' (\hat{\beta} - \beta) + (\hat{\beta} - \beta)' \mathbf{x}_i \mathbf{x}_i' (\hat{\beta} - \beta). \quad (7.23)$$

So when we take the average of the squared residuals we obtain the average of the squared errors, plus two terms which are (hopefully) asymptotically negligible.

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n e_i^2 - 2 \left( \frac{1}{n} \sum_{i=1}^n e_i \mathbf{x}_i' \right) (\hat{\beta} - \beta) \\ &\quad + (\hat{\beta} - \beta)' \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right) (\hat{\beta} - \beta). \end{aligned} \quad (7.24)$$

Indeed, the WLLN shows that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n e_i^2 &\xrightarrow{p} \sigma^2 \\ \frac{1}{n} \sum_{i=1}^n e_i \mathbf{x}_i' &\xrightarrow{p} \mathbb{E}(e_i \mathbf{x}_i') = 0 \\ \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' &\xrightarrow{p} \mathbb{E}(\mathbf{x}_i \mathbf{x}_i') = \mathbf{Q}_{xx} \end{aligned}$$

and Theorem 7.2.1 shows that  $\hat{\beta} \xrightarrow{p} \beta$ . Hence (7.24) converges in probability to  $\sigma^2$ , as desired.

Finally, since  $n/(n-k) \rightarrow 1$  as  $n \rightarrow \infty$ , it follows that

$$s^2 = \left( \frac{n}{n-k} \right) \hat{\sigma}^2 \xrightarrow{p} \sigma^2.$$

Thus both estimators are consistent.

**Theorem 7.5.1** *Under Assumption 7.1.1,  $\hat{\sigma}^2 \xrightarrow{p} \sigma^2$  and  $s^2 \xrightarrow{p} \sigma^2$  as  $n \rightarrow \infty$ .*

## 7.6 Homoskedastic Covariance Matrix Estimation

Theorem 7.3.2 shows that  $\sqrt{n}(\hat{\beta} - \beta)$  is asymptotically normal with asymptotic covariance matrix  $\mathbf{V}_\beta$ . For asymptotic inference (confidence intervals and tests) we need a consistent estimate of  $\mathbf{V}_\beta$ . Under homoskedasticity,  $\mathbf{V}_\beta$  simplifies to  $\mathbf{V}_\beta^0 = \mathbf{Q}_{xx}^{-1}\sigma^2$ , and in this section we consider the simplified problem of estimating  $\mathbf{V}_\beta^0$ .

The standard moment estimator of  $\mathbf{Q}_{xx}$  is  $\hat{\mathbf{Q}}_{xx}$  defined in (7.1), and thus an estimator for  $\mathbf{Q}_{xx}^{-1}$  is  $\hat{\mathbf{Q}}_{xx}^{-1}$ . Also, the standard estimator of  $\sigma^2$  is the unbiased estimator  $s^2$  defined in (4.30). Thus a natural plug-in estimator for  $\mathbf{V}_\beta^0 = \mathbf{Q}_{xx}^{-1}\sigma^2$  is  $\hat{\mathbf{V}}_\beta^0 = \hat{\mathbf{Q}}_{xx}^{-1}s^2$ .

Consistency of  $\hat{\mathbf{V}}_\beta^0$  for  $\mathbf{V}_\beta^0$  follows from consistency of the moment estimates  $\hat{\mathbf{Q}}_{xx}$  and  $s^2$ , and an application of the continuous mapping theorem. Specifically, Theorem 7.2.1 established that  $\hat{\mathbf{Q}}_{xx} \xrightarrow{p} \mathbf{Q}_{xx}$ , and Theorem 7.5.1 established  $s^2 \xrightarrow{p} \sigma^2$ . The function  $\mathbf{V}_\beta^0 = \mathbf{Q}_{xx}^{-1}\sigma^2$  is a continuous function of  $\mathbf{Q}_{xx}$  and  $\sigma^2$  so long as  $\mathbf{Q}_{xx} > 0$ , which holds true under Assumption 7.1.1.4. It follows by the CMT that

$$\hat{\mathbf{V}}_\beta^0 = \hat{\mathbf{Q}}_{xx}^{-1}s^2 \xrightarrow{p} \mathbf{Q}_{xx}^{-1}\sigma^2 = \mathbf{V}_\beta^0$$

so that  $\hat{\mathbf{V}}_\beta^0$  is consistent for  $\mathbf{V}_\beta^0$ , as desired.

**Theorem 7.6.1** *Under Assumption 7.1.1,  $\hat{\mathbf{V}}_\beta^0 \xrightarrow{p} \mathbf{V}_\beta^0$  as  $n \rightarrow \infty$ .*

It is instructive to notice that Theorem 7.6.1 does not require the assumption of homoskedasticity. That is,  $\hat{\mathbf{V}}_\beta^0$  is consistent for  $\mathbf{V}_\beta^0$  regardless if the regression is homoskedastic or heteroskedastic. However,  $\mathbf{V}_\beta^0 = \mathbf{V}_\beta = \text{avar}(\hat{\beta})$  only under homoskedasticity. Thus in the general case,  $\hat{\mathbf{V}}_\beta^0$  is consistent for a well-defined but non-useful object.

## 7.7 Heteroskedastic Covariance Matrix Estimation

Theorems 7.3.2 established that the asymptotic covariance matrix of  $\sqrt{n}(\hat{\beta} - \beta)$  is  $\mathbf{V}_\beta = \mathbf{Q}_{xx}^{-1}\mathbf{\Omega}\mathbf{Q}_{xx}^{-1}$ . We now consider estimation of this covariance matrix without imposing homoskedasticity. The standard approach is to use a plug-in estimator which replaces the unknowns with sample moments.

As described in the previous section, a natural estimator for  $\mathbf{Q}_{xx}^{-1}$  is  $\hat{\mathbf{Q}}_{xx}^{-1}$ , where  $\hat{\mathbf{Q}}_{xx}$  defined in (7.1).

The moment estimator for  $\mathbf{\Omega}$  is

$$\hat{\mathbf{\Omega}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \hat{e}_i^2, \quad (7.25)$$

leading to the plug-in covariance matrix estimator

$$\hat{\mathbf{V}}_{\beta}^W = \hat{\mathbf{Q}}_{xx}^{-1} \hat{\mathbf{\Omega}} \hat{\mathbf{Q}}_{xx}^{-1}. \quad (7.26)$$

You can check that  $\hat{\mathbf{V}}_{\beta}^W = n \hat{\mathbf{V}}_{\hat{\beta}}^W$  where  $\hat{\mathbf{V}}_{\hat{\beta}}^W$  is the White covariance matrix estimator introduced in (4.37).

As shown in Theorem 7.2.1,  $\hat{\mathbf{Q}}_{xx}^{-1} \xrightarrow{p} \mathbf{Q}_{xx}^{-1}$ , so we just need to verify the consistency of  $\hat{\mathbf{\Omega}}$ . The key is to replace the squared residual  $\hat{e}_i^2$  with the squared error  $e_i^2$ , and then show that the difference is asymptotically negligible.

Specifically, observe that

$$\begin{aligned} \hat{\mathbf{\Omega}} &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \hat{e}_i^2 \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' e_i^2 + \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' (\hat{e}_i^2 - e_i^2). \end{aligned} \quad (7.27)$$

The first term is an average of the iid random variables  $\mathbf{x}_i \mathbf{x}_i' e_i^2$ , and therefore by the WLLN converges in probability to its expectation, namely,

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' e_i^2 \xrightarrow{p} \mathbb{E}(\mathbf{x}_i \mathbf{x}_i' e_i^2) = \mathbf{\Omega}.$$

Technically, this requires that  $\mathbf{\Omega}$  has finite elements, which was shown in (7.9).

So to establish that  $\hat{\mathbf{\Omega}}$  is consistent for  $\mathbf{\Omega}$  it remains to show that

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' (\hat{e}_i^2 - e_i^2) \xrightarrow{p} 0. \quad (7.28)$$

There are multiple ways to do this. A reasonable straightforward yet slightly tedious derivation is to start by applying the Triangle Inequality (A.26) using a matrix norm:

$$\begin{aligned} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' (\hat{e}_i^2 - e_i^2) \right\| &\leq \frac{1}{n} \sum_{i=1}^n \left\| \mathbf{x}_i \mathbf{x}_i' (\hat{e}_i^2 - e_i^2) \right\| \\ &= \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2 |\hat{e}_i^2 - e_i^2|. \end{aligned} \quad (7.29)$$

Then recalling the expression for the squared residual (7.23), apply the Triangle Inequality and then the Schwarz Inequality (A.20) twice

$$\begin{aligned} |\hat{e}_i^2 - e_i^2| &\leq 2 \left| e_i \mathbf{x}_i' (\hat{\beta} - \beta) \right| + \left| (\hat{\beta} - \beta)' \mathbf{x}_i \mathbf{x}_i' (\hat{\beta} - \beta) \right| \\ &= 2 |e_i| \left| \mathbf{x}_i' (\hat{\beta} - \beta) \right| + \left| (\hat{\beta} - \beta)' \mathbf{x}_i \right|^2 \\ &\leq 2 |e_i| \|\mathbf{x}_i\| \|\hat{\beta} - \beta\| + \|\mathbf{x}_i\|^2 \|\hat{\beta} - \beta\|^2. \end{aligned} \quad (7.30)$$

Combining (7.29) and (7.30), we find

$$\begin{aligned} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' (\hat{e}_i^2 - e_i^2) \right\| &\leq 2 \left( \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^3 |e_i| \right) \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\| \\ &\quad + \left( \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^4 \right) \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|^2 \\ &= o_p(1). \end{aligned} \tag{7.31}$$

The expression is  $o_p(1)$  because  $\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\| \xrightarrow{p} 0$  and both averages in parenthesis are averages of random variables with finite mean under Assumption 7.1.2 (and are thus  $O_p(1)$ ). Indeed, by Hölder's Inequality (B.9)

$$\begin{aligned} \mathbb{E} \left( \|\mathbf{x}_i\|^3 |e_i| \right) &\leq \left( \mathbb{E} \left( \|\mathbf{x}_i\|^3 \right)^{4/3} \right)^{3/4} \left( \mathbb{E} (e_i^4) \right)^{1/4} \\ &= \left( \mathbb{E} \left( \|\mathbf{x}_i\|^4 \right) \right)^{3/4} \left( \mathbb{E} (e_i^4) \right)^{1/4} < \infty. \end{aligned}$$

We have established (7.28), as desired.

**Theorem 7.7.1** *Under Assumption 7.1.2, as  $n \rightarrow \infty$ ,  $\hat{\boldsymbol{\Omega}} \xrightarrow{p} \boldsymbol{\Omega}$  and  $\hat{\mathbf{V}}_{\boldsymbol{\beta}}^W \xrightarrow{p} \mathbf{V}_{\boldsymbol{\beta}}$ .*

For an alternative proof of this result, see Section 7.21.

## 7.8 Summary of Covariance Matrix Notation

The notation we have introduced may be somewhat confusing so it is helpful to write it down in one place. The exact variance of  $\hat{\boldsymbol{\beta}}$  (under the assumptions of the linear regression model) and the asymptotic variance of  $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$  (under the more general assumptions of the linear projection model) are

$$\begin{aligned} \mathbf{V}_{\hat{\boldsymbol{\beta}}} &= \text{var}(\hat{\boldsymbol{\beta}} | \mathbf{X}) = (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{D}\mathbf{X}) (\mathbf{X}'\mathbf{X})^{-1} \\ \mathbf{V}_{\boldsymbol{\beta}} &= \text{avar}(\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})) = \mathbf{Q}_{xx}^{-1} \boldsymbol{\Omega} \mathbf{Q}_{xx}^{-1}. \end{aligned}$$

The White estimates of these two covariance matrices are

$$\begin{aligned} \hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}^W &= (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \hat{e}_i^2 \right) (\mathbf{X}'\mathbf{X})^{-1} \\ \hat{\mathbf{V}}_{\boldsymbol{\beta}}^W &= \hat{\mathbf{Q}}_{xx}^{-1} \hat{\boldsymbol{\Omega}} \hat{\mathbf{Q}}_{xx}^{-1} \end{aligned}$$

and satisfy the simple relationship

$$\hat{\mathbf{V}}_{\boldsymbol{\beta}}^W = n \hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}^W.$$

Similarly, under the assumption of homoskedasticity the exact and asymptotic variances simplify to

$$\begin{aligned} \mathbf{V}_{\hat{\boldsymbol{\beta}}}^0 &= (\mathbf{X}'\mathbf{X})^{-1} \sigma^2 \\ \mathbf{V}_{\boldsymbol{\beta}}^0 &= \mathbf{Q}_{xx}^{-1} \sigma^2 \end{aligned}$$

and their standard estimators are

$$\begin{aligned}\widehat{\mathbf{V}}_{\beta}^0 &= (\mathbf{X}'\mathbf{X})^{-1} s^2 \\ \widehat{\mathbf{V}}_{\beta}^0 &= \widehat{\mathbf{Q}}_{\mathbf{x}\mathbf{x}}^{-1} s^2\end{aligned}$$

which also satisfy the relationship

$$\widehat{\mathbf{V}}_{\beta}^0 = n \widehat{\mathbf{V}}_{\beta}^0.$$

The exact formula and estimates are useful when constructing test statistics and standard errors. However, for theoretical purposes the asymptotic formula (variances and their estimates) are more useful, as these retain non-degenerate limits as the sample sizes diverge. That is why both sets of notation are useful.

## 7.9 Alternative Covariance Matrix Estimators\*

In Section 7.7 we introduced  $\widehat{\mathbf{V}}_{\beta}^W$  as an estimator of  $\mathbf{V}_{\beta}$ .  $\widehat{\mathbf{V}}_{\beta}^W$  is a scaled version of  $\widehat{\mathbf{V}}_{\beta}^W$  from Section 4.13, where we also introduced the alternative heteroskedasticity-robust covariance matrix estimators  $\widehat{\mathbf{V}}_{\beta}$ ,  $\widetilde{\mathbf{V}}_{\beta}$  and  $\overline{\mathbf{V}}_{\beta}$ . We now discuss the consistency properties of these estimators.

To do so we introduce their scaled versions, e.g.  $\widehat{\mathbf{V}}_{\beta} = n \widehat{\mathbf{V}}_{\beta}$ ,  $\widetilde{\mathbf{V}}_{\beta} = n \widetilde{\mathbf{V}}_{\beta}$ , and  $\overline{\mathbf{V}}_{\beta} = n \overline{\mathbf{V}}_{\beta}$ . These are (alternative) estimates of the asymptotic covariance matrix  $\mathbf{V}_{\beta}$ .

First, consider  $\widehat{\mathbf{V}}_{\beta}$ . Notice that  $\widehat{\mathbf{V}}_{\beta} = n \widehat{\mathbf{V}}_{\beta} = \frac{n}{n-k} \widehat{\mathbf{V}}_{\beta}^W$  where  $\widehat{\mathbf{V}}_{\beta}^W$  was defined in (7.26) and shown consistent for  $\mathbf{V}_{\beta}$  in Theorem 7.7.1. If  $k$  is fixed as  $n \rightarrow \infty$ , then  $\frac{n}{n-k} \rightarrow 1$  and thus

$$\widehat{\mathbf{V}}_{\beta} = (1 + o(1)) \widehat{\mathbf{V}}_{\beta}^W \xrightarrow{p} \mathbf{V}_{\beta}.$$

Thus  $\widehat{\mathbf{V}}_{\beta}$  is consistent for  $\mathbf{V}_{\beta}$ .

The alternative estimators  $\widetilde{\mathbf{V}}_{\beta}$  and  $\overline{\mathbf{V}}_{\beta}$  take the form (7.26) but with  $\widehat{\mathbf{\Omega}}$  replaced by

$$\widetilde{\mathbf{\Omega}} = \frac{1}{n} \sum_{i=1}^n (1 - h_{ii})^{-2} \mathbf{x}_i \mathbf{x}_i' \widehat{e}_i^2$$

and

$$\overline{\mathbf{\Omega}} = \frac{1}{n} \sum_{i=1}^n (1 - h_{ii})^{-1} \mathbf{x}_i \mathbf{x}_i' \widehat{e}_i^2,$$

respectively. To show that these estimators also consistent for  $\mathbf{V}_{\beta}$ , given  $\widehat{\mathbf{\Omega}} \xrightarrow{p} \mathbf{\Omega}$ , it is sufficient to show that the differences  $\widetilde{\mathbf{\Omega}} - \widehat{\mathbf{\Omega}}$  and  $\overline{\mathbf{\Omega}} - \widehat{\mathbf{\Omega}}$  converge in probability to zero as  $n \rightarrow \infty$ .

The trick is to use the fact that the leverage values are asymptotically negligible:

$$h_n^* = \max_{1 \leq i \leq n} h_{ii} = o_p(1). \quad (7.32)$$

(See Theorem 7.22.1 in Section 7.22.) Then using the Triangle Inequality

$$\begin{aligned}\|\overline{\mathbf{\Omega}} - \widehat{\mathbf{\Omega}}\| &\leq \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i \mathbf{x}_i'\| \widehat{e}_i^2 \left| (1 - h_{ii})^{-1} - 1 \right| \\ &\leq \left( \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2 \widehat{e}_i^2 \right) \left| (1 - h_n^*)^{-1} - 1 \right|.\end{aligned}$$

The sum in parenthesis can be shown to be  $O_p(1)$  under Assumption 7.1.2 by the same argument as in the proof of Theorem 7.7.1. (In fact, it can be shown to converge in probability to



$\mathbb{E} \left( \|\mathbf{x}_i\|^2 e_i^2 \right)$ . The term in absolute values is  $o_p(1)$  by (7.32). Thus the product is  $o_p(1)$ , which means that  $\bar{\boldsymbol{\Omega}} = \hat{\boldsymbol{\Omega}} + o_p(1) \longrightarrow \boldsymbol{\Omega}$ .

Similarly,

$$\begin{aligned} \|\tilde{\boldsymbol{\Omega}} - \hat{\boldsymbol{\Omega}}\| &\leq \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i \mathbf{x}_i'\| \hat{e}_i^2 \left| (1 - h_{ii})^{-2} - 1 \right| \\ &\leq \left( \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2 \hat{e}_i^2 \right) \left| (1 - h_n^*)^{-2} - 1 \right| \\ &= o_p(1). \end{aligned}$$

**Theorem 7.9.1** *Under Assumption 7.1.2, as  $n \rightarrow \infty$ ,  $\tilde{\boldsymbol{\Omega}} \xrightarrow{p} \boldsymbol{\Omega}$ ,  $\bar{\boldsymbol{\Omega}} \xrightarrow{p} \boldsymbol{\Omega}$ ,  $\hat{\mathbf{V}}_{\boldsymbol{\beta}} \xrightarrow{p} \mathbf{V}_{\boldsymbol{\beta}}$ ,  $\tilde{\mathbf{V}}_{\boldsymbol{\beta}} \xrightarrow{p} \mathbf{V}_{\boldsymbol{\beta}}$ , and  $\bar{\mathbf{V}}_{\boldsymbol{\beta}} \xrightarrow{p} \mathbf{V}_{\boldsymbol{\beta}}$ .*

Theorem 7.9.1 shows that the alternative covariance matrix estimators are also consistent for the asymptotic covariance matrix.

## 7.10 Functions of Parameters

In most serious applications the researcher is actually interested in a specific transformation of the coefficient vector  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_k)$ . For example, he or she may be interested in a single coefficient  $\beta_j$ , or a ratio  $\beta_j/\beta_l$ . More generally, interest may focus on a quantity such as consumer surplus which could be a complicated function of the coefficients. In any of these cases we can write the parameter of interest  $\boldsymbol{\theta}$  as a function of the coefficients, e.g.  $\boldsymbol{\theta} = \mathbf{r}(\boldsymbol{\beta})$  for some function  $\mathbf{r} : \mathbb{R}^k \rightarrow \mathbb{R}^q$ . The estimate of  $\boldsymbol{\theta}$  is

$$\hat{\boldsymbol{\theta}} = \mathbf{r}(\hat{\boldsymbol{\beta}}).$$

By the continuous mapping theorem (Theorem 6.11.1) and the fact  $\hat{\boldsymbol{\beta}} \xrightarrow{p} \boldsymbol{\beta}$  we can deduce that  $\hat{\boldsymbol{\theta}}$  is consistent for  $\boldsymbol{\theta}$  (if the function  $\mathbf{r}(\cdot)$  is continuous).

**Theorem 7.10.1** *Under Assumption 7.1.1, if  $\mathbf{r}(\boldsymbol{\beta})$  is continuous at the true value of  $\boldsymbol{\beta}$ , then as  $n \rightarrow \infty$ ,  $\hat{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}$ .*

Furthermore, if the transformation is sufficiently smooth, by the Delta Method (Theorem 6.12.3) we can show that  $\hat{\boldsymbol{\theta}}$  is asymptotically normal.

**Assumption 7.10.1**  $\mathbf{r}(\boldsymbol{\beta}) : \mathbb{R}^k \rightarrow \mathbb{R}^q$  is continuously differentiable at the true value of  $\boldsymbol{\beta}$  and  $\mathbf{R} = \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{r}(\boldsymbol{\beta})'$  has rank  $q$ .

**Theorem 7.10.2 Asymptotic Distribution of Functions of Parameters**

Under Assumptions 7.1.2 and 7.10.1, as  $n \rightarrow \infty$ ,

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_{\boldsymbol{\theta}}) \quad (7.33)$$

where

$$\mathbf{V}_{\boldsymbol{\theta}} = \mathbf{R}' \mathbf{V}_{\boldsymbol{\beta}} \mathbf{R} \quad (7.34)$$

In many cases, the function  $\mathbf{r}(\boldsymbol{\beta})$  is linear:

$$\mathbf{r}(\boldsymbol{\beta}) = \mathbf{R}' \boldsymbol{\beta}$$

for some  $k \times q$  matrix  $\mathbf{R}$ . In particular, if  $\mathbf{R}$  is a “selector matrix”

$$\mathbf{R} = \begin{pmatrix} \mathbf{I} \\ \mathbf{0} \end{pmatrix} \quad (7.35)$$

then we can partition  $\boldsymbol{\beta} = (\boldsymbol{\beta}'_1, \boldsymbol{\beta}'_2)'$  so that  $\mathbf{R}' \boldsymbol{\beta} = \boldsymbol{\beta}_1$  for  $\boldsymbol{\beta} = (\boldsymbol{\beta}'_1, \boldsymbol{\beta}'_2)'$ . Then

$$\mathbf{V}_{\boldsymbol{\theta}} = \begin{pmatrix} \mathbf{I} & \mathbf{0} \end{pmatrix} \mathbf{V}_{\boldsymbol{\beta}} \begin{pmatrix} \mathbf{I} \\ \mathbf{0} \end{pmatrix} = \mathbf{V}_{11},$$

the upper-left sub-matrix of  $\mathbf{V}_{11}$  given in (7.20). In this case (7.33) states that

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_1) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_{11}).$$

That is, subsets of  $\hat{\boldsymbol{\beta}}$  are approximately normal with variances given by the conformable subcomponents of  $\mathbf{V}$ .

To illustrate the case of a nonlinear transformation, take the example  $\theta = \beta_j/\beta_l$  for  $j \neq l$ . Then

$$\mathbf{R} = \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{r}(\boldsymbol{\beta}) = \begin{pmatrix} \frac{\partial}{\partial \beta_1} (\beta_j/\beta_l) \\ \vdots \\ \frac{\partial}{\partial \beta_j} (\beta_j/\beta_l) \\ \vdots \\ \frac{\partial}{\partial \beta_l} (\beta_j/\beta_l) \\ \vdots \\ \frac{\partial}{\partial \beta_k} (\beta_j/\beta_l) \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 1/\beta_l \\ \vdots \\ -\beta_j/\beta_l^2 \\ \vdots \\ 0 \end{pmatrix} \quad (7.36)$$

so

$$\mathbf{V}_{\boldsymbol{\theta}} = \mathbf{V}_{jj}/\beta_l^2 + \mathbf{V}_{ll}\beta_j^2/\beta_l^4 - 2\mathbf{V}_{jl}\beta_j/\beta_l^3$$

where  $\mathbf{V}_{ab}$  denotes the  $ab^{th}$  element of  $\mathbf{V}_{\boldsymbol{\beta}}$ .

For inference we need an estimate of the asymptotic variance matrix  $\mathbf{V}_{\boldsymbol{\theta}} = \mathbf{R}' \mathbf{V}_{\boldsymbol{\beta}} \mathbf{R}$ , and for this it is typical to use a plug-in estimator. The natural estimator of  $\mathbf{R}$  is the derivative evaluated at the point estimates

$$\hat{\mathbf{R}} = \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{r}(\hat{\boldsymbol{\beta}})'. \quad (7.37)$$

The derivative in (7.37) may be calculated analytically or numerically. By analytically, we mean working out for the formula for the derivative and replacing the unknowns by point estimates. For

example, if  $\theta = \beta_j/\beta_l$ , then  $\frac{\partial}{\partial \beta} \mathbf{r}(\beta)$  is (7.36). However in some cases the function  $\mathbf{r}(\beta)$  may be extremely complicated and a formula for the analytic derivative may not be easily available. In this case calculation by numerical differentiation may be preferable. Let  $\delta_l = (0 \cdots 1 \cdots 0)'$  be the unit vector with the “1” in the  $l^{\text{th}}$  place. Then the  $jl^{\text{th}}$  element of a numerical derivative  $\hat{\mathbf{R}}$  is

$$\hat{\mathbf{R}}_{jl} = \frac{\mathbf{r}_j(\hat{\beta} + \delta_l \varepsilon) - \mathbf{r}_j(\hat{\beta})}{\varepsilon}$$

for some small  $\varepsilon$ .

The estimate of  $\mathbf{V}_\theta$  is

$$\hat{\mathbf{V}}_\theta = \hat{\mathbf{R}}' \hat{\mathbf{V}}_\beta \hat{\mathbf{R}}. \quad (7.38)$$

Alternatively,  $\hat{\mathbf{V}}_\beta^0$ ,  $\tilde{\mathbf{V}}_\beta$  or  $\overline{\mathbf{V}}_\beta$  may be used in place of  $\hat{\mathbf{V}}_\beta$ . For example, the homoskedastic covariance matrix estimator is

$$\hat{\mathbf{V}}_\theta^0 = \hat{\mathbf{R}}' \hat{\mathbf{V}}_\beta^0 \hat{\mathbf{R}} = \hat{\mathbf{R}}' \hat{\mathbf{Q}}_{xx}^{-1} \hat{\mathbf{R}} s^2 \quad (7.39)$$

Given (7.37), (7.38) and (7.39) are simple to calculate using matrix operations.

As the primary justification for  $\hat{\mathbf{V}}_\theta$  is the asymptotic approximation (7.33),  $\hat{\mathbf{V}}_\theta$  is often called an **asymptotic covariance matrix estimator**.

The estimator  $\hat{\mathbf{V}}_\theta$  is consistent for  $\mathbf{V}_\theta$  under the conditions of Theorem 7.10.2 since  $\hat{\mathbf{V}}_\beta \xrightarrow{p} \mathbf{V}_\beta$  by Theorem 7.7.1, and

$$\hat{\mathbf{R}} = \frac{\partial}{\partial \beta} \mathbf{r}(\hat{\beta})' \xrightarrow{p} \frac{\partial}{\partial \beta} \mathbf{r}(\beta)' = \mathbf{R}$$

since  $\hat{\beta} \xrightarrow{p} \beta$  and the function  $\frac{\partial}{\partial \beta} \mathbf{r}(\beta)'$  is continuous in  $\beta$ .

**Theorem 7.10.3** Under Assumptions 7.1.2 and 7.10.1, as  $n \rightarrow \infty$ ,

$$\hat{\mathbf{V}}_\theta \xrightarrow{p} \mathbf{V}_\theta.$$

Theorem 7.10.3 shows that  $\hat{\mathbf{V}}_\theta$  is consistent for  $\mathbf{V}_\theta$  and thus may be used for asymptotic inference. In practice, we may set

$$\hat{\mathbf{V}}_{\hat{\theta}} = \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\hat{\beta}} \hat{\mathbf{R}} = n^{-1} \hat{\mathbf{R}}' \hat{\mathbf{V}}_\beta \hat{\mathbf{R}} \quad (7.40)$$

as an estimate of the variance of  $\hat{\theta}$ , or substitute an alternative covariance estimator such as  $\overline{\mathbf{V}}_{\hat{\beta}}$ .

## 7.11 Asymptotic Standard Errors

As described in Section 4.14, a standard error is an estimate of the standard deviation of the distribution of an estimator. Thus if  $\hat{\mathbf{V}}_{\hat{\beta}}$  is an estimate of the covariance matrix of  $\hat{\beta}$ , then standard errors are the square roots of the diagonal elements of this matrix. These take the form

$$s(\hat{\beta}_j) = \sqrt{\hat{\mathbf{V}}_{\hat{\beta}_j}} = \sqrt{[\hat{\mathbf{V}}_{\hat{\beta}}]_{jj}}.$$

Standard errors for  $\hat{\theta}$  are constructed similarly. Supposing that  $q = 1$  (so  $h(\beta)$  is real-valued), then the standard error for  $\hat{\theta}$  is the square root of (7.40)

$$s(\hat{\theta}) = \sqrt{\hat{\mathbf{R}}' \hat{\mathbf{V}}_{\hat{\beta}} \hat{\mathbf{R}}} = \sqrt{n^{-1} \hat{\mathbf{R}}' \hat{\mathbf{V}}_\beta \hat{\mathbf{R}}}.$$

When the justification is based on asymptotic theory we call  $s(\hat{\beta}_j)$  or  $s(\hat{\theta})$  an **asymptotic standard error** for  $\hat{\beta}_j$  or  $\hat{\theta}$ . When reporting your results, it is good practice to report standard errors for each reported estimate, and this includes functions and transformations of your parameter estimates. This helps users of the work (including yourself) assess the estimation precision.

We illustrate using the log wage regression

$$\log(\text{Wage}) = \beta_1 \text{ education} + \beta_2 \text{ experience} + \beta_3 \text{ experience}^2/100 + \beta_4 + e.$$

Consider the following three parameters of interest.

1. Percentage return to education:

$$\theta_1 = 100\beta_1$$

(100 times the partial derivative of the conditional expectation of log wages with respect to *education*.)

2. Percentage return to experience for individuals with 10 years of experience:

$$\theta_2 = 100\beta_2 + 20\beta_3$$

(100 times the partial derivative of the conditional expectation of log wages with respect to *experience*, evaluated at *experience* = 10.)

3. Experience level which maximizes expected log wages:

$$\theta_3 = -50\beta_2/\beta_3$$

(The level of *experience* at which the partial derivative of the conditional expectation of log wages with respect to *experience* equals 0.)

The  $4 \times 1$  vector  $\mathbf{R}$  for these three parameters is

$$\mathbf{R} = \begin{pmatrix} 100 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \quad \begin{pmatrix} 0 \\ 100 \\ 20 \\ 0 \end{pmatrix}, \quad \begin{pmatrix} 0 \\ -50/\beta_3 \\ 50\beta_2/\beta_3^2 \\ 0 \end{pmatrix},$$

respectively.

We use the subsample of married black women (all experience levels), which has 982 observations. The point estimates and standard errors are

$$\begin{aligned} \log(\widehat{\text{Wage}}) = & \begin{array}{cccc} 0.118 & 0.016 & -0.022 & 0.947 \\ (0.008) & (0.006) & (0.012) & (0.157) \end{array} \cdot \text{education} + \text{experience} - \text{experience}^2/100 + \text{constant} \\ & (7.41) \end{aligned}$$

The standard errors are the square roots of the Horn-Horn-Duncan covariance matrix estimate

$$\overline{\mathbf{V}}_{\hat{\beta}} = \begin{pmatrix} 0.632 & 0.131 & -0.143 & -11.1 \\ 0.131 & 0.390 & -0.731 & -6.25 \\ -0.143 & -0.731 & 1.48 & 9.43 \\ -11.1 & -6.25 & 9.43 & 246 \end{pmatrix} \times 10^{-4}. \quad (7.42)$$

We calculate that

$$\begin{aligned} \hat{\theta}_1 &= 100\hat{\beta}_1 \\ &= 100 \times 0.118 \\ &= 11.8 \end{aligned}$$

$$\begin{aligned} s(\hat{\theta}_1) &= \sqrt{100^2 \times 0.632 \times 10^{-4}} \\ &= 0.8 \end{aligned}$$

$$\begin{aligned} \hat{\theta}_2 &= 100\hat{\beta}_2 + 20\hat{\beta}_3 \\ &= 100 \times 0.016 - 20 \times 0.022 \\ &= 1.16 \end{aligned}$$

$$\begin{aligned} s(\hat{\theta}_2) &= \sqrt{\begin{pmatrix} 100 & 20 \end{pmatrix} \begin{pmatrix} 0.390 & -0.731 \\ -0.731 & 1.48 \end{pmatrix} \begin{pmatrix} 100 \\ 20 \end{pmatrix} \times 10^{-4}} \\ &= 0.55 \end{aligned}$$

$$\begin{aligned} \hat{\theta}_3 &= -50\hat{\beta}_2/\hat{\beta}_3 \\ &= 50 \times 0.016/0.022 \\ &= 35.2 \end{aligned}$$

$$\begin{aligned} s(\hat{\theta}_3) &= \sqrt{\begin{pmatrix} -50/\hat{\beta}_3 & 50\hat{\beta}_2/\hat{\beta}_3^2 \end{pmatrix} \begin{pmatrix} 0.390 & -0.731 \\ -0.731 & 1.48 \end{pmatrix} \begin{pmatrix} -50/\hat{\beta}_3 \\ 50\hat{\beta}_2/\hat{\beta}_3^2 \end{pmatrix} \times 10^{-4}} \\ &= 7.0. \end{aligned}$$

The calculations show that the estimate of the percentage return to education (for married black women) is about 12% per year, with a standard error of 0.8. The estimate of the percentage return to experience for those with 10 years of experience is 1.2% per year, with a standard error of 0.6. And the estimate of the experience level which maximizes expected log wages is 35 years, with a standard error of 7.

## 7.12 t-statistic

Let  $\theta = r(\beta) : \mathbb{R}^k \rightarrow \mathbb{R}$  be a parameter of interest,  $\hat{\theta}$  its estimate and  $s(\hat{\theta})$  its asymptotic standard error. Consider the statistic

$$T(\theta) = \frac{\hat{\theta} - \theta}{s(\hat{\theta})}. \quad (7.43)$$

Different writers have called (7.43) a **t-statistic**, a **t-ratio**, a **z-statistic** or a **studentized statistic**, sometimes using the different labels to distinguish between finite-sample and asymptotic inference. As the statistics themselves are always (7.43) we won't make this distinction, and will simply refer to  $T(\theta)$  as a t-statistic or a t-ratio. We also often suppress the parameter dependence, writing it as  $T$ . The t-statistic is a simple function of the estimate, its standard error, and the parameter.

By Theorems 7.10.2 and 7.10.3,  $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, V_\theta)$  and  $\hat{V}_\theta \xrightarrow{p} V_\theta$ . Thus

$$\begin{aligned} T(\theta) &= \frac{\hat{\theta} - \theta}{s(\hat{\theta})} \\ &= \frac{\sqrt{n}(\hat{\theta} - \theta)}{\sqrt{\hat{V}_\theta}} \\ &\xrightarrow{d} \frac{N(0, V_\theta)}{\sqrt{V_\theta}} \\ &= Z \sim N(0, 1). \end{aligned}$$

The last equality is by the property that affine functions of normal distributions are normal (Theorem 5.2.3).

Thus the asymptotic distribution of the t-ratio  $T(\theta)$  is the standard normal. Since this distribution does not depend on the parameters, we say that  $T(\theta)$  is **asymptotically pivotal**. In finite samples  $T(\theta)$  is not necessarily pivotal (as in the normal regression model) but the property states that the dependence on unknowns diminishes as  $n$  increases.

As we will see in the next section, it is also useful to consider the distribution of the **absolute t-ratio**  $|T(\theta)|$ . Since  $T(\theta) \xrightarrow{d} Z$ , the continuous mapping theorem yields  $|T(\theta)| \xrightarrow{d} |Z|$ . Letting  $\Phi(u) = \Pr(Z \leq u)$  denote the standard normal distribution function, we can calculate that the distribution function of  $|Z|$  is

$$\begin{aligned} \Pr(|Z| \leq u) &= \Pr(-u \leq Z \leq u) \\ &= \Pr(Z \leq u) - \Pr(Z < -u) \\ &= \Phi(u) - \Phi(-u) \\ &= 2\Phi(u) - 1 \end{aligned} \tag{7.44}$$

**Theorem 7.12.1** Under Assumptions 7.1.2 and 7.10.1,  $T(\theta) \xrightarrow{d} Z \sim N(0, 1)$  and  $|t_n(\theta)| \xrightarrow{d} |Z|$ .

The asymptotic normality of Theorem 7.12.1 is used to justify confidence intervals and tests for the parameters.

## 7.13 Confidence Intervals

The estimate  $\hat{\theta}$  is a **point estimate** for  $\theta$ , meaning that  $\hat{\theta}$  is a single value in  $\mathbb{R}^q$ . A broader concept is a **set estimate**  $\hat{C}$  which is a collection of values in  $\mathbb{R}^q$ . When the parameter  $\theta$  is real-valued then it is common to focus on sets of the form  $\hat{C} = [\hat{L}, \hat{U}]$  which is called an **interval estimate** for  $\theta$ .

An interval estimate  $\hat{C}$  is a function of the data and hence is random. The **coverage probability** of the interval  $\hat{C} = [\hat{L}, \hat{U}]$  is  $\Pr(\theta \in \hat{C})$ . The randomness comes from  $\hat{C}$  as the parameter  $\theta$  is treated as fixed. In Section 5.12 we introduced confidence intervals for the normal regression model, which used the finite sample distribution of the t-statistic to construct exact confidence intervals for the regression coefficients. When we are outside the normal regression model we cannot rely on the exact normal distribution theory, but instead use asymptotic approximations. A benefit is that we can construct confidence intervals for general parameters of interest  $\theta$ , not just regression coefficients.

An interval estimate  $\hat{C}$  is called a **confidence interval** when the goal is to set the coverage probability to equal a pre-specified target such as 90% or 95%.  $\hat{C}$  is called a  $1 - \alpha$  confidence interval if  $\inf_{\theta} \Pr_{\theta}(\theta \in \hat{C}) = 1 - \alpha$ .

When  $\hat{\theta}$  is asymptotically normal with standard error  $s(\hat{\theta})$ , the conventional confidence interval for  $\theta$  takes the form

$$\hat{C} = [\hat{\theta} - c \cdot s(\hat{\theta}), \quad \hat{\theta} + c \cdot s(\hat{\theta})] \tag{7.45}$$

where  $c$  equals the  $1 - \alpha$  quantile of the distribution of  $|Z|$ . Using (7.44) we calculate that  $c$  is equivalently the  $1 - \alpha/2$  quantile of the standard normal distribution. Thus,  $c$  solves

$$2\Phi(c) - 1 = 1 - \alpha.$$

This can be computed by, for example, `norminv(1- $\alpha$ /2)` in MATLAB. The confidence interval (7.45) is symmetric about the point estimate  $\hat{\theta}$ , and its length is proportional to the standard error  $s(\hat{\theta})$ .

Equivalently, (7.45) is the set of parameter values for  $\theta$  such that the t-statistic  $T(\theta)$  is smaller (in absolute value) than  $c$ , that is

$$\hat{C} = \{\theta : |T(\theta)| \leq c\} = \left\{ \theta : -c \leq \frac{\hat{\theta} - \theta}{s(\hat{\theta})} \leq c \right\}.$$

The coverage probability of this confidence interval is

$$\Pr(\theta \in \hat{C}) = \Pr(|T(\theta)| \leq c) \rightarrow \Pr(|Z| \leq c) = 1 - \alpha$$

where the limit is taken as  $n \rightarrow \infty$ , and holds since  $T(\theta)$  is asymptotically  $|Z|$  by Theorem 7.12.1. We call the limit the **asymptotic coverage probability**, and call  $\hat{C}$  an asymptotic  $1 - \alpha\%$  confidence interval for  $\theta$ . Since the t-ratio is asymptotically pivotal, the asymptotic coverage probability is independent of the parameter  $\theta$ .

It is useful to contrast the confidence interval (7.45) with (5.12) for the normal regression model. They are similar, but there are differences. The normal regression interval (5.12) only applies to regression coefficients  $\beta$ , not to functions  $\theta$  of the coefficients. The normal interval (5.12) also is constructed with the homoskedastic standard error, while (7.45) can be constructed with a heteroskedastic-robust standard error. Furthermore, the constants  $c$  in (5.12) are calculated using the student  $t$  distribution, while  $c$  in (7.45) are calculated using the normal distribution. The difference between the student  $t$  and normal values are typically small in practice (since sample sizes are large in typical economic applications). However, since the student  $t$  values are larger, it results in slightly larger confidence intervals, which is probably reasonable. (A practical rule of thumb is that if the sample sizes are sufficiently small that it makes a difference, then probably neither (5.12) nor (7.45) should be trusted.) Despite these differences, the coincidence of the intervals means that inference on regression coefficients is generally robust to using either the exact normal sampling assumption or the asymptotic large sample approximation, at least in large samples.

In Stata, by default the program reports 95% confidence intervals for each coefficient where the critical values  $c$  are calculated using the  $t_{n-k}$  distribution. This is done for all standard error methods even though it is only justified for homoskedastic standard errors and under normality.

The standard coverage probability for confidence intervals is 95%, leading to the choice  $c = 1.96$  for the constant in (7.45). Rounding 1.96 to 2, we obtain what might be the most commonly used confidence interval in applied econometric practice

$$\hat{C} = \left[ \hat{\theta} - 2s(\hat{\theta}), \quad \hat{\theta} + 2s(\hat{\theta}) \right]. \quad (7.46)$$

This is a useful rule-of thumb. This asymptotic 95% confidence interval  $\hat{C}$  is simple to compute and can be roughly calculated from tables of coefficient estimates and standard errors. (Technically, it is an asymptotic 95.4% interval, due to the substitution of 2.0 for 1.96, but this distinction is overly precise.)

**Theorem 7.13.1** *Under Assumptions 7.1.2 and 7.10.1, for  $\hat{C}$  defined in (7.45), with  $c = \Phi^{-1}(1 - \alpha/2)$ ,  $\Pr(\theta \in \hat{C}) \rightarrow 1 - \alpha$ . For  $c = 1.96$ ,  $\Pr(\theta \in \hat{C}) \rightarrow 0.95$ .*

Confidence intervals are a simple yet effective tool to assess estimation uncertainty. When reading a set of empirical results, look at the estimated coefficient estimates and the standard errors. For a parameter of interest, compute the confidence interval  $C_n$  and consider the meaning of the spread of the suggested values. If the range of values in the confidence interval are too wide to learn about  $\theta$ , then do not jump to a conclusion about  $\theta$  based on the point estimate alone.

For illustration, consider the three examples presented in Section 7.11 based on the log wage regression for married black women.

Percentage return to education. A 95% asymptotic confidence interval is  $11.8 \pm 1.96 \times 0.8 = [10.2, 13.3]$ .

Percentage return to experience for individuals with 10 years experience. A 90% asymptotic confidence interval is  $1.1 \pm 1.645 \times 0.4 = [0.5, 1.8]$ .

Experience level which maximizes expected log wages. An 80% asymptotic confidence interval is  $35 \pm 1.28 \times 7 = [26, 44]$ .

## 7.14 Regression Intervals

In the linear regression model the conditional mean of  $y_i$  given  $\mathbf{x}_i = \mathbf{x}$  is

$$m(\mathbf{x}) = \mathbb{E}(y_i \mid \mathbf{x}_i = \mathbf{x}) = \mathbf{x}'\boldsymbol{\beta}.$$

In some cases, we want to estimate  $m(\mathbf{x})$  at a particular point  $\mathbf{x}$ . Notice that this is a linear function of  $\boldsymbol{\beta}$ . Letting  $r(\boldsymbol{\beta}) = \mathbf{x}'\boldsymbol{\beta}$  and  $\theta = r(\boldsymbol{\beta})$ , we see that  $\hat{m}(\mathbf{x}) = \hat{\theta} = \mathbf{x}'\hat{\boldsymbol{\beta}}$  and  $\mathbf{R} = \mathbf{x}$ , so  $s(\hat{\theta}) = \sqrt{\mathbf{x}'\hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}\mathbf{x}}$ . Thus an asymptotic 95% confidence interval for  $m(\mathbf{x})$  is

$$\left[ \mathbf{x}'\hat{\boldsymbol{\beta}} \pm 1.96\sqrt{\mathbf{x}'\hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}\mathbf{x}} \right].$$

It is interesting to observe that if this is viewed as a function of  $\mathbf{x}$ , the width of the confidence set is dependent on  $\mathbf{x}$ .

To illustrate, we return to the log wage regression (3.14) of Section 3.7. The estimated regression equation is

$$\log(\widehat{Wage}) = \mathbf{x}'\hat{\boldsymbol{\beta}} = 0.155x + 0.698$$

where  $x = \text{education}$ . The covariance matrix estimate from (4.44) is

$$\hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}} = \begin{pmatrix} 0.001 & -0.015 \\ -0.015 & 0.243 \end{pmatrix}.$$

Thus the 95% confidence interval for the regression takes the form

$$0.155x + 0.698 \pm 1.96\sqrt{0.001x^2 - 0.030x + 0.243}.$$

The estimated regression and 95% intervals are shown in Figure 7.6. Notice that the confidence bands take a hyperbolic shape. This means that the regression line is less precisely estimated for very large and very small values of *education*.

Plots of the estimated regression line and confidence intervals are especially useful when the regression includes nonlinear terms. To illustrate, consider the log wage regression (7.41) which includes experience and its square, with covariance matrix (7.42). We are interested in plotting the regression estimate and regression intervals as a function of *experience*. Since the regression also includes *education*, to plot the estimates in a simple graph we need to fix *education* at a



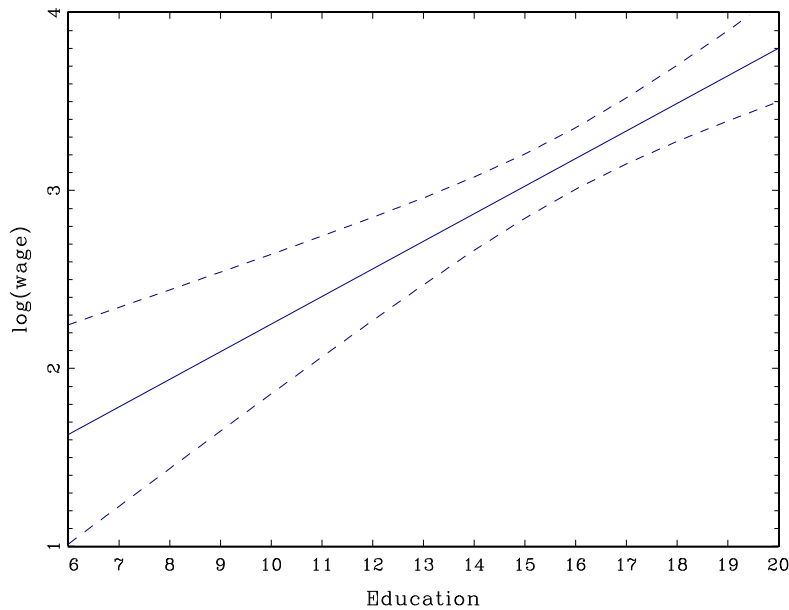


Figure 7.6: Wage on Education Regression Intervals

specific value. We select  $education=12$ . This only affects the level of the estimated regression, since  $education$  enters without an interaction. Define the points of evaluation

$$\mathbf{z}(x) = \begin{pmatrix} 12 \\ x \\ x^2/100 \\ 1 \end{pmatrix}$$

where  $x = experience$ .

Thus the 95% regression interval for  $education=12$ , as a function of  $x = experience$  is

$$\begin{aligned} & 0.118 \times 12 + 0.016 x - 0.022 x^2/100 + 0.947 \\ & \pm 1.96 \sqrt{\mathbf{z}(x)' \begin{pmatrix} 0.632 & 0.131 & -0.143 & -11.1 \\ 0.131 & 0.390 & -0.731 & -6.25 \\ -0.143 & -0.731 & 1.48 & 9.43 \\ -11.1 & -6.25 & 9.43 & 246 \end{pmatrix} \mathbf{z}(x) \times 10^{-4}} \\ & = 0.016 x - .00022 x^2 + 2.36 \\ & \pm 0.0196 \sqrt{70.608 - 9.356 x + 0.54428 x^2 - 0.01462 x^3 + 0.000148 x^4}. \end{aligned}$$

The estimated regression and 95% intervals are shown in Figure 7.7. The regression interval widens greatly for small and large values of experience, indicating considerable uncertainty about the effect of experience on mean wages for this population. The confidence bands take a more complicated shape than in Figure 7.6 due to the nonlinear specification.

## 7.15 Forecast Intervals

Suppose we are given a value of the regressor vector  $\mathbf{x}_{n+1}$  for an individual outside the sample, and we want to forecast (guess)  $y_{n+1}$  for this individual. This is equivalent to forecasting  $y_{n+1}$  given  $\mathbf{x}_{n+1} = \mathbf{x}$ , which will generally be a function of  $\mathbf{x}$ . A reasonable forecasting rule is the conditional mean  $m(\mathbf{x})$  as it is the mean-square-minimizing forecast. A point forecast is the estimated conditional mean  $\hat{m}(\mathbf{x}) = \mathbf{x}'\hat{\beta}$ . We would also like a measure of uncertainty for the forecast.

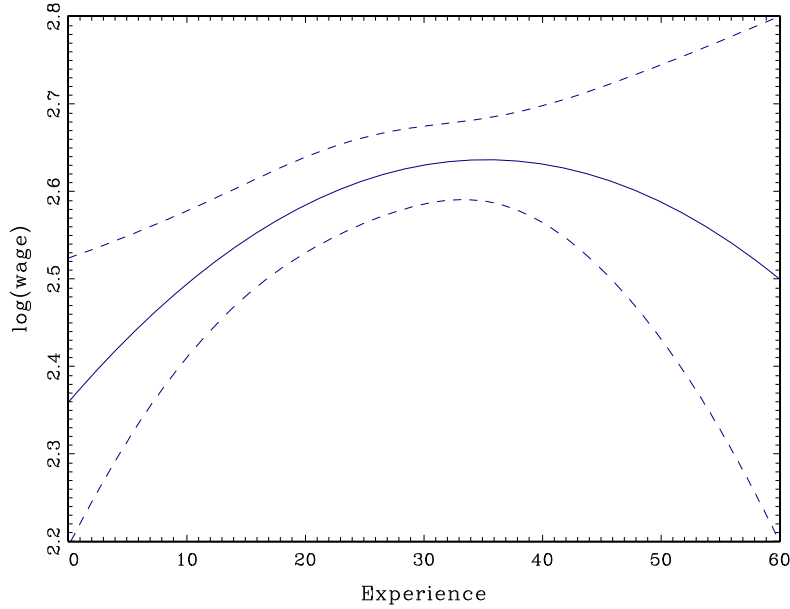


Figure 7.7: Wage on Experience Regression Intervals

The forecast error is  $\hat{e}_{n+1} = y_{n+1} - \hat{m}(\mathbf{x}) = e_{n+1} - \mathbf{x}'(\hat{\beta} - \beta)$ . As the out-of-sample error  $e_{n+1}$  is independent of the in-sample estimate  $\hat{\beta}$ , this has conditional variance

$$\begin{aligned} \mathbb{E}(\hat{e}_{n+1}^2 | \mathbf{x}_{n+1} = \mathbf{x}) &= \mathbb{E}\left(e_{n+1}^2 - 2\mathbf{x}'(\hat{\beta} - \beta)e_{n+1} + \mathbf{x}'(\hat{\beta} - \beta)(\hat{\beta} - \beta)\mathbf{x} | \mathbf{x}_{n+1} = \mathbf{x}\right) \\ &= \mathbb{E}(e_{n+1}^2 | \mathbf{x}_{n+1} = \mathbf{x}) + \mathbf{x}'\mathbb{E}(\hat{\beta} - \beta)(\hat{\beta} - \beta)'\mathbf{x} \\ &= \sigma^2(\mathbf{x}) + \mathbf{x}'\mathbf{V}_{\hat{\beta}}\mathbf{x}. \end{aligned}$$

Under homoskedasticity  $\mathbb{E}(e_{n+1}^2 | \mathbf{x}_{n+1}) = \sigma^2$ , the natural estimate of this variance is  $\hat{\sigma}^2 + \mathbf{x}'\hat{\mathbf{V}}_{\hat{\beta}}\mathbf{x}$ , so a standard error for the forecast is  $\hat{s}(\mathbf{x}) = \sqrt{\hat{\sigma}^2 + \mathbf{x}'\hat{\mathbf{V}}_{\hat{\beta}}\mathbf{x}}$ . Notice that this is different from the standard error for the conditional mean.

The conventional 95% forecast interval for  $y_{n+1}$  uses a normal approximation and sets

$$\left[\mathbf{x}'\hat{\beta} \pm 2\hat{s}(\mathbf{x})\right].$$

It is difficult, however, to fully justify this choice. It would be correct if we have a normal approximation to the ratio

$$\frac{e_{n+1} - \mathbf{x}'(\hat{\beta} - \beta)}{\hat{s}(\mathbf{x})}.$$

The difficulty is that the equation error  $e_{n+1}$  is generally non-normal, and asymptotic theory cannot be applied to a single observation. The only special exception is the case where  $e_{n+1}$  has the exact distribution  $N(0, \sigma^2)$ , which is generally invalid.

To get an accurate forecast interval, we need to estimate the conditional distribution of  $e_{n+1}$  given  $\mathbf{x}_{n+1} = \mathbf{x}$ , which is a much more difficult task. Perhaps due to this difficulty, many applied forecasters use the simple approximate interval  $\left[\mathbf{x}'\hat{\beta} \pm 2\hat{s}(\mathbf{x})\right]$  despite the lack of a convincing justification.

## 7.16 Wald Statistic

Let  $\boldsymbol{\theta} = \mathbf{r}(\boldsymbol{\beta}) : \mathbb{R}^k \rightarrow \mathbb{R}^q$  be any parameter vector of interest,  $\hat{\boldsymbol{\theta}}$  its estimate and  $\hat{\mathbf{V}}_{\hat{\boldsymbol{\theta}}}$  its covariance matrix estimator. Consider the quadratic form

$$W(\boldsymbol{\theta}) = (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})' \hat{\mathbf{V}}_{\hat{\boldsymbol{\theta}}}^{-1} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) = n (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})' \hat{\mathbf{V}}_{\hat{\boldsymbol{\theta}}}^{-1} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}). \quad (7.47)$$

where  $\hat{\mathbf{V}}_{\boldsymbol{\theta}} = n \hat{\mathbf{V}}_{\hat{\boldsymbol{\theta}}}$ . When  $q = 1$ , then  $W(\boldsymbol{\theta}) = T(\boldsymbol{\theta})^2$  is the square of the t-ratio. When  $q > 1$ ,  $W(\boldsymbol{\theta})$  is typically called a **Wald statistic**. We are interested in its sampling distribution.

The asymptotic distribution of  $W(\boldsymbol{\theta})$  is simple to derive given Theorem 7.10.2 and Theorem 7.10.3, which show that

$$\sqrt{n} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{d} Z \sim N(\mathbf{0}, \mathbf{V}_{\boldsymbol{\theta}})$$

and

$$\hat{\mathbf{V}}_{\boldsymbol{\theta}} \xrightarrow{p} \mathbf{V}_{\boldsymbol{\theta}}.$$

Note that  $\mathbf{V}_{\boldsymbol{\theta}} > 0$  since  $\mathbf{R}$  is full rank under Assumption 7.10.1. It follows that

$$W(\boldsymbol{\theta}) = \sqrt{n} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})' \hat{\mathbf{V}}_{\hat{\boldsymbol{\theta}}}^{-1} \sqrt{n} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{d} Z' \mathbf{V}_{\boldsymbol{\theta}}^{-1} Z \quad (7.48)$$

a quadratic in the normal random vector  $Z$ . As shown in Theorem 5.3.3, the distribution of this quadratic form is  $\chi_q^2$ , a chi-square random variable with  $q$  degrees of freedom.

**Theorem 7.16.1** Under Assumptions 7.1.2 and 7.10.1, as  $n \rightarrow \infty$ ,

$$W(\boldsymbol{\theta}) \xrightarrow{d} \chi_q^2.$$

Theorem 7.16.1 is used to justify multivariate confidence regions and multivariate hypothesis tests.

## 7.17 Homoskedastic Wald Statistic

Under the conditional homoskedasticity assumption  $\mathbb{E}(e_i^2 | \mathbf{x}_i) = \sigma^2$  we can construct the Wald statistic using the homoskedastic covariance matrix estimator  $\hat{\mathbf{V}}_{\boldsymbol{\theta}}^0$  defined in (7.39). This yields a homoskedastic Wald statistic

$$W^0(\boldsymbol{\theta}) = (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})' (\hat{\mathbf{V}}_{\hat{\boldsymbol{\theta}}}^0)^{-1} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) = n (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})' (\hat{\mathbf{V}}_{\hat{\boldsymbol{\theta}}}^0)^{-1} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}). \quad (7.49)$$

Under the additional assumption of conditional homoskedasticity, it has the same asymptotic distribution as  $W(\boldsymbol{\theta})$ .

**Theorem 7.17.1** Under Assumptions 7.1.2 and 7.10.1, and  $\mathbb{E}(e_i^2 | \mathbf{x}_i) = \sigma^2$ , as  $n \rightarrow \infty$ ,

$$W^0(\boldsymbol{\theta}) \xrightarrow{d} \chi_q^2.$$

## 7.18 Confidence Regions

A confidence region  $\widehat{C}$  is a set estimator for  $\boldsymbol{\theta} \in \mathbb{R}^q$  when  $q > 1$ . A confidence region  $\widehat{C}$  is a set in  $\mathbb{R}^q$  intended to cover the true parameter value with a pre-selected probability  $1 - \alpha$ . Thus an ideal confidence region has the coverage probability  $\Pr(\boldsymbol{\theta} \in \widehat{C}) = 1 - \alpha$ . In practice it is typically not possible to construct a region with exact coverage, but we can calculate its asymptotic coverage.

When the parameter estimate satisfies the conditions of Theorem 7.16.1, a good choice for a confidence region is the ellipse

$$\widehat{C} = \{\boldsymbol{\theta} : W(\boldsymbol{\theta}) \leq c_{1-\alpha}\}.$$

with  $c_{1-\alpha}$  the  $1 - \alpha$  quantile of the  $\chi_q^2$  distribution. (Thus  $F_q(c_{1-\alpha}) = 1 - \alpha$ .) It can be computed by, for example, `chi2inv(1- $\alpha$ ,q)` in MATLAB.

Theorem 7.16.1 implies

$$\Pr(\boldsymbol{\theta} \in \widehat{C}) \rightarrow \Pr(\chi_q^2 \leq c_{1-\alpha}) = 1 - \alpha$$

which shows that  $\widehat{C}$  has asymptotic coverage  $1 - \alpha$ .

To illustrate the construction of a confidence region, consider the estimated regression (7.41) of the model

$$\widehat{\log(\text{Wage})} = \beta_1 \text{education} + \beta_2 \text{experience} + \beta_3 \text{experience}^2/100 + \beta_4.$$

Suppose that the two parameters of interest are the percentage return to education  $\theta_1 = 100\beta_1$  and the percentage return to experience for individuals with 10 years experience  $\theta_2 = 100\beta_2 + 20\beta_3$ . These two parameters are a linear transformation of the regression parameters with point estimates

$$\widehat{\boldsymbol{\theta}} = \begin{pmatrix} 100 & 0 & 0 & 0 \\ 0 & 100 & 20 & 0 \end{pmatrix} \widehat{\boldsymbol{\beta}} = \begin{pmatrix} 11.8 \\ 1.2 \end{pmatrix},$$

and have the covariance matrix estimate

$$\begin{aligned} \widehat{\mathbf{V}}_{\widehat{\boldsymbol{\theta}}} &= \begin{pmatrix} 0 & 100 & 0 & 0 \\ 0 & 0 & 100 & 20 \end{pmatrix} \widehat{\mathbf{V}}_{\widehat{\boldsymbol{\beta}}} \begin{pmatrix} 0 & 0 \\ 100 & 0 \\ 0 & 100 \\ 0 & 20 \end{pmatrix} \\ &= \begin{pmatrix} 0.632 & 0.103 \\ 0.103 & 0.157 \end{pmatrix} \end{aligned}$$

with inverse

$$\widehat{\mathbf{V}}_{\widehat{\boldsymbol{\theta}}}^{-1} = \begin{pmatrix} 1.77 & -1.16 \\ -1.16 & 7.13 \end{pmatrix}.$$

Thus the Wald statistic is

$$\begin{aligned} W(\boldsymbol{\theta}) &= (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta})' \widehat{\mathbf{V}}_{\widehat{\boldsymbol{\theta}}}^{-1} (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \\ &= \begin{pmatrix} 11.8 - \theta_1 \\ 1.2 - \theta_2 \end{pmatrix}' \begin{pmatrix} 1.77 & -1.16 \\ -1.16 & 7.13 \end{pmatrix} \begin{pmatrix} 11.8 - \theta_1 \\ 1.2 - \theta_2 \end{pmatrix} \\ &= 1.77(11.8 - \theta_1)^2 - 2.32(11.8 - \theta_1)(1.2 - \theta_2) + 7.13(1.2 - \theta_2)^2. \end{aligned}$$

The 90% quantile of the  $\chi_2^2$  distribution is 4.605 (we use the  $\chi_2^2$  distribution as the dimension of  $\boldsymbol{\theta}$  is two), so an asymptotic 90% confidence region for the two parameters is the interior of the ellipse  $W(\boldsymbol{\theta}) = 4.605$  which is displayed in Figure 7.8. Since the estimated correlation of the two coefficient estimates is modest (about 0.3) the region is modestly elliptical.

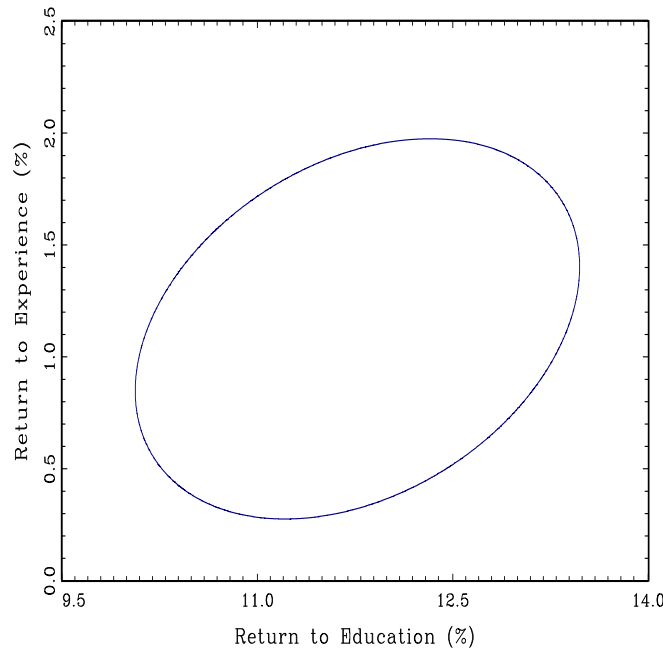


Figure 7.8: Confidence Region for Return to Experience and Return to Education

## 7.19 Semiparametric Efficiency in the Projection Model

In Section 4.7 we presented the Gauss-Markov theorem, which stated that in the homoskedastic CEF model, in the class of linear unbiased estimators the one with the smallest variance is least-squares. As we noted in that section, the restriction to linear unbiased estimators is unsatisfactory as it leaves open the possibility that an alternative (non-linear) estimator could have a smaller asymptotic variance. In addition, the restriction to the homoskedastic CEF model is also unsatisfactory as the projection model is more relevant for empirical application. The question remains: what is the most efficient estimator of the projection coefficient  $\beta$  (or functions  $\theta = h(\beta)$ ) in the projection model?

It turns out that it is straightforward to show that the projection model falls in the estimator class considered in Proposition 6.15.2. It follows that the least-squares estimator is semiparametrically efficient in the sense that it has the smallest asymptotic variance in the class of semiparametric estimators of  $\beta$ . This is a more powerful and interesting result than the Gauss-Markov theorem.

To see this, it is worth rephrasing Proposition 6.15.2 with amended notation. Suppose that a parameter of interest is  $\theta = g(\mu)$  where  $\mu = \mathbb{E}(z_i)$ , for which the moment estimators are  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n z_i$  and  $\hat{\theta} = g(\hat{\mu})$ . Let

$$\mathcal{L}_2(g) = \left\{ F : \mathbb{E} \|z\|^2 < \infty, \ g(u) \text{ is continuously differentiable at } u = \mathbb{E}(z) \right\}$$

be the set of distributions for which  $\hat{\theta}$  satisfies the central limit theorem.

**Proposition 7.19.1** *In the class of distributions  $F \in \mathcal{L}_2(g)$ ,  $\hat{\theta}$  is semiparametrically efficient for  $\theta$  in the sense that its asymptotic variance equals the semiparametric efficiency bound.*

Proposition 7.19.1 says that under the minimal conditions in which  $\hat{\theta}$  is asymptotically normal, then no semiparametric estimator can have a smaller asymptotic variance than  $\hat{\theta}$ .

To show that an estimator is semiparametrically efficient it is sufficient to show that it falls in the class covered by this Proposition. To show that the projection model falls in this class, we write  $\beta = Q_{xx}^{-1}Q_{xy} = g(\mu)$  where  $\mu = \mathbb{E}(z_i)$  and  $z_i = (x_i x_i', x_i y_i)$ . The class  $\mathcal{L}_2(g)$  equals the class of distributions

$$\mathcal{L}_4(\beta) = \left\{ F : \mathbb{E}(y^4) < \infty, \mathbb{E}\|x\|^4 < \infty, \mathbb{E}(x_i x_i') > 0 \right\}.$$

**Proposition 7.19.2** *In the class of distributions  $F \in \mathcal{L}_4(\beta)$ , the least-squares estimator  $\hat{\beta}$  is semiparametrically efficient for  $\beta$ .*

The least-squares estimator is an asymptotically efficient estimator of the projection coefficient because the latter is a smooth function of sample moments and the model implies no further restrictions. However, if the class of permissible distributions is restricted to a strict subset of  $\mathcal{L}_4(\beta)$  then least-squares can be inefficient. For example, the linear CEF model with heteroskedastic errors is a strict subset of  $\mathcal{L}_4(\beta)$ , and the GLS estimator has a smaller asymptotic variance than OLS. In this case, the knowledge that true conditional mean is linear allows for more efficient estimation of the unknown parameter.

From Proposition 7.19.1 we can also deduce that plug-in estimators  $\hat{\theta} = h(\hat{\beta})$  are semiparametrically efficient estimators of  $\theta = h(\beta)$  when  $h$  is continuously differentiable. We can also deduce that other parameters estimators are semiparametrically efficient, such as  $\hat{\sigma}^2$  for  $\sigma^2$ . To see this, note that we can write

$$\begin{aligned} \sigma^2 &= \mathbb{E} \left( (y_i - x_i' \beta)^2 \right) \\ &= \mathbb{E}(y_i^2) - 2\mathbb{E}(y_i x_i') \beta + \beta' \mathbb{E}(x_i x_i') \beta \\ &= Q_{yy} - Q_{yx} Q_{xx}^{-1} Q_{xy} \end{aligned}$$

which is a smooth function of the moments  $Q_{yy}$ ,  $Q_{yx}$  and  $Q_{xx}$ . Similarly the estimator  $\hat{\sigma}^2$  equals

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n \hat{e}_i^2 \\ &= \hat{Q}_{yy} - \hat{Q}_{yx} \hat{Q}_{xx}^{-1} \hat{Q}_{xy}. \end{aligned}$$

Since the variables  $y_i^2$ ,  $y_i x_i'$  and  $x_i x_i'$  all have finite variances when  $F \in \mathcal{L}_4(\beta)$ , the conditions of Proposition 7.19.1 are satisfied. We conclude:

**Proposition 7.19.3** *In the class of distributions  $F \in \mathcal{L}_4(\beta)$ ,  $\hat{\sigma}^2$  is semiparametrically efficient for  $\sigma^2$ .*

## 7.20 Semiparametric Efficiency in the Homoskedastic Regression Model\*

In Section 7.19 we showed that the OLS estimator is semiparametrically efficient in the projection model. What if we restrict attention to the classical homoskedastic regression model? Is OLS still efficient in this class? In this section we derive the asymptotic semiparametric efficiency bound

for this model, and show that it is the same as that obtained by the OLS estimator. Therefore it turns out that least-squares is efficient in this class as well.

Recall that in the homoskedastic regression model the asymptotic variance of the OLS estimator  $\hat{\beta}$  for  $\beta$  is  $\mathbf{V}_{\beta}^0 = \mathbf{Q}_{xx}^{-1}\sigma^2$ . Therefore, as described in Section 6.15, it is sufficient to find a parametric submodel whose Cramer-Rao bound for estimation of  $\beta$  is  $\mathbf{V}_{\beta}^0$ . This would establish that  $\mathbf{V}_{\beta}^0$  is the semiparametric variance bound and the OLS estimator  $\hat{\beta}$  is semiparametrically efficient for  $\beta$ .

Let the joint density of  $y$  and  $\mathbf{x}$  be written as  $f(y, \mathbf{x}) = f_1(y | \mathbf{x}) f_2(\mathbf{x})$ , the product of the conditional density of  $y$  given  $\mathbf{x}$  and the marginal density of  $\mathbf{x}$ . Now consider the parametric submodel

$$f(y, \mathbf{x} | \boldsymbol{\theta}) = f_1(y | \mathbf{x}) (1 + (y - \mathbf{x}'\beta)(\mathbf{x}'\boldsymbol{\theta})/\sigma^2) f_2(\mathbf{x}). \quad (7.50)$$

You can check that in this submodel the marginal density of  $\mathbf{x}$  is  $f_2(\mathbf{x})$  and the conditional density of  $y$  given  $\mathbf{x}$  is  $f_1(y | \mathbf{x}) (1 + (y - \mathbf{x}'\beta)(\mathbf{x}'\boldsymbol{\theta})/\sigma^2)$ . To see that the latter is a valid conditional density, observe that the regression assumption implies that  $\int y f_1(y | \mathbf{x}) dy = \mathbf{x}'\beta$  and therefore

$$\begin{aligned} & \int f_1(y | \mathbf{x}) (1 + (y - \mathbf{x}'\beta)(\mathbf{x}'\boldsymbol{\theta})/\sigma^2) dy \\ &= \int f_1(y | \mathbf{x}) dy + \int f_1(y | \mathbf{x}) (y - \mathbf{x}'\beta) dy (\mathbf{x}'\boldsymbol{\theta})/\sigma^2 \\ &= 1. \end{aligned}$$

In this parametric submodel the conditional mean of  $y$  given  $\mathbf{x}$  is

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\theta}}(y | \mathbf{x}) &= \int y f_1(y | \mathbf{x}) (1 + (y - \mathbf{x}'\beta)(\mathbf{x}'\boldsymbol{\theta})/\sigma^2) dy \\ &= \int y f_1(y | \mathbf{x}) dy + \int y f_1(y | \mathbf{x}) (y - \mathbf{x}'\beta)(\mathbf{x}'\boldsymbol{\theta})/\sigma^2 dy \\ &= \int y f_1(y | \mathbf{x}) dy + \int (y - \mathbf{x}'\beta)^2 f_1(y | \mathbf{x}) (\mathbf{x}'\boldsymbol{\theta})/\sigma^2 dy \\ &\quad + \int (y - \mathbf{x}'\beta) f_1(y | \mathbf{x}) dy (\mathbf{x}'\beta)(\mathbf{x}'\boldsymbol{\theta})/\sigma^2 \\ &= \mathbf{x}'(\beta + \boldsymbol{\theta}), \end{aligned}$$

using the homoskedasticity assumption  $\int (y - \mathbf{x}'\beta)^2 f_1(y | \mathbf{x}) dy = \sigma^2$ . This means that in this parametric submodel, the conditional mean is linear in  $\mathbf{x}$  and the regression coefficient is  $\beta(\boldsymbol{\theta}) = \beta + \boldsymbol{\theta}$ .

We now calculate the score for estimation of  $\boldsymbol{\theta}$ . Since

$$\frac{\partial}{\partial \boldsymbol{\theta}} \log f(y, \mathbf{x} | \boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \log (1 + (y - \mathbf{x}'\beta)(\mathbf{x}'\boldsymbol{\theta})/\sigma^2) = \frac{\mathbf{x}(y - \mathbf{x}'\beta)/\sigma^2}{1 + (y - \mathbf{x}'\beta)(\mathbf{x}'\boldsymbol{\theta})/\sigma^2}$$

the score is

$$\mathbf{s} = \frac{\partial}{\partial \boldsymbol{\theta}} \log f(y, \mathbf{x} | \boldsymbol{\theta}_0) = \mathbf{x}e/\sigma^2.$$

The Cramer-Rao bound for estimation of  $\boldsymbol{\theta}$  (and therefore  $\beta(\boldsymbol{\theta})$  as well) is

$$(\mathbb{E}(\mathbf{s}\mathbf{s}'))^{-1} = (\sigma^{-4}\mathbb{E}((\mathbf{x}e)(\mathbf{x}e)'))^{-1} = \sigma^2\mathbf{Q}_{xx}^{-1} = \mathbf{V}_{\beta}^0.$$

We have shown that there is a parametric submodel (7.50) whose Cramer-Rao bound for estimation of  $\beta$  is identical to the asymptotic variance of the least-squares estimator, which therefore is the semiparametric variance bound.

**Theorem 7.20.1** *In the homoskedastic regression model, the semiparametric variance bound for estimation of  $\beta$  is  $\mathbf{V}^0 = \sigma^2 \mathbf{Q}_{xx}^{-1}$  and the OLS estimator is semiparametrically efficient.*

This result is similar to the Gauss-Markov theorem, in that it asserts the efficiency of the least-squares estimator in the context of the homoskedastic regression model. The difference is that the Gauss-Markov theorem states that OLS has the smallest variance among the set of unbiased linear estimators, while Theorem 7.20.1 states that OLS has the smallest asymptotic variance among all regular estimators. This is a much more powerful statement.

## 7.21 Uniformly Consistent Residuals\*

It seems natural to view the residuals  $\hat{e}_i$  as estimates of the unknown errors  $e_i$ . Are they consistent estimates? In this section we develop an appropriate convergence result. This is not a widely-used technique, and can safely be skipped by most readers.

Notice that we can write the residual as

$$\begin{aligned}\hat{e}_i &= y_i - \mathbf{x}_i' \hat{\beta} \\ &= e_i + \mathbf{x}_i' \beta - \mathbf{x}_i' \hat{\beta} \\ &= e_i - \mathbf{x}_i' (\hat{\beta} - \beta).\end{aligned}\tag{7.51}$$

Since  $\hat{\beta} - \beta \xrightarrow{p} \mathbf{0}$  it seems reasonable to guess that  $\hat{e}_i$  will be close to  $e_i$  if  $n$  is large.

We can bound the difference in (7.51) using the Schwarz inequality (A.20) to find

$$|\hat{e}_i - e_i| = |\mathbf{x}_i' (\hat{\beta} - \beta)| \leq \|\mathbf{x}_i\| \|\hat{\beta} - \beta\|.\tag{7.52}$$

To bound (7.52) we can use  $\|\hat{\beta} - \beta\| = O_p(n^{-1/2})$  from Theorem 7.3.2, but we also need to bound the random variable  $\|\mathbf{x}_i\|$ . If the regressor is bounded, that is,  $\|\mathbf{x}_i\| \leq B < \infty$ , then  $|\hat{e}_i - e_i| \leq B \|\hat{\beta} - \beta\| = O_p(n^{-1/2})$ . However if the regressor does not have bounded support then we have to be more careful.

The key is Theorem 6.14.1 which shows that  $\mathbb{E} \|\mathbf{x}_i\|^r < \infty$  implies  $\mathbf{x}_i = o_p(n^{1/r})$  uniformly in  $i$ , or

$$n^{-1/r} \max_{1 \leq i \leq n} \|\mathbf{x}_i\| \xrightarrow{p} 0.$$

Applied to (7.52) we obtain

$$\begin{aligned}\max_{1 \leq i \leq n} |\hat{e}_i - e_i| &\leq \max_{1 \leq i \leq n} \|\mathbf{x}_i\| \|\hat{\beta} - \beta\| \\ &= o_p(n^{-1/2+1/r}).\end{aligned}$$

We have shown the following.

**Theorem 7.21.1** *Under Assumption 7.1.2 and  $\mathbb{E} \|\mathbf{x}_i\|^r < \infty$ , then uniformly in  $1 \leq i \leq n$*

$$\hat{e}_i = e_i + o_p(n^{-1/2+1/r}).\tag{7.53}$$



The rate of convergence in (7.53) depends on  $r$ . Assumption 7.1.2 requires  $r \geq 4$ , so the rate of convergence is at least  $o_p(n^{-1/4})$ . As  $r$  increases, the rate improves. As a limiting case, from Theorem 6.14.1 we see that if  $\mathbb{E}(\exp(\mathbf{t}'\mathbf{x}_i)) < \infty$  for some  $\mathbf{t} \neq 0$  then  $\mathbf{x}_i = o_p((\log n)^{1+\eta})$  uniformly in  $i$ , and thus  $\hat{e}_i = e_i + o_p(n^{-1/2}(\log n)^{1+\eta})$ .

We mentioned in Section 7.7 that there are multiple ways to prove the consistent of the covariance matrix estimator  $\hat{\mathbf{\Omega}}$ . We now show that Theorem 7.21.1 provides one simple method to establish (7.31) and thus Theorem 7.7.1. Let  $q_n = \max_{1 \leq i \leq n} |\hat{e}_i - e_i| = o_p(n^{-1/4})$ . Since

$$\hat{e}_i^2 - e_i^2 = 2e_i(\hat{e}_i - e_i) + (\hat{e}_i - e_i)^2,$$

then

$$\begin{aligned} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' (\hat{e}_i^2 - e_i^2) \right\| &\leq \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i \mathbf{x}_i'\| |\hat{e}_i^2 - e_i^2| \\ &\leq \frac{2}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2 |e_i| |\hat{e}_i - e_i| + \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2 |\hat{e}_i - e_i|^2 \\ &\leq \frac{2}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2 |e_i| q_n + \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2 q_n^2 \\ &\leq o_p(n^{-1/4}). \end{aligned}$$

## 7.22 Asymptotic Leverage\*

Recall the definition of leverage from (3.25)

$$h_{ii} = \mathbf{x}_i' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i.$$

These are the diagonal elements of the projection matrix  $\mathbf{P}$  and appear in the formula for leave-one-out prediction errors and several covariance matrix estimators. We can show that under iid sampling the leverage values are uniformly asymptotically small.

Let  $\lambda_{\min}(\mathbf{A})$  and  $\lambda_{\max}(\mathbf{A})$  denote the smallest and largest eigenvalues of a symmetric square matrix  $\mathbf{A}$ , and note that  $\lambda_{\max}(\mathbf{A}^{-1}) = (\lambda_{\min}(\mathbf{A}))^{-1}$ .

Since  $\frac{1}{n} \mathbf{X}'\mathbf{X} \xrightarrow{p} \mathbf{Q}_{xx} > 0$  then by the CMT,  $\lambda_{\min}(\frac{1}{n} \mathbf{X}'\mathbf{X}) \xrightarrow{p} \lambda_{\min}(\mathbf{Q}_{xx}) > 0$ . (The latter is positive since  $\mathbf{Q}_{xx}$  is positive definite and thus all its eigenvalues are positive.) Then by the Quadratic Inequality (A.28)

$$\begin{aligned} h_{ii} &= \mathbf{x}_i' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i \\ &\leq \lambda_{\max}((\mathbf{X}'\mathbf{X})^{-1}) (\mathbf{x}_i' \mathbf{x}_i) \\ &= \left( \lambda_{\min}\left(\frac{1}{n} \mathbf{X}'\mathbf{X}\right) \right)^{-1} \frac{1}{n} \|\mathbf{x}_i\|^2 \\ &\leq (\lambda_{\min}(\mathbf{Q}_{xx}) + o_p(1))^{-1} \frac{1}{n} \max_{1 \leq i \leq n} \|\mathbf{x}_i\|^2. \end{aligned} \tag{7.54}$$

Theorem 6.14.1 shows that  $\mathbb{E} \|\mathbf{x}_i\|^r < \infty$  implies  $\max_{1 \leq i \leq n} \|\mathbf{x}_i\|^2 = (\max_{1 \leq i \leq n} \|\mathbf{x}_i\|)^2 = o_p(n^{2/r})$  and thus (7.54) is  $o_p(n^{2/r-1})$ .

**Theorem 7.22.1** *If  $\mathbf{x}_i$  is independent and identically distributed and  $\mathbb{E} \|\mathbf{x}_i\|^r < \infty$  for some  $r \geq 2$ , then uniformly in  $1 \leq i \leq n$ ,  $h_{ii} = o_p(n^{2/r-1})$ .*

For any  $r \geq 2$  then  $h_{ii} = o_p(1)$  (uniformly in  $i \leq n$ ). Larger  $r$  implies a stronger rate of convergence, for example  $r = 4$  implies  $h_{ii} = o_p(n^{-1/2})$ .

Theorem (7.22.1) implies that under random sampling with finite variances and large samples, no individual observation should have a large leverage value. Consequently individual observations should not be influential, unless one of these conditions is violated.

## Exercises

**Exercise 7.1** Take the model  $y_i = \mathbf{x}'_{1i}\boldsymbol{\beta}_1 + \mathbf{x}'_{2i}\boldsymbol{\beta}_2 + e_i$  with  $\mathbb{E}(\mathbf{x}_i e_i) = \mathbf{0}$ . Suppose that  $\boldsymbol{\beta}_1$  is estimated by regressing  $y_i$  on  $\mathbf{x}_{1i}$  only. Find the probability limit of this estimator. In general, is it consistent for  $\boldsymbol{\beta}_1$ ? If not, under what conditions is this estimator consistent for  $\boldsymbol{\beta}_1$ ?

**Exercise 7.2** Let  $\mathbf{y}$  be  $n \times 1$ ,  $\mathbf{X}$  be  $n \times k$  (rank  $k$ ).  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$  with  $\mathbb{E}(\mathbf{x}_i e_i) = 0$ . Define the *ridge regression* estimator

$$\hat{\boldsymbol{\beta}} = \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}'_i + \lambda \mathbf{I}_k \right)^{-1} \left( \sum_{i=1}^n \mathbf{x}_i y_i \right) \quad (7.55)$$

here  $\lambda > 0$  is a fixed constant. Find the probability limit of  $\hat{\boldsymbol{\beta}}$  as  $n \rightarrow \infty$ . Is  $\hat{\boldsymbol{\beta}}$  consistent for  $\boldsymbol{\beta}$ ?

**Exercise 7.3** For the ridge regression estimator (7.55), set  $\lambda = cn$  where  $c > 0$  is fixed as  $n \rightarrow \infty$ . Find the probability limit of  $\hat{\boldsymbol{\beta}}$  as  $n \rightarrow \infty$ .

**Exercise 7.4** Verify some of the calculations reported in Section 7.4. Specifically, suppose that  $x_{1i}$  and  $x_{2i}$  only take the values  $\{-1, +1\}$ , symmetrically, with

$$\begin{aligned} \Pr(x_{1i} = x_{2i} = 1) &= \Pr(x_{1i} = x_{2i} = -1) = 3/8 \\ \Pr(x_{1i} = 1, x_{2i} = -1) &= \Pr(x_{1i} = -1, x_{2i} = 1) = 1/8 \\ \mathbb{E}(e_i^2 \mid x_{1i} = x_{2i}) &= \frac{5}{4} \\ \mathbb{E}(e_i^2 \mid x_{1i} \neq x_{2i}) &= \frac{1}{4}. \end{aligned}$$

Verify the following:

- (a)  $\mathbb{E}(x_{1i}) = 0$
- (b)  $\mathbb{E}(x_{1i}^2) = 1$
- (c)  $\mathbb{E}(x_{1i}x_{2i}) = \frac{1}{2}$
- (d)  $\mathbb{E}(e_i^2) = 1$
- (e)  $\mathbb{E}(x_{1i}^2 e_i^2) = 1$
- (f)  $\mathbb{E}(x_{1i}x_{2i}e_i^2) = \frac{7}{8}$ .

**Exercise 7.5** Show (7.19)-(7.22).

**Exercise 7.6** The model is

$$\begin{aligned} y_i &= \mathbf{x}'_i \boldsymbol{\beta} + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= \mathbf{0} \\ \boldsymbol{\Omega} &= \mathbb{E}(\mathbf{x}_i \mathbf{x}'_i e_i^2). \end{aligned}$$

Find the method of moments estimators  $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\Omega}})$  for  $(\boldsymbol{\beta}, \boldsymbol{\Omega})$ .

- (a) In this model, are  $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\Omega}})$  efficient estimators of  $(\boldsymbol{\beta}, \boldsymbol{\Omega})$ ?
- (b) If so, in what sense are they efficient?

**Exercise 7.7** Of the variables  $(y_i^*, y_i, \mathbf{x}_i)$  only the pair  $(y_i, \mathbf{x}_i)$  are observed. In this case, we say that  $y_i^*$  is a *latent* variable. Suppose

$$\begin{aligned} y_i^* &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= \mathbf{0} \\ y_i &= y_i^* + u_i \end{aligned}$$

where  $u_i$  is a measurement error satisfying

$$\begin{aligned} \mathbb{E}(\mathbf{x}_i u_i) &= \mathbf{0} \\ \mathbb{E}(y_i^* u_i) &= 0 \end{aligned}$$

Let  $\hat{\boldsymbol{\beta}}$  denote the OLS coefficient from the regression of  $y_i$  on  $\mathbf{x}_i$ .

- (a) Is  $\boldsymbol{\beta}$  the coefficient from the linear projection of  $y_i$  on  $\mathbf{x}_i$ ?
- (b) Is  $\hat{\boldsymbol{\beta}}$  consistent for  $\boldsymbol{\beta}$  as  $n \rightarrow \infty$ ?
- (c) Find the asymptotic distribution of  $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$  as  $n \rightarrow \infty$ .

**Exercise 7.8** Find the asymptotic distribution of  $\sqrt{n}(\hat{\sigma}^2 - \sigma^2)$  as  $n \rightarrow \infty$ .

**Exercise 7.9** The model is

$$\begin{aligned} y_i &= x_i \beta + e_i \\ \mathbb{E}(e_i | x_i) &= 0 \end{aligned}$$

where  $x_i \in \mathbb{R}$ . Consider the two estimators

$$\begin{aligned} \hat{\beta} &= \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} \\ \tilde{\beta} &= \frac{1}{n} \sum_{i=1}^n \frac{y_i}{x_i}. \end{aligned}$$

- (a) Under the stated assumptions, are both estimators consistent for  $\beta$ ?
- (b) Are there conditions under which either estimator is efficient?

**Exercise 7.10** In the homoskedastic regression model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$  with  $\mathbb{E}(e_i | \mathbf{x}_i) = 0$  and  $\mathbb{E}(e_i^2 | \mathbf{x}_i) = \sigma^2$ , suppose  $\hat{\boldsymbol{\beta}}$  is the OLS estimate of  $\boldsymbol{\beta}$  with covariance matrix estimate  $\hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}$ , based on a sample of size  $n$ . Let  $\hat{\sigma}^2$  be the estimate of  $\sigma^2$ . You wish to forecast an out-of-sample value of  $y_{n+1}$  given that  $\mathbf{x}_{n+1} = \mathbf{x}$ . Thus the available information is the sample  $(\mathbf{y}, \mathbf{X})$ , the estimates  $(\hat{\boldsymbol{\beta}}, \hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}, \hat{\sigma}^2)$ , the residuals  $\hat{\mathbf{e}}$ , and the out-of-sample value of the regressors,  $\mathbf{x}_{n+1}$ .

- (a) Find a point forecast of  $y_{n+1}$ .
- (b) Find an estimate of the variance of this forecast.

**Exercise 7.11** Take a regression model with i.i.d. observations  $(y_i, x_i)$  and scalar  $x_i$

$$\begin{aligned} y_i &= x_i \beta + e_i \\ \mathbb{E}(e_i | x_i) &= 0 \\ \Omega &= \mathbb{E}(x_i^2 e_i^2) \end{aligned}$$

Let  $\hat{\beta}$  be the OLS estimate of  $\beta$  with residuals  $\hat{e}_i = y_i - x_i\hat{\beta}$ . Consider the estimates of  $\Omega$

$$\begin{aligned}\tilde{\Omega} &= \frac{1}{n} \sum_{i=1}^n x_i^2 e_i^2 \\ \hat{\Omega} &= \frac{1}{n} \sum_{i=1}^n x_i^2 \hat{e}_i^2\end{aligned}$$

- (a) Find the asymptotic distribution of  $\sqrt{n}(\tilde{\Omega} - \Omega)$  as  $n \rightarrow \infty$ .
- (b) Find the asymptotic distribution of  $\sqrt{n}(\hat{\Omega} - \Omega)$  as  $n \rightarrow \infty$ .
- (c) How do you use the regression assumption  $\mathbb{E}(e_i | x_i) = 0$  in your answer to (b)?

**Exercise 7.12** Consider the model

$$\begin{aligned}y_i &= \alpha + \beta x_i + e_i \\ \mathbb{E}(e_i) &= 0 \\ \mathbb{E}(x_i e_i) &= 0\end{aligned}$$

with both  $y_i$  and  $x_i$  scalar. Assuming  $\alpha > 0$  and  $\beta < 0$ , suppose the parameter of interest is the area under the regression curve (e.g. consumer surplus), which is  $A = -\alpha^2/2\beta$ .

Let  $\hat{\theta} = (\hat{\alpha}, \hat{\beta})'$  be the least-squares estimates of  $\theta = (\alpha, \beta)'$  so that  $\sqrt{n}(\hat{\theta} - \theta) \rightarrow_d N(0, \mathbf{V}_\theta)$  and let  $\hat{\mathbf{V}}_\theta$  be a standard consistent estimate for  $\mathbf{V}_\theta$ .

- (a) Given the above, describe an estimator of  $A$ .
- (b) Construct an asymptotic  $(1 - \eta)$  confidence interval for  $A$ .

**Exercise 7.13** Consider an iid sample  $\{y_i, x_i\}$   $i = 1, \dots, n$  where  $y_i$  and  $x_i$  are scalar. Consider the reverse projection model

$$\begin{aligned}x_i &= y_i\gamma + u_i \\ \mathbb{E}(y_i u_i) &= 0\end{aligned}$$

and define the parameter of interest as  $\theta = 1/\gamma$

- (a) Propose an estimator  $\hat{\gamma}$  of  $\gamma$ .
- (b) Propose an estimator  $\hat{\theta}$  of  $\theta$ .
- (c) Find the asymptotic distribution of  $\hat{\theta}$ .
- (d) Find an asymptotic standard error for  $\hat{\theta}$ .

**Exercise 7.14** Take the model

$$\begin{aligned}y_i &= x_{1i}\beta_1 + x_{2i}\beta_2 + e_i \\ \mathbb{E}(x_i e_i) &= 0\end{aligned}$$

with both  $\beta_1 \in \mathbb{R}$  and  $\beta_2 \in \mathbb{R}$ , and define the parameter

$$\theta = \beta_1\beta_2$$

- (a) What is the appropriate estimator  $\hat{\theta}$  for  $\theta$ ?
- (b) Find the asymptotic distribution of  $\hat{\theta}$  under standard regularity conditions.
- (c) Show how to calculate an asymptotic 95% confidence interval for  $\theta$ .

**Exercise 7.15** Take the linear model

$$y_i = x_i\beta + e_i$$

$$\mathbb{E}(e_i \mid x_i) = 0$$

with  $n$  observations and  $x_i$  is scalar (real-valued). Consider the estimator

$$\hat{\beta} = \frac{\sum_{i=1}^n x_i^3 y_i}{\sum_{i=1}^n x_i^4}$$

Find the asymptotic distribution of  $\sqrt{n}(\hat{\beta} - \beta)$  as  $n \rightarrow \infty$ .

**Exercise 7.16** Out of an iid sample  $(y_i, \mathbf{x}_i)$  of size  $n$ , you randomly take half the observations and estimate the least-squares regression of  $y_i$  on  $\mathbf{x}_i$  using only this sub-sample.

$$y_i = \mathbf{x}_i' \hat{\boldsymbol{\beta}} + \hat{e}_i$$

Is the estimated slope coefficient  $\hat{\boldsymbol{\beta}}$  consistent for the population projection coefficient? Explain your reasoning.

**Exercise 7.17** An economist reports a set of parameter estimates, including the coefficient estimates  $\hat{\beta}_1 = 1.0$ ,  $\hat{\beta}_2 = 0.8$ , and standard errors  $s(\hat{\beta}_1) = 0.07$  and  $s(\hat{\beta}_2) = 0.07$ . The author writes “The estimates show that  $\beta_1$  is larger than  $\beta_2$ .”

- (a) Write down the formula for an asymptotic 95% confidence interval for  $\theta = \beta_1 - \beta_2$ , expressed as a function of  $\hat{\beta}_1$ ,  $\hat{\beta}_2$ ,  $s(\hat{\beta}_1)$ ,  $s(\hat{\beta}_2)$  and  $\hat{\rho}$ , where  $\hat{\rho}$  is the estimated correlation between  $\hat{\beta}_1$  and  $\hat{\beta}_2$ .
- (b) Can  $\hat{\rho}$  be calculated from the reported information?
- (c) Is the author correct? Does the reported information support the author’s claim?

**Exercise 7.18** Suppose an economic model suggests

$$g(x) = \mathbb{E}(y_i \mid x_i = x) = \beta_0 + \beta_1 x + \beta_2 x^2$$

where  $x_i \in \mathbb{R}$ . You have a random sample  $(y_i, x_i)$ ,  $i = 1, \dots, n$ .

- (a) Describe how to estimate  $g(x)$  at a given value  $x$ .
- (b) Describe (be specific) an appropriate confidence interval for  $g(x)$ .

**Exercise 7.19** Take the model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(\mathbf{x}_i e_i) = \mathbf{0}$$

and suppose you have observations  $i = 1, \dots, 2n$ . (The number of observations is  $2n$ .) You randomly split the sample in half, (each has  $n$  observations), calculate  $\hat{\boldsymbol{\beta}}_1$  by least-squares on the first sample, and  $\hat{\boldsymbol{\beta}}_2$  by least-squares on the second sample. What is the asymptotic distribution of  $\sqrt{n}(\hat{\boldsymbol{\beta}}_1 - \hat{\boldsymbol{\beta}}_2)$ ?

**Exercise 7.20** The data  $\{y_i, \mathbf{x}_i, w_i\}$  is from a random sample,  $i = 1, \dots, n$ . The parameter  $\beta$  is estimated by minimizing the criterion function

$$S(\beta) = \sum_{i=1}^n w_i (y_i - \mathbf{x}_i' \beta)^2$$

That is  $\hat{\beta} = \operatorname{argmin}_{\beta} S(\beta)$ .

- Find an explicit expression for  $\hat{\beta}$ .
- What population parameter  $\beta$  is  $\hat{\beta}$  estimating? (Be explicit about any assumptions you need to impose. But don't make more assumptions than necessary.)
- Find the probability limit for  $\hat{\beta}$  as  $n \rightarrow \infty$ .
- Find the asymptotic distribution of  $\sqrt{n}(\hat{\beta} - \beta)$  as  $n \rightarrow \infty$ .

**Exercise 7.21** Take the model

$$\begin{aligned} y_i &= \mathbf{x}_i' \beta + e_i \\ \mathbb{E}(e_i \mid \mathbf{x}_i) &= 0 \\ \mathbb{E}(e_i^2 \mid \mathbf{x}_i) &= \sigma_i^2 = \mathbf{z}_i' \gamma \end{aligned}$$

where  $\mathbf{z}_i$  is a (vector) function of  $\mathbf{x}_i$ . The sample is  $i = 1, \dots, n$  with iid observations. For simplicity, assume that  $\mathbf{z}_i' \gamma > 0$  for all  $\mathbf{z}_i$ . Suppose you are interested in forecasting  $y_{n+1}$  given  $\mathbf{x}_{n+1} = \mathbf{x}$  and  $\mathbf{z}_{n+1} = \mathbf{z}$  for some out-of-sample observation  $n+1$ . Describe how you would construct a point forecast and a forecast interval for  $y_{n+1}$ .

**Exercise 7.22** Take the model

$$\begin{aligned} y_i &= \mathbf{x}_i' \beta + e_i \\ \mathbb{E}(e_i \mid \mathbf{x}_i) &= 0 \\ z_i &= (\mathbf{x}_i' \beta) \gamma + u_i \\ \mathbb{E}(u_i \mid \mathbf{x}_i) &= 0 \end{aligned}$$

Your goal is to estimate  $\gamma$ . (Note that  $\gamma$  is scalar.) You use a two-step estimator:

- Estimate  $\hat{\beta}$  by least-squares of  $y_i$  on  $\mathbf{x}_i$ .
  - Estimate  $\hat{\gamma}$  by least-squares of  $z_i$  on  $\mathbf{x}_i' \hat{\beta}$ .
- Show that  $\hat{\gamma}$  is consistent for  $\gamma$ .
  - Find the asymptotic distribution of  $\hat{\gamma}$  when  $\gamma = 0$ .

**Exercise 7.23** The model is

$$\begin{aligned} y_i &= x_i \beta + e_i \\ \mathbb{E}(e_i \mid x_i) &= 0 \end{aligned}$$

where  $x_i \in R$ . Consider the estimator

$$\tilde{\beta} = \frac{1}{n} \sum_{i=1}^n \frac{y_i}{x_i}.$$

Find conditions under which  $\tilde{\beta}$  is consistent for  $\beta$  as  $n \rightarrow \infty$ .

**Exercise 7.24** Of the random variables  $(y_i^*, y_i, \mathbf{x}_i)$  only the pair  $(y_i, \mathbf{x}_i)$  are observed. (In this case, we say that  $y_i^*$  is a *latent* variable.) Suppose  $\mathbb{E}(y_i^* | \mathbf{x}_i) = \mathbf{x}_i' \boldsymbol{\beta}$  and  $y = y_i^* + u_i$ , where  $u_i$  is a measurement error satisfying  $\mathbb{E}(u_i | y_i^*, \mathbf{x}_i) = 0$ . Let  $\hat{\boldsymbol{\beta}}$  denote the OLS coefficient from the regression of  $y_i$  on  $\mathbf{x}_i$ .

- (a) Find  $\mathbb{E}(y_i | \mathbf{x}_i)$ .
- (b) Is  $\hat{\boldsymbol{\beta}}$  consistent for  $\boldsymbol{\beta}$  as  $n \rightarrow \infty$ ?
- (c) Find the asymptotic distribution of  $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$  as  $n \rightarrow \infty$ .

**Exercise 7.25** The parameter of  $\beta$  is defined in the model

$$y_i = x_i^* \beta + e_i$$

where  $e_i$  is independent of  $x_i^*$ ,  $\mathbb{E}(e_i) = 0$ ,  $\mathbb{E}(e_i^2) = \sigma^2$ . The observables are  $(y_i, x_i)$  where

$$x_i = x_i^* v_i$$

and  $v_i > 0$  is random measurement error. Assume that  $v_i$  is independent of  $x_i^*$  and  $e_i$ . Also assume that  $x_i$  and  $x_i^*$  are non-negative and real-valued. Consider the least-squares estimator  $\hat{\beta}$  for  $\beta$ .

- (a) Find the plim of  $\hat{\beta}$ , expressed in terms of  $\beta$  and moments of  $(x_i, v_i, e_i)$
- (b) Can you find a non-trivial condition under which  $\hat{\beta}$  is consistent for  $\beta$ ? (By non-trivial, we mean something other than  $v_i = 1$ .)

**Exercise 7.26** Take the standard model

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= \mathbf{0} \end{aligned}$$

For a positive function  $w(\mathbf{x})$ , let  $w_i = w(\mathbf{x}_i)$ . Consider the estimator

$$\tilde{\boldsymbol{\beta}} = \left( \sum_{i=1}^n w_i \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left( \sum_{i=1}^n w_i \mathbf{x}_i y_i \right).$$

Find the probability limit (as  $n \rightarrow \infty$ ) of  $\tilde{\boldsymbol{\beta}}$ . (Do you need to add an assumption?) Is  $\tilde{\boldsymbol{\beta}}$  consistent for  $\boldsymbol{\beta}$ ? If not, under what assumption is  $\tilde{\boldsymbol{\beta}}$  consistent for  $\boldsymbol{\beta}$ ?

**Exercise 7.27** Take the regression model

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(e_i | \mathbf{x}_i) &= 0 \\ \mathbb{E}(e_i^2 | \mathbf{x}_i) &= \sigma_i^2 \end{aligned}$$

with  $\mathbf{x}_i \in R^k$ . Assume that  $\Pr(e_i = 0) = 0$ . Consider the infeasible estimator

$$\tilde{\boldsymbol{\beta}} = \left( \sum_{i=1}^n e_i^{-2} \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left( \sum_{i=1}^n e_i^{-2} \mathbf{x}_i y_i \right).$$

This is a WLS estimator using the weights  $e_i^{-2}$ .

- (a) Find the asymptotic distribution of  $\tilde{\boldsymbol{\beta}}$



- (b) Contrast your result with the asymptotic distribution of infeasible GLS.

**Exercise 7.28** The model is

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(e_i \mid \mathbf{x}_i) = 0.$$

An econometrician is worried about the impact of some unusually large values of the regressors. The model is thus estimated on the subsample for which  $|\mathbf{x}_i| \leq c$ , for some fixed  $c$ . Let  $\tilde{\boldsymbol{\beta}}$  denote the OLS estimator on this subsample. It equals

$$\tilde{\boldsymbol{\beta}} = \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' 1(|\mathbf{x}_i| \leq c) \right)^{-1} \left( \sum_{i=1}^n \mathbf{x}_i y_i 1(|\mathbf{x}_i| \leq c) \right)$$

where  $1(\cdot)$  denotes the indicator function.

- (a) Show that  $\tilde{\boldsymbol{\beta}} \rightarrow_p \boldsymbol{\beta}$ .
- (b) Find the asymptotic distribution of  $\sqrt{n}(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta})$

**Exercise 7.29** As in Exercise 3.24, use the CPS dataset and the subsample of white male Hispanics. Estimate the regression

$$\widehat{\log(Wage)} = \beta_1 \text{education} + \beta_2 \text{experience} + \beta_3 \text{experience}^2/100 + \beta_4.$$

- (a) Report the coefficients and robust standard errors.
- (b) Let  $\theta$  be the ratio of the return to one year of education to the return to one year of experience. Write  $\theta$  as a function of the regression coefficients and variables. Compute  $\hat{\theta}$  from the estimated model.
- (c) Write out the formula for the asymptotic standard error for  $\hat{\theta}$  as a function of the covariance matrix for  $\hat{\boldsymbol{\beta}}$ . Compute  $\widehat{s}(\hat{\theta})$  from the estimated model.
- (d) Construct a 90% asymptotic confidence interval for  $\theta$  from the estimated model.
- (e) Compute the regression function at  $edu = 12$  and  $experience = 20$ . Compute a 95% confidence interval for the regression function at this point.
- (f) Consider an out-of-sample individual with 16 years of education and 5 years experience. Construct an 80% forecast interval for their log wage and wage. [To obtain the forecast interval for the wage, apply the exponential function to both endpoints.]

# Chapter 8

## Restricted Estimation

### 8.1 Introduction

In the linear projection model

$$\begin{aligned}y_i &= \mathbf{x}'_i \boldsymbol{\beta} + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= \mathbf{0}\end{aligned}$$

a common task is to impose a constraint on the coefficient vector  $\boldsymbol{\beta}$ . For example, partitioning  $\mathbf{x}'_i = (\mathbf{x}'_{1i}, \mathbf{x}'_{2i})$  and  $\boldsymbol{\beta}' = (\boldsymbol{\beta}'_1, \boldsymbol{\beta}'_2)$ , a typical constraint is an exclusion restriction of the form  $\boldsymbol{\beta}_2 = \mathbf{0}$ . In this case the constrained model is

$$\begin{aligned}y_i &= \mathbf{x}'_{1i} \boldsymbol{\beta}_1 + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= \mathbf{0}\end{aligned}$$

At first glance this appears the same as the linear projection model, but there is one important difference: the error  $e_i$  is uncorrelated with the entire regressor vector  $\mathbf{x}'_i = (\mathbf{x}'_{1i}, \mathbf{x}'_{2i})$  not just the included regressor  $\mathbf{x}_{1i}$ .

In general, a set of  $q$  linear constraints on  $\boldsymbol{\beta}$  takes the form

$$\mathbf{R}'\boldsymbol{\beta} = \mathbf{c} \tag{8.1}$$

where  $\mathbf{R}$  is  $k \times q$ ,  $\text{rank}(\mathbf{R}) = q < k$  and  $\mathbf{c}$  is  $q \times 1$ . The assumption that  $\mathbf{R}$  is full rank means that the constraints are linearly independent (there are no redundant or contradictory constraints). We can define the restricted parameter space  $\mathbf{B}_R$  as the set of values of  $\boldsymbol{\beta}$  which satisfy (8.1), that is

$$\mathbf{B}_R = \{\boldsymbol{\beta} : \mathbf{R}'\boldsymbol{\beta} = \mathbf{c}\}$$

The constraint  $\boldsymbol{\beta}_2 = \mathbf{0}$  discussed above is a special case of the constraint (8.1) with

$$\mathbf{R} = \begin{pmatrix} \mathbf{0} \\ \mathbf{I} \end{pmatrix}, \tag{8.2}$$

a selector matrix, and  $\mathbf{c} = \mathbf{0}$ .

Another common restriction is that a set of coefficients sum to a known constant, i.e.  $\beta_1 + \beta_2 = 1$ . This constraint arises in a constant-return-to-scale production function. Other common restrictions include the equality of coefficients  $\beta_1 = \beta_2$ , and equal and offsetting coefficients  $\beta_1 = -\beta_2$ .

A typical reason to impose a constraint is that we believe (or have information) that the constraint is true. By imposing the constraint we hope to improve estimation efficiency. The goal is to obtain consistent estimates with reduced variance relative to the unconstrained estimator.

The questions then arise: How should we estimate the coefficient vector  $\boldsymbol{\beta}$  imposing the linear restriction (8.1)? If we impose such constraints, what is the sampling distribution of the resulting estimator? How should we calculate standard errors? These are the questions explored in this chapter.

## 8.2 Constrained Least Squares

An intuitively appealing method to estimate a constrained linear projection is to minimize the least-squares criterion subject to the constraint  $\mathbf{R}'\boldsymbol{\beta} = \mathbf{c}$ .

The constrained least-squares estimator is

$$\tilde{\boldsymbol{\beta}}_{\text{cls}} = \underset{\mathbf{R}'\boldsymbol{\beta}=\mathbf{c}}{\operatorname{argmin}} SSE(\boldsymbol{\beta}) \quad (8.3)$$

where

$$SSE(\boldsymbol{\beta}) = \sum_{i=1}^n (y_i - \mathbf{x}_i'\boldsymbol{\beta})^2 = \mathbf{y}'\mathbf{y} - 2\mathbf{y}'\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}. \quad (8.4)$$

The estimator  $\tilde{\boldsymbol{\beta}}_{\text{cls}}$  minimizes the sum of squared errors over all  $\boldsymbol{\beta}$  such that  $\boldsymbol{\beta} \in \mathbf{B}_{\mathbf{R}}$ , or equivalently such that the restriction (8.1) holds. We call  $\tilde{\boldsymbol{\beta}}_{\text{cls}}$  the **constrained least-squares** (CLS) estimator. We follow the convention of using a tilde “~” rather than a hat “^” to indicate that  $\tilde{\boldsymbol{\beta}}_{\text{cls}}$  is a restricted estimator in contrast to the unrestricted least-squares estimator  $\hat{\boldsymbol{\beta}}$ , and write it as  $\tilde{\boldsymbol{\beta}}_{\text{cls}}$  to be clear that the estimation method is CLS.

One method to find the solution to (8.3) uses the technique of Lagrange multipliers. The problem (8.3) is equivalent to the minimization of the Lagrangian

$$\mathcal{L}(\boldsymbol{\beta}, \boldsymbol{\lambda}) = \frac{1}{2}SSE(\boldsymbol{\beta}) + \boldsymbol{\lambda}'(\mathbf{R}'\boldsymbol{\beta} - \mathbf{c}) \quad (8.5)$$

over  $(\boldsymbol{\beta}, \boldsymbol{\lambda})$ , where  $\boldsymbol{\lambda}$  is an  $s \times 1$  vector of Lagrange multipliers. The first-order conditions for minimization of (8.5) are

$$\frac{\partial}{\partial \boldsymbol{\beta}} \mathcal{L}(\tilde{\boldsymbol{\beta}}_{\text{cls}}, \tilde{\boldsymbol{\lambda}}_{\text{cls}}) = -\mathbf{X}'\mathbf{y} + \mathbf{X}'\mathbf{X}\tilde{\boldsymbol{\beta}}_{\text{cls}} + \mathbf{R}\tilde{\boldsymbol{\lambda}}_{\text{cls}} = \mathbf{0} \quad (8.6)$$

and

$$\frac{\partial}{\partial \boldsymbol{\lambda}} \mathcal{L}(\tilde{\boldsymbol{\beta}}_{\text{cls}}, \tilde{\boldsymbol{\lambda}}_{\text{cls}}) = \mathbf{R}'\tilde{\boldsymbol{\beta}}_{\text{cls}} - \mathbf{c} = \mathbf{0}. \quad (8.7)$$

Premultiplying (8.6) by  $\mathbf{R}'(\mathbf{X}'\mathbf{X})^{-1}$  we obtain

$$-\mathbf{R}'\hat{\boldsymbol{\beta}} + \mathbf{R}'\tilde{\boldsymbol{\beta}}_{\text{cls}} + \mathbf{R}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}\tilde{\boldsymbol{\lambda}}_{\text{cls}} = \mathbf{0} \quad (8.8)$$

where  $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$  is the unrestricted least-squares estimator. Imposing  $\mathbf{R}'\tilde{\boldsymbol{\beta}}_{\text{cls}} - \mathbf{c} = \mathbf{0}$  from (8.7) and solving for  $\tilde{\boldsymbol{\lambda}}_{\text{cls}}$  we find

$$\tilde{\boldsymbol{\lambda}}_{\text{cls}} = \left[ \mathbf{R}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R} \right]^{-1} (\mathbf{R}'\hat{\boldsymbol{\beta}} - \mathbf{c}).$$

Notice that  $(\mathbf{X}'\mathbf{X})^{-1} > 0$  and  $\mathbf{R}$  full rank imply that  $\mathbf{R}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R} > 0$  and is hence invertible. (See Section A.9.)

Substituting this expression into (8.6) and solving for  $\tilde{\boldsymbol{\beta}}_{\text{cls}}$  we find the solution to the constrained minimization problem (8.3)

$$\tilde{\boldsymbol{\beta}}_{\text{cls}} = \hat{\boldsymbol{\beta}} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{R} \left[ \mathbf{R}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R} \right]^{-1} (\mathbf{R}'\hat{\boldsymbol{\beta}} - \mathbf{c}). \quad (8.9)$$

(See Exercise 8.5 to verify that (8.9) satisfies (8.1).)

This is a general formula for the CLS estimator. It also can be written as

$$\tilde{\boldsymbol{\beta}}_{\text{cls}} = \hat{\boldsymbol{\beta}} - \hat{\mathbf{Q}}_{\mathbf{x}\mathbf{x}}^{-1}\mathbf{R} \left( \mathbf{R}'\hat{\mathbf{Q}}_{\mathbf{x}\mathbf{x}}^{-1}\mathbf{R} \right)^{-1} (\mathbf{R}'\hat{\boldsymbol{\beta}} - \mathbf{c}). \quad (8.10)$$

The CLS residuals are

$$\tilde{e}_i = y_i - \mathbf{x}_i'\tilde{\boldsymbol{\beta}}_{\text{cls}}$$

and the  $n \times 1$  vector of residuals are written in vector notation as  $\tilde{\mathbf{e}}$ .

In Stata, constrained least squares is implemented using the **cnsreg** command.

### 8.3 Exclusion Restriction

While (8.9) is a general formula for the CLS estimator, in most cases the estimator can be found by applying least-squares to a reparameterized equation. To illustrate, let us return to the first example presented at the beginning of the chapter – a simple exclusion restriction. Recall the unconstrained model is

$$y_i = \mathbf{x}'_{1i}\boldsymbol{\beta}_1 + \mathbf{x}'_{2i}\boldsymbol{\beta}_2 + e_i \quad (8.11)$$

the exclusion restriction is  $\boldsymbol{\beta}_2 = \mathbf{0}$ , and the constrained equation is

$$y_i = \mathbf{x}'_{1i}\boldsymbol{\beta}_1 + e_i. \quad (8.12)$$

In this setting the CLS estimator is OLS of  $y_i$  on  $x_{1i}$ . (See Exercise 8.1.) We can write this as

$$\tilde{\boldsymbol{\beta}}_1 = \left( \sum_{i=1}^n \mathbf{x}_{1i} \mathbf{x}'_{1i} \right)^{-1} \left( \sum_{i=1}^n \mathbf{x}_{1i} y_i \right). \quad (8.13)$$

The CLS estimator of the entire vector  $\boldsymbol{\beta}' = (\boldsymbol{\beta}'_1, \boldsymbol{\beta}'_2)$  is

$$\tilde{\boldsymbol{\beta}} = \begin{pmatrix} \tilde{\boldsymbol{\beta}}_1 \\ \mathbf{0} \end{pmatrix}. \quad (8.14)$$

It is not immediately obvious, but (8.9) and (8.14) are algebraically (and numerically) equivalent. To see this, the first component of (8.9) with (8.2) is

$$\tilde{\boldsymbol{\beta}}_1 = (\mathbf{I} \quad \mathbf{0}) \left[ \hat{\boldsymbol{\beta}} - \hat{\mathbf{Q}}_{xx}^{-1} \begin{pmatrix} \mathbf{0} \\ \mathbf{I} \end{pmatrix} \left[ (\mathbf{0} \quad \mathbf{I}) \hat{\mathbf{Q}}_{xx}^{-1} \begin{pmatrix} \mathbf{0} \\ \mathbf{I} \end{pmatrix} \right]^{-1} (\mathbf{0} \quad \mathbf{I}) \hat{\boldsymbol{\beta}} \right].$$

Using (3.39) this equals

$$\begin{aligned} \tilde{\boldsymbol{\beta}}_1 &= \hat{\boldsymbol{\beta}}_1 - \hat{\mathbf{Q}}^{12} \left( \hat{\mathbf{Q}}^{22} \right)^{-1} \hat{\boldsymbol{\beta}}_2 \\ &= \hat{\boldsymbol{\beta}}_1 + \hat{\mathbf{Q}}_{11 \cdot 2}^{-1} \hat{\mathbf{Q}}_{12} \hat{\mathbf{Q}}_{22}^{-1} \hat{\mathbf{Q}}_{22 \cdot 1} \hat{\boldsymbol{\beta}}_2 \\ &= \hat{\mathbf{Q}}_{11 \cdot 2}^{-1} \left( \hat{\mathbf{Q}}_{1y} - \hat{\mathbf{Q}}_{12} \hat{\mathbf{Q}}_{22}^{-1} \hat{\mathbf{Q}}_{2y} \right) \\ &\quad + \hat{\mathbf{Q}}_{11 \cdot 2}^{-1} \hat{\mathbf{Q}}_{12} \hat{\mathbf{Q}}_{22}^{-1} \hat{\mathbf{Q}}_{22 \cdot 1} \hat{\mathbf{Q}}_{22 \cdot 1}^{-1} \left( \hat{\mathbf{Q}}_{2y} - \hat{\mathbf{Q}}_{21} \hat{\mathbf{Q}}_{11}^{-1} \hat{\mathbf{Q}}_{1y} \right) \\ &= \hat{\mathbf{Q}}_{11 \cdot 2}^{-1} \left( \hat{\mathbf{Q}}_{1y} - \hat{\mathbf{Q}}_{12} \hat{\mathbf{Q}}_{22}^{-1} \hat{\mathbf{Q}}_{21} \hat{\mathbf{Q}}_{11}^{-1} \hat{\mathbf{Q}}_{1y} \right) \\ &= \hat{\mathbf{Q}}_{11 \cdot 2}^{-1} \left( \hat{\mathbf{Q}}_{11} - \hat{\mathbf{Q}}_{12} \hat{\mathbf{Q}}_{22}^{-1} \hat{\mathbf{Q}}_{21} \right) \hat{\mathbf{Q}}_{11}^{-1} \hat{\mathbf{Q}}_{1y} \\ &= \hat{\mathbf{Q}}_{11}^{-1} \hat{\mathbf{Q}}_{1y} \end{aligned}$$

which is (8.14) as originally claimed.

### 8.4 Finite Sample Properties

In this section we explore some of the properties of the CLS estimator in the linear regression model

$$y_i = \mathbf{x}'_i \boldsymbol{\beta} + e_i \quad (8.15)$$

$$\mathbb{E}(e_i \mid \mathbf{x}_i) = 0. \quad (8.16)$$

First, it is useful to write the estimator, and the residuals, as linear functions of the error vector. These are algebraic relationships and do not rely on the linear regression assumptions.

**Theorem 8.4.1** Define  $\mathbf{P} = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$  and

$$\mathbf{A} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R} \left( \mathbf{R}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R} \right)^{-1} \mathbf{R}' (\mathbf{X}'\mathbf{X})^{-1}.$$

Then

1.  $\mathbf{R}'\hat{\boldsymbol{\beta}} - \mathbf{c} = \mathbf{R}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \mathbf{e}$
2.  $\tilde{\boldsymbol{\beta}}_{\text{cls}} - \boldsymbol{\beta} = \left( (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' - \mathbf{A}\mathbf{X}' \right) \mathbf{e}$
3.  $\tilde{\mathbf{e}} = (\mathbf{I} - \mathbf{P} + \mathbf{X}\mathbf{A}\mathbf{X}') \mathbf{e}$
4.  $\mathbf{I} - \mathbf{P} + \mathbf{X}\mathbf{A}\mathbf{X}'$  is symmetric and idempotent
5.  $\text{tr}(\mathbf{I} - \mathbf{P} + \mathbf{X}\mathbf{A}\mathbf{X}') = n - k + q.$

See Exercise 8.6.

Given the linearity of Theorem 8.4.1.2, it is not hard to show that the CLS estimator is unbiased for  $\boldsymbol{\beta}$

**Theorem 8.4.2** In the linear regression model (8.15)-(8.16) under 8.6.1,  $\mathbb{E}(\tilde{\boldsymbol{\beta}}_{\text{cls}} | \mathbf{X}) = \boldsymbol{\beta}.$

See Exercise 8.7.

Given the linearity we can also calculate the variance matrix of  $\tilde{\boldsymbol{\beta}}_{\text{cls}}$ . For this we will add the assumption of conditional homoskedasticity to simplify the expression.

**Theorem 8.4.3** In the homoskedastic linear regression model (8.15)-(8.16) with  $\mathbb{E}(e_i^2 | \mathbf{x}_i) = \sigma^2$ , under 8.6.1,

$$\begin{aligned} \mathbf{V}_{\tilde{\boldsymbol{\beta}}}^0 &= \text{var}(\tilde{\boldsymbol{\beta}}_{\text{cls}} | \mathbf{X}) \\ &= \left( (\mathbf{X}'\mathbf{X})^{-1} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R} \left( \mathbf{R}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R} \right)^{-1} \mathbf{R}' (\mathbf{X}'\mathbf{X})^{-1} \right) \sigma^2 \end{aligned}$$

See Exercise 8.8. We use the  $\mathbf{V}_{\tilde{\boldsymbol{\beta}}}^0$  notation to emphasize that this is the variance matrix under the assumption of conditional homoskedasticity.

For inference we need an estimate of  $\mathbf{V}_{\tilde{\boldsymbol{\beta}}}^0$ . A natural estimator is

$$\hat{\mathbf{V}}_{\tilde{\boldsymbol{\beta}}}^0 = \left( (\mathbf{X}'\mathbf{X})^{-1} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R} \left( \mathbf{R}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R} \right)^{-1} \mathbf{R}' (\mathbf{X}'\mathbf{X})^{-1} \right) s_{\text{cls}}^2$$

where

$$s_{\text{cls}}^2 = \frac{1}{n - k + q} \sum_{i=1}^n \tilde{e}_i^2 \quad (8.17)$$

is a biased-corrected estimator of  $\sigma^2$ . Standard errors for the components of  $\beta$  are then found by taking the squares roots of the diagonal elements of  $\hat{\mathbf{V}}_{\tilde{\beta}}$ , for example

$$s(\hat{\beta}_j) = \sqrt{[\hat{\mathbf{V}}_{\tilde{\beta}}^0]_{jj}}.$$

The estimator (8.17) has the property that it is unbiased for  $\sigma^2$  under conditional homoskedasticity. To see this, using the properties of Theorem 8.4.1,

$$\begin{aligned} (n - k + q) s_{\text{cls}}^2 &= \tilde{\mathbf{e}}' \tilde{\mathbf{e}} \\ &= \mathbf{e}' (\mathbf{I} - \mathbf{P} + \mathbf{XAX}') (\mathbf{I} - \mathbf{P} + \mathbf{XAX}') \mathbf{e} \\ &= \mathbf{e}' (\mathbf{I} - \mathbf{P} + \mathbf{XAX}') \mathbf{e}. \end{aligned} \quad (8.18)$$

We defer the remainder of the proof to Exercise 8.9.

**Theorem 8.4.4** *In the homoskedastic linear regression model (8.15)-(8.16) with  $\mathbb{E}(e_i^2 | \mathbf{x}_i) = \sigma^2$ , under 8.6.1,  $\mathbb{E}(s_{\text{cls}}^2 | \mathbf{X}) = \sigma^2$  and  $\mathbb{E}(\hat{\mathbf{V}}_{\tilde{\beta}}^0 | \mathbf{X}) = \mathbf{V}_{\tilde{\beta}}^0$ .*

Now consider the distributional properties in the normal regression model

$$\begin{aligned} y_i &= \mathbf{x}_i' \beta + e_i \\ e_i &\sim N(0, \sigma^2). \end{aligned}$$

By the linearity of Theorem 8.4.1.2, conditional on  $\mathbf{X}$ ,  $\tilde{\beta}_{\text{cls}} - \beta$  is normal. Given Theorems 8.4.2 and 8.4.3, we deduce that  $\tilde{\beta}_{\text{cls}} \sim N(\beta, \mathbf{V}_{\tilde{\beta}}^0)$ .

Similarly, from Exercise 8.4.1 we know  $\tilde{\mathbf{e}} = (\mathbf{I} - \mathbf{P} + \mathbf{XAX}') \mathbf{e}$  is linear in  $\mathbf{e}$  so is also conditionally normal. Furthermore, since  $(\mathbf{I} - \mathbf{P} + \mathbf{XAX}') (\mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} - \mathbf{XA}) = 0$ ,  $\tilde{\mathbf{e}}$  and  $\tilde{\beta}_{\text{cls}}$  are uncorrelated and thus independent. Thus  $s_{\text{cls}}^2$  and  $\tilde{\beta}_{\text{cls}}$  are independent.

From (8.18) and the fact that  $\mathbf{I} - \mathbf{P} + \mathbf{XAX}'$  is idempotent with rank  $n - k + q$ , it follows that

$$s_{\text{cls}}^2 \sim \sigma^2 \chi_{n-k+q}^2 / (n - k + q).$$

It follows that the t-statistic has the exact distribution

$$\begin{aligned} T &= \frac{\hat{\beta}_j - \beta_j}{s(\hat{\beta}_j)} \\ &\sim \frac{N(0, 1)}{\sqrt{\chi_{n-k+q}^2 / (n - k + q)}} \\ &\sim t_{n-k+q} \end{aligned}$$

a student  $t$  distribution with  $n - k + q$  degrees of freedom.

The relevance of this calculation is that the “degrees of freedom” for a CLS regression problem equal  $n - k + q$  rather than  $n - k$  as in the OLS regression problem. Essentially, the model has  $k - q$  free parameters instead of  $k$ . Another way of thinking about this is that estimation of a model with  $k$  coefficients and  $q$  restrictions is equivalent to estimation with  $k - q$  coefficients.

We summarize the properties of the normal regression model

**Theorem 8.4.5** *In the normal linear regression model linear regression model (8.15)-(8.16), under 8.6.1,*

$$\begin{aligned}\tilde{\beta}_{\text{cls}} &\sim N(\beta, \mathbf{V}_{\tilde{\beta}}^0) \\ \frac{(n-k+q)s_{\text{cls}}^2}{\sigma^2} &\sim \chi_{n-k+q}^2 \\ T &\sim t_{n-k+q}\end{aligned}$$

An interesting relationship is that in the homoskedastic regression model

$$\begin{aligned}(\hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{cls}}, \tilde{\beta}_{\text{cls}}) &= \mathbb{E} \left( (\hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{cls}}) (\tilde{\beta}_{\text{cls}} - \beta)' \right) \\ &= \mathbb{E} \left( (\mathbf{A}\mathbf{X}') \left( \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} - \mathbf{X}\mathbf{A} \right) \right) \sigma^2 = 0\end{aligned}$$

so  $\hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{cls}}$  and  $\tilde{\beta}_{\text{cls}}$  are uncorrelated and hence independent. One corollary is

$$\text{cov}(\hat{\beta}_{\text{ols}}, \tilde{\beta}_{\text{cls}}) = \text{var}(\tilde{\beta}_{\text{cls}})$$

A second corollary is

$$\begin{aligned}\text{var}(\hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{cls}}) &= \text{var}(\hat{\beta}_{\text{ols}}) - \text{var}(\tilde{\beta}_{\text{cls}}) \\ &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R} \left( \mathbf{R}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R} \right)^{-1} \mathbf{R}' (\mathbf{X}'\mathbf{X})^{-1} \sigma^2.\end{aligned}\tag{8.19}$$

This also shows us the difference between the CLS and OLS variances

$$\text{var}(\hat{\beta}_{\text{ols}}) - \text{var}(\tilde{\beta}_{\text{cls}}) = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R} \left( \mathbf{R}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R} \right)^{-1} \mathbf{R}' (\mathbf{X}'\mathbf{X})^{-1} \sigma^2 \geq 0$$

the final equality meaning positive semi-definite. It follows that  $\text{var}(\hat{\beta}_{\text{ols}}) \geq \text{var}(\tilde{\beta}_{\text{cls}})$  in the positive definite sense, and thus CLS is more efficient than OLS. Both estimators are unbiased (in the linear regression model), and CLS has a lower variance matrix (in the linear homoskedastic regression model).

The relationship (8.19) is rather interesting and will appear again. The expression says that the variance of the difference between the estimators is equal to the difference between the variances. This is rather special. It occurs (generically) when we are comparing an efficient and an inefficient estimator. We call (8.19) the **Hausmann Equality** as it was first pointed out in econometrics by Hausman (1978).

## 8.5 Minimum Distance

The previous section explored the finite sample distribution theory under the assumptions of the linear regression model, homoskedastic regression model, and normal regression model. We now return to the general projection model where we do not impose linearity, homoskedasticity, nor normality. We are interested in the question: Can we do better than CLS in this setting?

A minimum distance estimator tries to find a parameter value which satisfies the constraint which is as close as possible to the unconstrained estimate. Let  $\hat{\beta}$  be the unconstrained least-squares estimator, and for some  $k \times k$  positive definite weight matrix  $\widehat{\mathbf{W}} > 0$  define the quadratic criterion function

$$J(\beta) = n (\hat{\beta} - \beta)' \widehat{\mathbf{W}} (\hat{\beta} - \beta).\tag{8.20}$$

This is a (squared) weighted Euclidean distance between  $\hat{\beta}$  and  $\beta$ .  $J(\beta)$  is small if  $\beta$  is close to  $\hat{\beta}$ , and is minimized at zero only if  $\beta = \hat{\beta}$ . A **minimum distance estimator**  $\tilde{\beta}_{\text{md}}$  for  $\beta$  minimizes  $J(\beta)$  subject to the constraint (8.1), that is,

$$\tilde{\beta}_{\text{md}} = \underset{\mathbf{R}'\beta = \mathbf{c}}{\operatorname{argmin}} J(\beta). \quad (8.21)$$

The CLS estimator is the special case when  $\widehat{\mathbf{W}} = \widehat{\mathbf{Q}}_{xx}$ , and we write this criterion function as

$$J^0(\beta) = n(\hat{\beta} - \beta)' \widehat{\mathbf{Q}}_{xx}(\hat{\beta} - \beta). \quad (8.22)$$

To see the equality of CLS and minimum distance, rewrite the least-squares criterion as follows. Write the unconstrained least-squares fitted equation as  $y_i = \mathbf{x}_i' \hat{\beta} + \hat{e}_i$  and substitute this equation into  $SSE(\beta)$  to obtain

$$\begin{aligned} SSE(\beta) &= \sum_{i=1}^n (y_i - \mathbf{x}_i' \beta)^2 \\ &= \sum_{i=1}^n (\mathbf{x}_i' \hat{\beta} + \hat{e}_i - \mathbf{x}_i' \beta)^2 \\ &= \sum_{i=1}^n \hat{e}_i^2 + (\hat{\beta} - \beta)' \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right) (\hat{\beta} - \beta) \\ &= n\hat{\sigma}^2 + J^0(\beta) \end{aligned} \quad (8.23)$$

where the third equality uses the fact that  $\sum_{i=1}^n \mathbf{x}_i \hat{e}_i = 0$ , and the last line uses  $\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' = n\widehat{\mathbf{Q}}_{xx}$ . The expression (8.23) only depends on  $\beta$  through  $J^0(\beta)$ . Thus minimization of  $SSE(\beta)$  and  $J^0(\beta)$  are equivalent, and hence  $\tilde{\beta}_{\text{md}} = \tilde{\beta}_{\text{cls}}$  when  $\widehat{\mathbf{W}} = \widehat{\mathbf{Q}}_{xx}$ .

We can solve for  $\tilde{\beta}_{\text{md}}$  explicitly by the method of Lagrange multipliers. The Lagrangian is

$$\mathcal{L}(\beta, \lambda) = \frac{1}{2} J(\beta, \widehat{\mathbf{W}}) + \lambda' (\mathbf{R}'\beta - \mathbf{c})$$

which is minimized over  $(\beta, \lambda)$ . The solution is

$$\tilde{\lambda}_{\text{md}} = n \left( \mathbf{R}' \widehat{\mathbf{W}}^{-1} \mathbf{R} \right)^{-1} \left( \mathbf{R}' \hat{\beta} - \mathbf{c} \right) \quad (8.24)$$

$$\tilde{\beta}_{\text{md}} = \hat{\beta} - \widehat{\mathbf{W}}^{-1} \mathbf{R} \left( \mathbf{R}' \widehat{\mathbf{W}}^{-1} \mathbf{R} \right)^{-1} \left( \mathbf{R}' \hat{\beta} - \mathbf{c} \right). \quad (8.25)$$

(See Exercise 8.10.) Comparing (8.25) with (8.10) we can see that  $\tilde{\beta}_{\text{md}}$  specializes to  $\tilde{\beta}_{\text{cls}}$  when we set  $\widehat{\mathbf{W}} = \widehat{\mathbf{Q}}_{xx}$ .

An obvious question is which weight matrix  $\widehat{\mathbf{W}}$  is best. We will address this question after we derive the asymptotic distribution for a general weight matrix.

## 8.6 Asymptotic Distribution

We first show that the class of minimum distance estimators are consistent for the population parameters when the constraints are valid.

**Assumption 8.6.1**  $\mathbf{R}'\beta = \mathbf{c}$  where  $\mathbf{R}$  is  $k \times q$  with  $\operatorname{rank}(\mathbf{R}) = q$ .



**Assumption 8.6.2**  $\widehat{\mathbf{W}} \xrightarrow{p} \mathbf{W} > 0$ .

**Theorem 8.6.1 Consistency**

Under Assumptions 7.1.1, 8.6.1, and 8.6.2,  $\tilde{\beta}_{\text{md}} \xrightarrow{p} \beta$  as  $n \rightarrow \infty$ .

For a proof, see Exercise 8.11.

Theorem 8.6.1 shows that consistency holds for any weight matrix with a positive definite limit, so the result includes the CLS estimator.

Similarly, the constrained estimators are asymptotically normally distributed.

**Theorem 8.6.2 Asymptotic Normality**

Under Assumptions 7.1.2, 8.6.1, and 8.6.2,

$$\sqrt{n} \left( \tilde{\beta}_{\text{md}} - \beta \right) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_{\beta}(\mathbf{W})) \quad (8.26)$$

as  $n \rightarrow \infty$ , where

$$\begin{aligned} \mathbf{V}_{\beta}(\mathbf{W}) = & \mathbf{V}_{\beta} - \mathbf{W}^{-1} \mathbf{R} (\mathbf{R}' \mathbf{W}^{-1} \mathbf{R})^{-1} \mathbf{R}' \mathbf{V}_{\beta} \\ & - \mathbf{V}_{\beta} \mathbf{R} (\mathbf{R}' \mathbf{W}^{-1} \mathbf{R})^{-1} \mathbf{R}' \mathbf{W}^{-1} \\ & + \mathbf{W}^{-1} \mathbf{R} (\mathbf{R}' \mathbf{W}^{-1} \mathbf{R})^{-1} \mathbf{R}' \mathbf{V}_{\beta} \mathbf{R} (\mathbf{R}' \mathbf{W}^{-1} \mathbf{R})^{-1} \mathbf{R}' \mathbf{W}^{-1} \end{aligned} \quad (8.27)$$

and  $\mathbf{V}_{\beta} = \mathbf{Q}_{xx}^{-1} \Omega \mathbf{Q}_{xx}^{-1}$

For a proof, see Exercise 8.12.

Theorem 8.6.2 shows that the minimum distance estimator is asymptotically normal for all positive definite weight matrices. The asymptotic variance depends on  $\mathbf{W}$ . The theorem includes the CLS estimator as a special case by setting  $\mathbf{W} = \mathbf{Q}_{xx}$ .

**Theorem 8.6.3 Asymptotic Distribution of CLS Estimator**

Under Assumptions 7.1.2 and 8.6.1, as  $n \rightarrow \infty$

$$\sqrt{n} \left( \tilde{\beta}_{\text{cls}} - \beta \right) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_{\text{cls}})$$

where

$$\begin{aligned} \mathbf{V}_{\text{cls}} = & \mathbf{V}_{\beta} - \mathbf{Q}_{xx}^{-1} \mathbf{R} (\mathbf{R}' \mathbf{Q}_{xx}^{-1} \mathbf{R})^{-1} \mathbf{R}' \mathbf{V}_{\beta} \\ & - \mathbf{V}_{\beta} \mathbf{R} (\mathbf{R}' \mathbf{Q}_{xx}^{-1} \mathbf{R})^{-1} \mathbf{R}' \mathbf{Q}_{xx}^{-1} \\ & + \mathbf{Q}_{xx}^{-1} \mathbf{R} (\mathbf{R}' \mathbf{Q}_{xx}^{-1} \mathbf{R})^{-1} \mathbf{R}' \mathbf{V}_{\beta} \mathbf{R} (\mathbf{R}' \mathbf{Q}_{xx}^{-1} \mathbf{R})^{-1} \mathbf{R}' \mathbf{Q}_{xx}^{-1} \end{aligned}$$

For a proof, see Exercise 8.13.

## 8.7 Variance Estimation and Standard Errors

Earlier we introduced the covariance matrix estimator under the assumption of conditional homoskedasticity. We now introduce an estimator which does not impose homoskedasticity.

The asymptotic covariance matrix  $\mathbf{V}_{\text{cls}}$  may be estimated by replacing  $\mathbf{V}_\beta$  with a consistent estimates such as  $\hat{\mathbf{V}}_\beta$ . A more efficient estimate is obtained by using the restricted estimates. Given the constrained least-squares residuals  $\tilde{e}_i = y_i - \mathbf{x}'_i \tilde{\boldsymbol{\beta}}_{\text{cls}}$  we can estimate the matrix  $\boldsymbol{\Omega} = \mathbb{E}(\mathbf{x}_i \mathbf{x}'_i e_i^2)$  by

$$\tilde{\boldsymbol{\Omega}} = \frac{1}{n - k + q} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}'_i \tilde{e}_i^2.$$

Notice that we have defined  $\tilde{\boldsymbol{\Omega}}$  using an adjusted degrees of freedom. This is an ad hoc adjustment designed to mimic that used for estimation of the error variance  $\sigma^2$ . Given  $\tilde{\boldsymbol{\Omega}}$  the moment estimator of  $\mathbf{V}_\beta$  is

$$\tilde{\mathbf{V}}_\beta = \hat{\mathbf{Q}}_{xx}^{-1} \tilde{\boldsymbol{\Omega}} \hat{\mathbf{Q}}_{xx}^{-1}$$

and that for  $\mathbf{V}_{\text{cls}}$  is

$$\begin{aligned} \tilde{\mathbf{V}}_{\text{cls}} &= \tilde{\mathbf{V}}_\beta - \hat{\mathbf{Q}}_{xx}^{-1} \mathbf{R} \left( \mathbf{R}' \hat{\mathbf{Q}}_{xx}^{-1} \mathbf{R} \right)^{-1} \mathbf{R}' \tilde{\mathbf{V}}_\beta \\ &\quad - \tilde{\mathbf{V}}_\beta \mathbf{R} \left( \mathbf{R}' \hat{\mathbf{Q}}_{xx}^{-1} \mathbf{R} \right)^{-1} \mathbf{R}' \hat{\mathbf{Q}}_{xx}^{-1} \\ &\quad + \hat{\mathbf{Q}}_{xx}^{-1} \mathbf{R} \left( \mathbf{R}' \hat{\mathbf{Q}}_{xx}^{-1} \mathbf{R} \right)^{-1} \mathbf{R}' \tilde{\mathbf{V}}_\beta \mathbf{R} \left( \mathbf{R}' \hat{\mathbf{Q}}_{xx}^{-1} \mathbf{R} \right)^{-1} \mathbf{R}' \hat{\mathbf{Q}}_{xx}^{-1}. \end{aligned}$$

We can calculate standard errors for any linear combination  $\mathbf{h}' \tilde{\boldsymbol{\beta}}_{\text{cls}}$  so long as  $\mathbf{h}$  does not lie in the range space of  $\mathbf{R}$ . A standard error for  $\mathbf{h}' \tilde{\boldsymbol{\beta}}$  is

$$s(\mathbf{h}' \tilde{\boldsymbol{\beta}}_{\text{cls}}) = \left( n^{-1} \mathbf{h}' \tilde{\mathbf{V}}_{\text{cls}} \mathbf{h} \right)^{1/2}.$$

## 8.8 Efficient Minimum Distance Estimator

Theorem 8.6.2 shows that the minimum distance estimators, which include CLS as a special case, are asymptotically normal with an asymptotic covariance matrix which depends on the weight matrix  $\mathbf{W}$ . The asymptotically optimal weight matrix is the one which minimizes the asymptotic variance  $\mathbf{V}_\beta(\mathbf{W})$ . This turns out to be  $\mathbf{W} = \mathbf{V}_\beta^{-1}$  as is shown in Theorem 8.8.1 below. Since  $\mathbf{V}_\beta^{-1}$  is unknown this weight matrix cannot be used for a feasible estimator, but we can replace  $\mathbf{V}_\beta^{-1}$  with a consistent estimate  $\hat{\mathbf{V}}_\beta^{-1}$  and the asymptotic distribution (and efficiency) are unchanged. We call the minimum distance estimator setting  $\hat{\mathbf{W}} = \hat{\mathbf{V}}_\beta^{-1}$  the **efficient minimum distance estimator** and takes the form

$$\tilde{\boldsymbol{\beta}}_{\text{emd}} = \hat{\boldsymbol{\beta}} - \hat{\mathbf{V}}_\beta \mathbf{R} \left( \mathbf{R}' \hat{\mathbf{V}}_\beta \mathbf{R} \right)^{-1} \left( \mathbf{R}' \hat{\boldsymbol{\beta}} - \mathbf{c} \right). \quad (8.28)$$

The asymptotic distribution of (8.28) can be deduced from Theorem 8.6.2. (See Exercises 8.14 and 8.15.)

**Theorem 8.8.1 Efficient Minimum Distance Estimator***Under Assumptions 7.1.2 and 8.6.1,*

$$\sqrt{n} \left( \tilde{\beta}_{\text{emd}} - \beta \right) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_{\beta, \text{emd}})$$

*as  $n \rightarrow \infty$ , where*

$$\mathbf{V}_{\beta, \text{emd}} = \mathbf{V}_{\beta} - \mathbf{V}_{\beta} \mathbf{R} (\mathbf{R}' \mathbf{V}_{\beta} \mathbf{R})^{-1} \mathbf{R}' \mathbf{V}_{\beta}. \quad (8.29)$$

*Since*

$$\mathbf{V}_{\beta, \text{emd}} \leq \mathbf{V}_{\beta} \quad (8.30)$$

*the estimator (8.28) has lower asymptotic variance than the unrestricted estimator. Furthermore, for any  $\mathbf{W}$ ,*

$$\mathbf{V}_{\beta, \text{emd}} \leq \mathbf{V}_{\beta}(\mathbf{W}) \quad (8.31)$$

*so (8.28) is asymptotically efficient in the class of minimum distance estimators.*

Theorem 8.8.1 shows that the minimum distance estimator with the smallest asymptotic variance is (8.28). One implication is that the constrained least squares estimator is generally inefficient. The interesting exception is the case of conditional homoskedasticity, in which case the optimal weight matrix is  $\mathbf{W} = (\mathbf{V}_{\beta}^0)^{-1}$  so in this case CLS is an efficient minimum distance estimator. Otherwise when the error is conditionally heteroskedastic, there are asymptotic efficiency gains by using minimum distance rather than least squares.

The fact that CLS is generally inefficient is counter-intuitive and requires some reflection to understand. Standard intuition suggests to apply the same estimation method (least squares) to the unconstrained and constrained models, and this is the most common empirical practice. But Theorem 8.8.1 shows that this is not the efficient estimation method. Instead, the efficient minimum distance estimator has a smaller asymptotic variance. Why? The reason is that the least-squares estimator does not make use of the regressor  $\mathbf{x}_{2i}$ . It ignores the information  $\mathbb{E}(\mathbf{x}_{2i}e_i) = \mathbf{0}$ . This information is relevant when the error is heteroskedastic and the excluded regressors are correlated with the included regressors.

Inequality (8.30) shows that the efficient minimum distance estimator  $\tilde{\beta}_{\text{emd}}$  has a smaller asymptotic variance than the unrestricted least squares estimator  $\hat{\beta}$ . This means that estimation is more efficient by imposing correct restrictions when we use the minimum distance method.

## 8.9 Exclusion Restriction Revisited

We return to the example of estimation with a simple exclusion restriction. The model is

$$y_i = \mathbf{x}'_{1i}\beta_1 + \mathbf{x}'_{2i}\beta_2 + e_i$$

with the exclusion restriction  $\beta_2 = \mathbf{0}$ . We have introduced three estimators of  $\beta_1$ . The first is unconstrained least-squares applied to (8.11), which can be written as

$$\hat{\beta}_1 = \hat{\mathbf{Q}}_{11 \cdot 2}^{-1} \hat{\mathbf{Q}}_{1y \cdot 2}.$$

From Theorem 7.33 and equation (7.20) its asymptotic variance is

$$\text{avar}(\hat{\beta}_1) = \mathbf{Q}_{11 \cdot 2}^{-1} (\mathbf{\Omega}_{11} - \mathbf{Q}_{12} \mathbf{Q}_{22}^{-1} \mathbf{\Omega}_{21} - \mathbf{\Omega}_{12} \mathbf{Q}_{22}^{-1} \mathbf{Q}_{21} + \mathbf{Q}_{12} \mathbf{Q}_{22}^{-1} \mathbf{\Omega}_{22} \mathbf{Q}_{22}^{-1} \mathbf{Q}_{21}) \mathbf{Q}_{11 \cdot 2}^{-1}.$$

The second estimator of  $\beta_1$  is the CLS estimator, which can be written as

$$\tilde{\beta}_{1,\text{cls}} = \hat{\mathbf{Q}}_{11}^{-1} \hat{\mathbf{Q}}_{1y}.$$

Its asymptotic variance can be deduced from Theorem 8.6.3, but it is simpler to apply the CLT directly to show that

$$\text{avar}(\tilde{\beta}_{1,\text{cls}}) = \mathbf{Q}_{11}^{-1} \Omega_{11} \mathbf{Q}_{11}^{-1}. \quad (8.32)$$

The third estimator of  $\beta_1$  is the efficient minimum distance estimator. Applying (8.28), it equals

$$\tilde{\beta}_{1,\text{md}} = \hat{\beta}_1 - \hat{\mathbf{V}}_{12} \hat{\mathbf{V}}_{22}^{-1} \hat{\beta}_2 \quad (8.33)$$

where we have partitioned

$$\hat{\mathbf{V}}_{\beta} = \begin{bmatrix} \hat{\mathbf{V}}_{11} & \hat{\mathbf{V}}_{12} \\ \hat{\mathbf{V}}_{21} & \hat{\mathbf{V}}_{22} \end{bmatrix}.$$

From Theorem 8.8.1 its asymptotic variance is

$$\text{avar}(\tilde{\beta}_{1,\text{md}}) = \mathbf{V}_{11} - \mathbf{V}_{12} \mathbf{V}_{22}^{-1} \mathbf{V}_{21}. \quad (8.34)$$

See Exercise 8.16 to verify equations (8.32), (8.33), and (8.34).

In general, the three estimators are different, and they have different asymptotic variances.

It is quite instructive to compare the asymptotic variances of the CLS and unconstrained least-squares estimators to assess whether or not the constrained estimator is necessarily more efficient than the unconstrained estimator.

First, consider the case of conditional homoskedasticity. In this case the two covariance matrices simplify to

$$\text{avar}(\hat{\beta}_1) = \sigma^2 \mathbf{Q}_{11 \cdot 2}^{-1}$$

and

$$\text{avar}(\tilde{\beta}_{1,\text{cls}}) = \sigma^2 \mathbf{Q}_{11}^{-1}.$$

If  $\mathbf{Q}_{12} = 0$  (so  $\mathbf{x}_{1i}$  and  $\mathbf{x}_{2i}$  are orthogonal) then these two variance matrices are equal and the two estimators have equal asymptotic efficiency. Otherwise, since  $\mathbf{Q}_{12} \mathbf{Q}_{22}^{-1} \mathbf{Q}_{21} \geq 0$ , then  $\mathbf{Q}_{11} \geq \mathbf{Q}_{11} - \mathbf{Q}_{12} \mathbf{Q}_{22}^{-1} \mathbf{Q}_{21}$ , and consequently

$$\mathbf{Q}_{11}^{-1} \sigma^2 \leq (\mathbf{Q}_{11} - \mathbf{Q}_{12} \mathbf{Q}_{22}^{-1} \mathbf{Q}_{21})^{-1} \sigma^2.$$

This means that under conditional homoskedasticity,  $\tilde{\beta}_{1,\text{cls}}$  has a lower asymptotic variance matrix than  $\hat{\beta}_1$ . Therefore in this context, constrained least-squares is more efficient than unconstrained least-squares. This is consistent with our intuition that imposing a correct restriction (excluding an irrelevant regressor) improves estimation efficiency.

However, in the general case of conditional heteroskedasticity this ranking is not guaranteed. In fact what is really amazing is that the variance ranking can be reversed. The CLS estimator can have a larger asymptotic variance than the unconstrained least squares estimator.

To see this let's use the simple heteroskedastic example from Section 7.4. In that example,  $\mathbf{Q}_{11} = \mathbf{Q}_{22} = 1$ ,  $\mathbf{Q}_{12} = \frac{1}{2}$ ,  $\Omega_{11} = \Omega_{22} = 1$ , and  $\Omega_{12} = \frac{7}{8}$ . We can calculate (see Exercise 8.17) that  $\mathbf{Q}_{11 \cdot 2} = \frac{3}{4}$  and

$$\text{avar}(\hat{\beta}_1) = \frac{2}{3} \quad (8.35)$$

$$\text{avar}(\tilde{\beta}_{1,\text{cls}}) = 1 \quad (8.36)$$

$$\text{avar}(\tilde{\beta}_{1,\text{md}}) = \frac{5}{8}. \quad (8.37)$$

Thus the restricted least-squares estimator  $\tilde{\beta}_{1,\text{cls}}$  has a larger variance than the unrestricted least-squares estimator  $\hat{\beta}_1$ ! The minimum distance estimator has the smallest variance of the three, as expected.

What we have found is that when the estimation method is least-squares, deleting the irrelevant variable  $x_{2i}$  can actually increase estimation variance, or equivalently, adding an irrelevant variable can actually decrease the estimation variance.

To repeat this unexpected finding, we have shown in a very simple example that it is possible for least-squares applied to the short regression (8.12) to be less efficient for estimation of  $\beta_1$  than least-squares applied to the long regression (8.11), even though the constraint  $\beta_2 = 0$  is valid! This result is strongly counter-intuitive. It seems to contradict our initial motivation for pursuing constrained estimation – to improve estimation efficiency.

It turns out that a more refined answer is appropriate. Constrained estimation is desirable, but not constrained least-squares estimation. While least-squares is asymptotically efficient for estimation of the unconstrained projection model, it is not an efficient estimator of the constrained projection model.

## 8.10 Variance and Standard Error Estimation

We have discussed covariance matrix estimation for the CLS estimator, but not yet for the EMD estimator.

The asymptotic covariance matrix (8.29) may be estimated by replacing  $\mathbf{V}_\beta$  with a consistent estimate. It is best to construct the variance estimate using  $\tilde{\beta}_{\text{emd}}$ . The EMD residuals are  $\tilde{e}_i = y_i - \mathbf{x}_i' \tilde{\beta}_{\text{emd}}$ . Using these we can estimate the matrix  $\Omega = \mathbb{E}(\mathbf{x}_i \mathbf{x}_i' e_i^2)$  by

$$\tilde{\Omega} = \frac{1}{n - k + q} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \tilde{e}_i^2.$$

Following the formula for CLS we recommend an adjusted degrees of freedom. Given  $\tilde{\Omega}$  the moment estimator of  $\mathbf{V}_\beta$  is

$$\tilde{\mathbf{V}}_\beta = \hat{\mathbf{Q}}_{xx}^{-1} \tilde{\Omega} \hat{\mathbf{Q}}_{xx}^{-1}$$

Given this, we construct the variance estimator

$$\tilde{\mathbf{V}}_{\beta,\text{emd}} = \tilde{\mathbf{V}}_\beta - \tilde{\mathbf{V}}_\beta \mathbf{R} \left( \mathbf{R}' \tilde{\mathbf{V}}_\beta \mathbf{R} \right)^{-1} \mathbf{R}' \tilde{\mathbf{V}}_\beta. \quad (8.38)$$

A standard error for  $\mathbf{h}' \tilde{\beta}$  is then

$$s(\mathbf{h}' \tilde{\beta}) = \left( n^{-1} \mathbf{h}' \tilde{\mathbf{V}}_{\beta,\text{emd}} \mathbf{h} \right)^{1/2}. \quad (8.39)$$

## 8.11 Hausman Equality

Form (8.28) we have

$$\begin{aligned} \sqrt{n} \left( \hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{emd}} \right) &= \hat{\mathbf{V}}_\beta \mathbf{R} \left( \mathbf{R}' \hat{\mathbf{V}}_\beta \mathbf{R} \right)^{-1} \sqrt{n} \left( \mathbf{R}' \hat{\beta}_{\text{ols}} - \mathbf{c} \right) \\ &\xrightarrow{d} \mathbf{N} \left( \mathbf{0}, \mathbf{V}_\beta \mathbf{R} \left( \mathbf{R}' \mathbf{V}_\beta \mathbf{R} \right)^{-1} \mathbf{R}' \mathbf{V}_\beta \right). \end{aligned}$$

It follows that the asymptotic variances of the estimators satisfy the relationship

$$\text{avar} \left( \hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{emd}} \right) = \text{avar} \left( \hat{\beta}_{\text{ols}} \right) - \text{avar} \left( \tilde{\beta}_{\text{emd}} \right). \quad (8.40)$$

We call (8.40) the Hausman Equality: the asymptotic variance of the difference between an efficient and inefficient estimator is the difference in the asymptotic variances.

## 8.12 Example: Mankiw, Romer and Weil (1992)

We illustrate the methods by replicating some of the estimates reported in a well-known paper by Mankiw, Romer, and Weil (1992). The paper investigates the implications of the Solow growth model using cross-country regressions. A key equation in their paper regresses the change between 1960 and 1985 in log GDP per capita on (1) log GDP in 1960, (2) the log of the ratio of aggregate investment to GDP, (3) the log of the sum of the population growth rate  $n$ , the technological growth rate  $g$ , and the rate of depreciation  $\delta$ , and (4) the log of the percentage of the working-age population that is in secondary school (*School*), the latter a proxy for human-capital accumulation.

The data is available on the textbook webpage in the file **MRW1992**.

The sample is 98 non-oil-producing countries, and the data was reported in the published paper. As  $g$  and  $\delta$  were unknown the authors set  $g + \delta = 0.05$ . We report least-squares estimates in the first column of the table below, using the authors' original data. The estimates are consistent with the Solow theory due to the positive coefficients on investment and human capital and negative coefficient for population growth. The estimates are also consistent with the convergence hypothesis (that income levels tend towards a common mean over time) as the coefficient on initial GDP is negative.

The authors show that in the Solow model the 2<sup>nd</sup>, 3<sup>rd</sup> and 4<sup>th</sup> coefficients sum to zero. They reestimated the equation imposing this constraint. We present constrained least-squares estimates in the second column, and efficient minimum distance estimates in the third column. Most of the coefficients and standard errors only exhibit small changes by imposing the constraint. The one exception is the coefficient on log population growth, which increases in magnitude and its standard error decreases substantially. The differences between the CLS and EMD estimates are modest but not inconsequential.

Table  
Estimates of Solow Growth Model  
Dependent Variable  $\log \frac{GDP_{1985}}{GDP_{1960}}$

	$\hat{\beta}_{OLS}$	$\hat{\beta}_{CLS}$	$\hat{\beta}_{EMD}$
$\log GDP_{1960}$	-0.29 (0.05)	-0.30 (0.05)	-0.30 (0.05)
$\log \frac{I}{GDP}$	0.52 (0.11)	0.50 (0.09)	0.46 (0.08)
$\log (n + g + \delta)$	-0.51 (0.25)	-0.74 (0.08)	-0.71 (0.08)
$\log School$	0.23 (0.07)	0.24 (0.07)	0.25 (0.07)
Intercept	3.02 (0.74)	2.46 (0.44)	2.48 (0.44)

Note: Standard errors are heteroskedasticity-consistent

We now present Stata, R and MATLAB code which implements these estimates.

You may notice that the Stata code has a section which uses the Mata matrix programming language. This is used because Stata does not implement the efficient minimum distance estimator, so needs to be separately programmed. As illustrated here, the Mata language allows a Stata user to implement methods using commands which are quite similar to MATLAB.

#### Stata do File

```

use "MRW1992.dta", clear
gen lndY = log(Y85)-log(Y60)
gen lnY60 = log(Y60)
gen lnI = log(invest/100)
gen lnG = log(pop_growth/100+0.05)
gen lnS = log(school/100)
// Unrestricted regression
reg lndY lnY60 lnI lnG lnS if N==1, r
// Store result for efficient minimum distance
mat b = e(b)'
scalar k = e(rank)
mat V = e(V)
// Constrained regression
constraint define 1 lnI+lnG+lnS=0
cnsreg lndY lnY60 lnI lnG lnS if N==1, constraints(1) r
// Efficient minimum distance
mata{
    data = st_data(.,("lnY60","lnI","lnG","lnS","lndY","N"))
    data_select = select(data,data[.,6]==1)
    y = data_select[.,5]
    n = rows(y)
    x = (data_select[.,1..4],J(n,1,1))
    k = cols(x)
    invx = invsym(x'*x)
    b_ols = st_matrix("b")
    V_ols = st_matrix("V")
    R = (0\1\1\1\0)
    b_emd = b_ols-V_ols*R*invsym(R'*V_ols*R)*R'*b_ols
    e_emd = J(1,k,y-x*b_emd)
    xe_emd = x:e_emd
    xe_emd'*xe_emd
    V2 = (n/(n-k+1))*invx*(xe_emd'*xe_emd)*invx
    V_emd = V2 - V2*R*invsym(R'*V2*R)*R'*V2
    se_emd = diagonal(sqrt(V_emd))
    st_matrix("b_emd",b_emd)
    st_matrix("se_emd",se_emd)}
mat list b_emd
mat list se_emd

```

**R Program File**

```

# Load the data and create variables
data <- read.table("MRW1992.txt",header=TRUE)
N <- matrix(data$N,ncol=1)
lnY <- matrix(log(data$Y85)-log(data$Y60),ncol=1)
lnY60 <- matrix(log(data$Y60),ncol=1)
lnI <- matrix(log(data$invest/100),ncol=1)
lnG <- matrix(log(data$pop_growth/100+0.05),ncol=1)
lnS <- matrix(log(data$school/100),ncol=1)
xx <- as.matrix(cbind(lnY60,lnI,lnG,lnS,matrix(1,nrow(lnY),1)))
x <- xx[N==1,]
y <- lnY[N==1]
n <- nrow(x)
k <- ncol(x)
# Unrestricted regression
invx <- solve(t(x)%*%x)
beta_ols <- invx%*%t(x)%*%y
e_ols <- rep((y-x%*%beta_ols),times=k)
xe_ols <- x*e_ols
V_ols <- (n/(n-k))*invx%*%(t(xe_ols)%*%xe_ols)%*%invx
se_ols <- sqrt(diag(V_ols))
print(beta_ols)
print(se_ols)
# Constrained regression
R <- c(0,1,1,1,0)
iR = invx%*%R%*%solve(t(R)%*%invx%*%R)%*%t(R)
b_cls <- b_ols - iR%*%b_ols
e_cls <- rep((y-x%*%b_cls),times=k)
xe_cls <- x*e_cls
V_tilde <- (n/(n-k+1))*invx%*%(t(xe_cls)%*%xe_cls)%*%invx
V_cls <- V_tilde - iR%*%V_tilde - V_tilde%*%t(iR) +
iR%*%V_tilde%*%t(iR)
print(b_cls)
print(se_cls)
# Efficient minimum distance
Vr = V_ols%*%R%*%solve(t(R)%*%V_ols%*%R)%*%t(R)
b_emd <- b_ols - Vr%*%b_ols
e_emd <- rep((y-x%*%b_emd),times=k)
xe_emd <- x*e_emd
V2 <- (n/(n-k+1))*invx%*%(t(xe_emd)%*%xe_emd)%*%invx
V_emd <- V2 - V2%*%R%*%solve(t(R)%*%V2%*%R)%*%t(R)%*%V2
se_emd <- sqrt(diag(V_emd))

```



**MATLAB Program File**

```

% Load the data and create variables
data = xlsread('MRW1992.xlsx');
N = data(:,1);
Y60 = data(:,4);
Y85 = data(:,5);
pop_growth = data(:,7);
invest = data(:,8);
school = data(:,9);
lnY = log(Y85)-log(Y60);
lnY60 = log(Y60);
lnI = log(invest/100);
lnG = log(pop_growth/100+0.05);
lnS = log(school/100);
xx = [lnY60,lnI,lnG,lnS,ones(size(lnY,1),1)];
x = xx(N==1,:);
y = lnY(N==1);
[n,k] = size(x);
% Unrestricted regression
invx = inv(x'*x);
beta_ols = invx*x'*y;
e_ols = repmat((y-x*beta_ols),1,k);
xe_ols = x.*e_ols;
V_ols = (n/(n-k))*invx*(xe_ols'*xe_ols)*invx;
se_ols = sqrt(diag(V_ols));
display(beta_ols);
display(se_ols);
% Constrained regression
R = [0;1;1;1;0];
iR = invx*R*inv(R'*invx*R)*R';
beta_cls = beta_ols - iR*beta_ols;
e_cls = repmat((y-x*beta_cls),1,k);
xe_cls = x.*e_cls;
V_tilde = (n/(n-k+1))*invx*(xe_cls'*xe_cls)*invx;
V_cls = V_tilde - iR*V_tilde - V_tilde*(iR')...
+ iR*V_tilde*(iR');
se_cls = sqrt(diag(V_cls));
display(beta_cls);
display(se_cls);
% (3) Efficient minimum distance
beta_emd = beta_ols - V_ols*R*inv(R'*V_ols*R)*R'*beta_ols;
e_emd = repmat((y-x*beta_emd),1,k);
xe_emd = x.*e_emd;
V2 = (n/(n-k+1))*invx*(xe_emd'*xe_emd)*invx;
V_emd = V2 - V2*R*inv(R'*V2*R)*R'*V2;
se_emd = sqrt(diag(V_emd));
display(beta_emd);display(se_emd);

```

### 8.13 Misspecification

What are the consequences for a constrained estimator  $\tilde{\beta}$  if the constraint (8.1) is incorrect? To be specific, suppose that

$$\mathbf{R}'\beta = \mathbf{c}^*$$

where  $\mathbf{c}^*$  is not necessarily equal to  $\mathbf{c}$ .

This situation is a generalization of the analysis of “omitted variable bias” from Section 2.23, where we found that the short regression (e.g. (8.13)) is estimating a different projection coefficient than the long regression (e.g. (8.11)).

One mechanical answer is that we can use the formula (8.25) for the minimum distance estimator to find that

$$\tilde{\beta}_{\text{md}} \xrightarrow{p} \beta_{\text{md}}^* = \beta - \mathbf{W}^{-1}\mathbf{R}(\mathbf{R}'\mathbf{W}^{-1}\mathbf{R})^{-1}(\mathbf{c}^* - \mathbf{c}). \quad (8.41)$$

The second term,  $\mathbf{W}^{-1}\mathbf{R}(\mathbf{R}'\mathbf{W}^{-1}\mathbf{R})^{-1}(\mathbf{c}^* - \mathbf{c})$ , shows that imposing an incorrect constraint leads to inconsistency – an asymptotic bias. We can call the limiting value  $\beta_{\text{md}}^*$  the minimum-distance projection coefficient or the pseudo-true value implied by the restriction.

However, we can say more.

For example, we can describe some characteristics of the approximating projections. The CLS estimator projection coefficient has the representation

$$\beta_{\text{cls}}^* = \underset{\mathbf{R}'\beta = \mathbf{c}}{\operatorname{argmin}} \mathbb{E} (y_i - \mathbf{x}_i'\beta)^2,$$

the best linear predictor subject to the constraint (8.1). The minimum distance estimator converges to

$$\beta_{\text{md}}^* = \underset{\mathbf{R}'\beta = \mathbf{c}}{\operatorname{argmin}} (\beta - \beta_0)' \mathbf{W} (\beta - \beta_0)$$

where  $\beta_0$  is the true coefficient. That is,  $\beta_{\text{md}}^*$  is the coefficient vector satisfying (8.1) closest to the true value in the weighted Euclidean norm. These calculations show that the constrained estimators are still reasonable in the sense that they produce good approximations to the true coefficient, conditional on being required to satisfy the constraint.

We can also show that  $\tilde{\beta}_{\text{md}}$  has an asymptotic normal distribution. The trick is to define the pseudo-true value

$$\beta_n^* = \beta - \widehat{\mathbf{W}}^{-1}\mathbf{R}(\mathbf{R}'\widehat{\mathbf{W}}^{-1}\mathbf{R})^{-1}(\mathbf{c}^* - \mathbf{c}). \quad (8.42)$$

(Note that (8.41) and (8.42) are different!) Then

$$\begin{aligned} \sqrt{n}(\tilde{\beta}_{\text{md}} - \beta_n^*) &= \sqrt{n}(\hat{\beta} - \beta) - \widehat{\mathbf{W}}^{-1}\mathbf{R}(\mathbf{R}'\widehat{\mathbf{W}}^{-1}\mathbf{R})^{-1}\sqrt{n}(\mathbf{R}'\hat{\beta} - \mathbf{c}^*) \\ &= \left(\mathbf{I} - \widehat{\mathbf{W}}^{-1}\mathbf{R}(\mathbf{R}'\widehat{\mathbf{W}}^{-1}\mathbf{R})^{-1}\mathbf{R}'\right)\sqrt{n}(\hat{\beta} - \beta) \\ &\xrightarrow{d} \left(\mathbf{I} - \mathbf{W}^{-1}\mathbf{R}(\mathbf{R}'\mathbf{W}^{-1}\mathbf{R})^{-1}\mathbf{R}'\right)\mathbf{N}(\mathbf{0}, \mathbf{V}_\beta) \\ &= \mathbf{N}(\mathbf{0}, \mathbf{V}_\beta(\mathbf{W})). \end{aligned} \quad (8.43)$$

In particular

$$\sqrt{n}(\tilde{\beta}_{\text{emd}} - \beta_n^*) \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{V}_\beta^*).$$

This means that even when the constraint (8.1) is misspecified, the conventional covariance matrix estimator (8.38) and standard errors (8.39) are appropriate measures of the sampling variance, though the distributions are centered at the pseudo-true values (or projections)  $\beta_n^*$  rather than  $\beta$ . The fact that the estimators are biased is an unavoidable consequence of misspecification.

An alternative approach to the asymptotic distribution theory under misspecification uses the concept of local alternatives. It is a technical device which might seem a bit artificial, but it is a powerful method to derive useful distributional approximations in a wide variety of contexts. The idea is to index the true coefficient  $\beta_n$  by  $n$  via the relationship

$$\mathbf{R}'\beta_n = \mathbf{c} + \delta n^{-1/2}. \quad (8.44)$$

Equation (8.44) specifies that  $\beta_n$  violates (8.1) and thus the constraint is misspecified. However, the constraint is “close” to correct, as the difference  $\mathbf{R}'\beta_n - \mathbf{c} = \delta n^{-1/2}$  is “small” in the sense that it decreases with the sample size  $n$ . We call (8.44) **local misspecification**.

The asymptotic theory is then derived as  $n \rightarrow \infty$  under the sequence of probability distributions with the coefficients  $\beta_n$ . The way to think about this is that the true value of the parameter is  $\beta_n$ , and it is “close” to satisfying (8.1). The reason why the deviation is proportional to  $n^{-1/2}$  is because this is the only choice under which the localizing parameter  $\delta$  appears in the asymptotic distribution but does not dominate it. The best way to see this is to work through the asymptotic approximation.

Since  $\beta_n$  is the true coefficient value, then  $y_i = \mathbf{x}_i'\beta_n + e_i$  and we have the standard representation for the unconstrained estimator, namely

$$\begin{aligned} \sqrt{n}(\hat{\beta} - \beta_n) &= \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i e_i \right) \\ &\xrightarrow{d} N(\mathbf{0}, \mathbf{V}_\beta). \end{aligned} \quad (8.45)$$

There is no difference under fixed (classical) or local asymptotics, since the right-hand-side is independent of the coefficient  $\beta_n$ .

A difference arises for the constrained estimator. Using (8.44),  $\mathbf{c} = \mathbf{R}'\beta_n - \delta n^{-1/2}$ , so

$$\mathbf{R}'\hat{\beta} - \mathbf{c} = \mathbf{R}'(\hat{\beta} - \beta_n) + \delta n^{-1/2}$$

and

$$\begin{aligned} \tilde{\beta}_{\text{md}} &= \hat{\beta} - \widehat{\mathbf{W}}^{-1} \mathbf{R} \left( \mathbf{R}' \widehat{\mathbf{W}}^{-1} \mathbf{R} \right)^{-1} \left( \mathbf{R}' \hat{\beta} - \mathbf{c} \right) \\ &= \hat{\beta} - \widehat{\mathbf{W}}^{-1} \mathbf{R} \left( \mathbf{R}' \widehat{\mathbf{W}}^{-1} \mathbf{R} \right)^{-1} \mathbf{R}' (\hat{\beta} - \beta_n) + \widehat{\mathbf{W}}^{-1} \mathbf{R} \left( \mathbf{R}' \widehat{\mathbf{W}}^{-1} \mathbf{R} \right)^{-1} \delta n^{-1/2}. \end{aligned}$$

It follows that

$$\begin{aligned} \sqrt{n}(\tilde{\beta}_{\text{md}} - \beta_n) &= \left( \mathbf{I} - \widehat{\mathbf{W}}^{-1} \mathbf{R} \left( \mathbf{R}' \widehat{\mathbf{W}}^{-1} \mathbf{R} \right)^{-1} \mathbf{R}' \right) \sqrt{n}(\hat{\beta} - \beta_n) \\ &\quad + \widehat{\mathbf{W}}^{-1} \mathbf{R} \left( \mathbf{R}' \widehat{\mathbf{W}}^{-1} \mathbf{R} \right)^{-1} \delta. \end{aligned}$$

The first term is asymptotically normal (from 8.45)). The second term converges in probability to a constant. This is because the  $n^{-1/2}$  local scaling in (8.44) is exactly balanced by the  $\sqrt{n}$  scaling of the estimator. No alternative rate would have produced this result.

Consequently, we find that the asymptotic distribution equals

$$\begin{aligned} \sqrt{n}(\tilde{\beta}_{\text{md}} - \beta_n) &\xrightarrow{d} N(\mathbf{0}, \mathbf{V}_\beta) + \mathbf{W}^{-1} \mathbf{R} \left( \mathbf{R}' \mathbf{W}^{-1} \mathbf{R} \right)^{-1} \delta \\ &= N(\delta^*, \mathbf{V}_\beta(\mathbf{W})) \end{aligned} \quad (8.46)$$

where

$$\delta^* = \mathbf{W}^{-1} \mathbf{R} \left( \mathbf{R}' \mathbf{W}^{-1} \mathbf{R} \right)^{-1} \delta.$$

The asymptotic distribution (8.46) is an approximation of the sampling distribution of the restricted estimator under misspecification. The distribution (8.46) contains an asymptotic bias component  $\delta^*$ . The approximation is not fundamentally different from (8.43) – they both have the same asymptotic variances, and both reflect the bias due to misspecification. The difference is that (8.43) puts the bias on the left-side of the convergence arrow, while (8.46) has the bias on the right-side. There is no substantive difference between the two, but (8.46) is more convenient for some purposes, such as the analysis of the power of tests, as we will explore in the next chapter.

## 8.14 Nonlinear Constraints

In some cases it is desirable to impose nonlinear constraints on the parameter vector  $\beta$ . They can be written as

$$\mathbf{r}(\beta) = \mathbf{0} \quad (8.47)$$

where  $\mathbf{r} : \mathbb{R}^k \rightarrow \mathbb{R}^q$ . This includes the linear constraints (8.1) as a special case. An example of (8.47) which cannot be written as (8.1) is  $\beta_1\beta_2 = 1$ , which is (8.47) with  $\mathbf{r}(\beta) = \beta_1\beta_2 - 1$ .

The constrained least-squares and minimum distance estimators of  $\beta$  subject to (8.47) solve the minimization problems

$$\tilde{\beta}_{\text{cls}} = \underset{\mathbf{r}(\beta)=\mathbf{0}}{\operatorname{argmin}} SSE(\beta) \quad (8.48)$$

$$\tilde{\beta}_{\text{md}} = \underset{\mathbf{r}(\beta)=\mathbf{0}}{\operatorname{argmin}} J(\beta) \quad (8.49)$$

where  $SSE(\beta)$  and  $J(\beta)$  are defined in (8.4) and (8.20), respectively. The solutions minimize the Lagrangians

$$\mathcal{L}(\beta, \lambda) = \frac{1}{2} SSE(\beta) + \lambda' \mathbf{r}(\beta) \quad (8.50)$$

or

$$\mathcal{L}(\beta, \lambda) = \frac{1}{2} J(\beta) + \lambda' \mathbf{r}(\beta) \quad (8.51)$$

over  $(\beta, \lambda)$ .

Computationally, there is no general closed-form solution for the estimator so they must be found numerically. Algorithms to numerically solve (8.48) and (8.49) are known as **constrained optimization** methods, and are available in programming languages including MATLAB, GAUSS and R.

**Assumption 8.14.1**  $\mathbf{r}(\beta) = \mathbf{0}$ ,  $\mathbf{r}(\beta)$  is continuously differentiable at the true  $\beta$ , and  $\operatorname{rank}(\mathbf{R}) = q$ , where  $\mathbf{R} = \frac{\partial}{\partial \beta} \mathbf{r}(\beta)'$ .

The asymptotic distribution is a simple generalization of the case of a linear constraint, but the proof is more delicate.

**Theorem 8.14.1** *Under Assumptions 7.1.2, 8.14.1, and 8.6.2, for  $\tilde{\beta} = \tilde{\beta}_{\text{md}}$  and  $\tilde{\beta} = \tilde{\beta}_{\text{cls}}$  defined in (8.48) and (8.49),*

$$\sqrt{n} \left( \tilde{\beta} - \beta \right) \xrightarrow{d} N(0, V_{\beta}(\mathbf{W}))$$

*as  $n \rightarrow \infty$ , where  $V_{\beta}(\mathbf{W})$  is defined in (8.27). For  $\tilde{\beta}_{\text{cls}}$ ,  $\mathbf{W} = \mathbf{Q}_{xx}$  and  $V_{\beta}(\mathbf{W}) = V_{\text{cls}}$  as defined in Theorem 8.6.3.  $V_{\beta}(\mathbf{W})$  is minimized with  $\mathbf{W} = V_{\beta}^{-1}$ , in which case the asymptotic variance is*

$$V_{\beta}^* = V_{\beta} - V_{\beta} \mathbf{R} (\mathbf{R}' V_{\beta} \mathbf{R})^{-1} \mathbf{R}' V_{\beta}.$$

The asymptotic variance matrix for the efficient minimum distance estimator can be estimated by

$$\hat{V}_{\beta}^* = \hat{V}_{\beta} - \hat{V}_{\beta} \hat{\mathbf{R}} \left( \hat{\mathbf{R}}' \hat{V}_{\beta} \hat{\mathbf{R}} \right)^{-1} \hat{\mathbf{R}}' \hat{V}_{\beta}$$

where

$$\hat{\mathbf{R}} = \frac{\partial}{\partial \beta} \mathbf{r}(\tilde{\beta}_{\text{md}})'. \quad (8.52)$$

Standard errors for the elements of  $\tilde{\beta}_{\text{md}}$  are the square roots of the diagonal elements of  $\hat{V}_{\beta}^* = n^{-1} \hat{V}_{\beta}^*$ .

## 8.15 Inequality Restrictions

Inequality constraints on the parameter vector  $\beta$  take the form

$$\mathbf{r}(\beta) \geq \mathbf{0} \quad (8.53)$$

for some function  $\mathbf{r} : \mathbb{R}^k \rightarrow \mathbb{R}^q$ . The most common example is a non-negative constraint

$$\beta_1 \geq 0.$$

The constrained least-squares and minimum distance estimators can be written as

$$\tilde{\beta}_{\text{cls}} = \underset{\mathbf{r}(\beta) \geq \mathbf{0}}{\operatorname{argmin}} SSE(\beta) \quad (8.54)$$

and

$$\tilde{\beta}_{\text{md}} = \underset{\mathbf{r}(\beta) \geq \mathbf{0}}{\operatorname{argmin}} J(\beta). \quad (8.55)$$

Except in special cases the constrained estimators do not have simple algebraic solutions. An important exception is when there is a single non-negativity constraint, e.g.  $\beta_1 \geq 0$  with  $q = 1$ . In this case the constrained estimator can be found by two-step approach. First compute the unconstrained estimator  $\hat{\beta}$ . If  $\hat{\beta}_1 \geq 0$  then  $\tilde{\beta} = \hat{\beta}$ . Second, if  $\hat{\beta}_1 < 0$  then impose  $\beta_1 = 0$  (eliminate the regressor  $X_1$ ) and re-estimate. This yields the constrained least-squares estimator. While this method works when there is a single non-negativity constraint, it does not immediately generalize to other contexts.

The computational problems (8.54) and (8.55) are examples of **quadratic programming** problems. Quick and easy computer algorithms are available in programming languages including MATLAB, GAUSS and R.

Inference on inequality-constrained estimators is unfortunately quite challenging. The conventional asymptotic theory gives rise to the following dichotomy. If the true parameter satisfies the strict inequality  $\mathbf{r}(\boldsymbol{\beta}) > \mathbf{0}$ , then asymptotically the estimator is not subject to the constraint and the inequality-constrained estimator has an asymptotic distribution equal to the unconstrained case. However if the true parameter is on the boundary, e.g.  $\mathbf{r}(\boldsymbol{\beta}) = \mathbf{0}$ , then the estimator has a truncated structure. This is easiest to see in the one-dimensional case. If we have an estimator  $\hat{\beta}$  which satisfies  $\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} Z = N(0, V_\beta)$  and  $\beta = 0$ , then the constrained estimator  $\tilde{\beta} = \max[\hat{\beta}, 0]$  will have the asymptotic distribution  $\sqrt{n}\tilde{\beta} \xrightarrow{d} \max[Z, 0]$ , a “half-normal” distribution.

## 8.16 Technical Proofs\*

**Proof of Theorem 8.8.1, Equation (8.31).** Let  $\mathbf{R}_\perp$  be a full rank  $k \times (k - q)$  matrix satisfying  $\mathbf{R}'_\perp \mathbf{V}_\beta \mathbf{R} = \mathbf{0}$  and then set  $\mathbf{C} = [\mathbf{R}, \mathbf{R}_\perp]$  which is full rank and invertible. Then we can calculate that

$$\begin{aligned} \mathbf{C}' \mathbf{V}_\beta^* \mathbf{C} &= \begin{bmatrix} \mathbf{R}' \mathbf{V}_\beta^* \mathbf{R} & \mathbf{R}' \mathbf{V}_\beta^* \mathbf{R}_\perp \\ \mathbf{R}'_\perp \mathbf{V}_\beta^* \mathbf{R} & \mathbf{R}'_\perp \mathbf{V}_\beta^* \mathbf{R}_\perp \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{R}'_\perp \mathbf{V}_\beta \mathbf{R}_\perp \end{bmatrix} \end{aligned}$$

and

$$\begin{aligned} &\mathbf{C}' \mathbf{V}_\beta(\mathbf{W}) \mathbf{C} \\ &= \begin{bmatrix} \mathbf{R}' \mathbf{V}_\beta^*(\mathbf{W}) \mathbf{R} & \mathbf{R}' \mathbf{V}_\beta^*(\mathbf{W}) \mathbf{R}_\perp \\ \mathbf{R}'_\perp \mathbf{V}_\beta^*(\mathbf{W}) \mathbf{R} & \mathbf{R}'_\perp \mathbf{V}_\beta^*(\mathbf{W}) \mathbf{R}_\perp \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{R}'_\perp \mathbf{V}_\beta \mathbf{R}_\perp + \mathbf{R}'_\perp \mathbf{W} \mathbf{R} (\mathbf{R}' \mathbf{W} \mathbf{R})^{-1} \mathbf{R}' \mathbf{V}_\beta \mathbf{R} (\mathbf{R}' \mathbf{W} \mathbf{R})^{-1} \mathbf{R}' \mathbf{W} \mathbf{R}_\perp \end{bmatrix}. \end{aligned}$$

Thus

$$\begin{aligned} &\mathbf{C}' (\mathbf{V}_\beta(\mathbf{W}) - \mathbf{V}_\beta^*) \mathbf{C} \\ &= \mathbf{C}' \mathbf{V}_\beta(\mathbf{W}) \mathbf{C} - \mathbf{C}' \mathbf{V}_\beta^* \mathbf{C} \\ &= \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{R}'_\perp \mathbf{W} \mathbf{R} (\mathbf{R}' \mathbf{W} \mathbf{R})^{-1} \mathbf{R}' \mathbf{V}_\beta \mathbf{R} (\mathbf{R}' \mathbf{W} \mathbf{R})^{-1} \mathbf{R}' \mathbf{W} \mathbf{R}_\perp \end{bmatrix} \\ &\geq \mathbf{0} \end{aligned}$$

Since  $\mathbf{C}$  is invertible it follows that  $\mathbf{V}_\beta(\mathbf{W}) - \mathbf{V}_\beta^* \geq \mathbf{0}$  which is (8.31).  $\blacksquare$

cls

**Proof of Theorem 8.14.1.** We show the result for the minimum distance estimator  $\tilde{\boldsymbol{\beta}} = \tilde{\boldsymbol{\beta}}_{\text{md}}$ , as the proof for the constrained least-squares estimator is similar. For simplicity we assume that the constrained estimator is consistent  $\tilde{\boldsymbol{\beta}} \xrightarrow{p} \boldsymbol{\beta}$ . This can be shown with more effort, but requires a deeper treatment than appropriate for this textbook.

For each element  $r_j(\boldsymbol{\beta})$  of the  $q$ -vector  $\mathbf{r}(\boldsymbol{\beta})$ , by the mean value theorem there exists a  $\boldsymbol{\beta}_j^*$  on the line segment joining  $\tilde{\boldsymbol{\beta}}$  and  $\boldsymbol{\beta}$  such that

$$\mathbf{r}_j(\tilde{\boldsymbol{\beta}}) = \mathbf{r}_j(\boldsymbol{\beta}) + \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{r}_j(\boldsymbol{\beta}_j^*)' (\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}). \quad (8.56)$$

Let  $\mathbf{R}_n^*$  be the  $k \times q$  matrix

$$\mathbf{R}^* = \begin{bmatrix} \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{r}_1(\boldsymbol{\beta}_1^*) & \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{r}_2(\boldsymbol{\beta}_2^*) & \cdots & \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{r}_q(\boldsymbol{\beta}_q^*) \end{bmatrix}.$$

Since  $\tilde{\boldsymbol{\beta}} \xrightarrow{p} \boldsymbol{\beta}$  it follows that  $\boldsymbol{\beta}_j^* \xrightarrow{p} \boldsymbol{\beta}$ , and by the CMT,  $\mathbf{R}^* \xrightarrow{p} \mathbf{R}$ . Stacking the (8.56), we obtain

$$\mathbf{r}(\tilde{\boldsymbol{\beta}}) = \mathbf{r}(\boldsymbol{\beta}) + \mathbf{R}^{*'} (\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}).$$

Since  $\mathbf{r}(\tilde{\boldsymbol{\beta}}) = \mathbf{0}$  by construction and  $\mathbf{r}(\boldsymbol{\beta}) = \mathbf{0}$  by Assumption 8.6.1, this implies

$$\mathbf{0} = \mathbf{R}^{*'} (\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}). \quad (8.57)$$

The first-order condition for (8.51) is

$$\widehat{\mathbf{W}} (\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}) = \widehat{\mathbf{R}} \tilde{\boldsymbol{\lambda}}.$$

where  $\widehat{\mathbf{R}}$  is defined in (8.52).

Premultiplying by  $\mathbf{R}^{*'} \widehat{\mathbf{W}}^{-1}$ , inverting, and using (8.57), we find

$$\tilde{\boldsymbol{\lambda}} = \left( \mathbf{R}^{*'} \widehat{\mathbf{W}}^{-1} \widehat{\mathbf{R}} \right)^{-1} \mathbf{R}^{*'} (\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}) = \left( \mathbf{R}^{*'} \widehat{\mathbf{W}}^{-1} \widehat{\mathbf{R}} \right)^{-1} \mathbf{R}^{*'} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}).$$

Thus

$$\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta} = \left( \mathbf{I} - \widehat{\mathbf{W}}^{-1} \widehat{\mathbf{R}} \left( \mathbf{R}_n^{*'} \widehat{\mathbf{W}}^{-1} \widehat{\mathbf{R}} \right)^{-1} \mathbf{R}_n^{*'} \right) (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}). \quad (8.58)$$

From Theorem 7.3.2 and Theorem 7.7.1 we find

$$\begin{aligned} \sqrt{n} (\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}) &= \left( \mathbf{I} - \widehat{\mathbf{W}}^{-1} \widehat{\mathbf{R}} \left( \mathbf{R}_n^{*'} \widehat{\mathbf{W}}^{-1} \widehat{\mathbf{R}} \right)^{-1} \mathbf{R}_n^{*'} \right) \sqrt{n} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \\ &\xrightarrow{d} \left( \mathbf{I} - \mathbf{W}^{-1} \mathbf{R} (\mathbf{R}' \mathbf{W}^{-1} \mathbf{R})^{-1} \mathbf{R}' \right) \mathbf{N}(\mathbf{0}, \mathbf{V}_{\boldsymbol{\beta}}) \\ &= \mathbf{N}(\mathbf{0}, \mathbf{V}_{\boldsymbol{\beta}}(\mathbf{W})). \end{aligned}$$

■

## Exercises

**Exercise 8.1** In the model  $\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{e}$ , show directly from definition (8.3) that the CLS estimate of  $\boldsymbol{\beta} = (\boldsymbol{\beta}_1, \boldsymbol{\beta}_2)$  subject to the constraint that  $\boldsymbol{\beta}_2 = \mathbf{0}$  is the OLS regression of  $\mathbf{y}$  on  $\mathbf{X}_1$ .

**Exercise 8.2** In the model  $\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{e}$ , show directly from definition (8.3) that the CLS estimate of  $\boldsymbol{\beta} = (\boldsymbol{\beta}_1, \boldsymbol{\beta}_2)$ , subject to the constraint that  $\boldsymbol{\beta}_1 = \mathbf{c}$  (where  $\mathbf{c}$  is some given vector) is the OLS regression of  $\mathbf{y} - \mathbf{X}_1\mathbf{c}$  on  $\mathbf{X}_2$ .

**Exercise 8.3** In the model  $\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{e}$ , with  $\mathbf{X}_1$  and  $\mathbf{X}_2$  each  $n \times k$ , find the CLS estimate of  $\boldsymbol{\beta} = (\boldsymbol{\beta}_1, \boldsymbol{\beta}_2)$ , subject to the constraint that  $\boldsymbol{\beta}_1 = -\boldsymbol{\beta}_2$ .

**Exercise 8.4** In the linear projection model  $y_i = \alpha + \mathbf{x}'_i\boldsymbol{\beta} + e_i$ , consider the restriction  $\boldsymbol{\beta} = \mathbf{0}$ .

- (a) Find the constrained least-squares (CLS) estimator of  $\alpha$  under the restriction  $\boldsymbol{\beta} = \mathbf{0}$ .
- (b) Find an expression for the efficient minimum distance estimator of  $\alpha$  under the restriction  $\boldsymbol{\beta} = \mathbf{0}$ .

**Exercise 8.5** Verify that for  $\tilde{\boldsymbol{\beta}}_{\text{cls}}$  defined in (8.9) that  $\mathbf{R}'\tilde{\boldsymbol{\beta}}_{\text{cls}} = \mathbf{c}$ .

**Exercise 8.6** Prove Theorem 8.4.1

**Exercise 8.7** Prove Theorem 8.4.2, that is,  $\mathbb{E}(\tilde{\boldsymbol{\beta}}_{\text{cls}} | \mathbf{X}) = \boldsymbol{\beta}$ , under the assumptions of the linear regression model and (8.1).

Hint: Use Theorem 8.4.1.

**Exercise 8.8** Prove Theorem 8.4.3.

**Exercise 8.9** Prove Theorem 8.4.4, that is,  $\mathbb{E}(s_{\text{cls}}^2 | \mathbf{X}) = \sigma^2$ , under the assumptions of the homoskedastic regression model and (8.1).

**Exercise 8.10** Verify (8.24) and (8.25), and that the minimum distance estimator  $\tilde{\boldsymbol{\beta}}_{\text{md}}$  with  $\widehat{\mathbf{W}} = \widehat{\mathbf{Q}}_{xx}$  equals the CLS estimator.

**Exercise 8.11** Prove Theorem 8.6.1.

**Exercise 8.12** Prove Theorem 8.6.2.

**Exercise 8.13** Prove Theorem 8.6.3. (Hint: Use that CLS is a special case of Theorem 8.6.2.)

**Exercise 8.14** Verify that (8.29) is  $\mathbf{V}_{\boldsymbol{\beta}}(\mathbf{W})$  with  $\mathbf{W} = \mathbf{V}_{\boldsymbol{\beta}}^{-1}$ .

**Exercise 8.15** Prove (8.30). Hint: Use (8.29).

**Exercise 8.16** Verify (8.32), (8.33) and (8.34)

**Exercise 8.17** Verify (8.35), (8.36), and (8.37).



**Exercise 8.18** Suppose you have two independent samples

$$y_{1i} = \mathbf{x}'_{1i}\boldsymbol{\beta}_1 + e_{1i}$$

and

$$y_{2i} = \mathbf{x}'_{2i}\boldsymbol{\beta}_2 + e_{2i}$$

both of sample size  $n$ , and both  $\mathbf{x}_{1i}$  and  $\mathbf{x}_{2i}$  are  $k \times 1$ . You estimate  $\boldsymbol{\beta}_1$  and  $\boldsymbol{\beta}_2$  by OLS on each sample,  $\hat{\boldsymbol{\beta}}_1$  and  $\hat{\boldsymbol{\beta}}_2$ , say, with asymptotic covariance matrix estimators  $\hat{\mathbf{V}}_{\boldsymbol{\beta}_1}$  and  $\hat{\mathbf{V}}_{\boldsymbol{\beta}_2}$  (which are consistent for the asymptotic covariance matrices  $\mathbf{V}_{\boldsymbol{\beta}_1}$  and  $\mathbf{V}_{\boldsymbol{\beta}_2}$ ). Consider efficient minimum distance estimation under the restriction  $\boldsymbol{\beta}_1 = \boldsymbol{\beta}_2$ .

- Find the estimator  $\tilde{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta} = \boldsymbol{\beta}_1 = \boldsymbol{\beta}_2$
- Find the asymptotic distribution of  $\tilde{\boldsymbol{\beta}}$ .
- How would you approach the problem if the sample sizes are different, say  $n_1$  and  $n_2$ ?

**Exercise 8.19** As in Exercise 7.29 and 3.24, use the CPS dataset and the subsample of white male Hispanics.

- Estimate the regression

$$\begin{aligned} \log(\widehat{Wage}) = & \beta_1 \text{education} + \beta_2 \text{experience} + \beta_3 \text{experience}^2/100 + \beta_4 \text{Married}_1 \\ & + \beta_5 \text{Married}_2 + \beta_6 \text{Married}_3 + \beta_7 \text{Widowed} + \beta_8 \text{Divorced} + \beta_9 \text{Separated} + \beta_{10} \end{aligned}$$

where  $\text{Married}_1$ ,  $\text{Married}_2$ , and  $\text{Married}_3$  are the first three marital status codes as listed in Section 3.19.

- Estimate the equation using constrained least-squares, imposing the constraints  $\beta_4 = \beta_7$  and  $\beta_8 = \beta_9$ , and report the estimates and standard errors
- Estimate the equation using efficient minimum distance, imposing the same constraints, and report the estimates and standard errors
- Under what constraint on the coefficients is the wage equation non-decreasing in experience for experience up to 50?
- Estimate the equation imposing  $\beta_4 = \beta_7$ ,  $\beta_8 = \beta_9$ , and the inequality from part (d).

**Exercise 8.20** Take the model

$$\begin{aligned} y_i &= m(x_i) + e_i \\ m(x) &= \beta_0 + \beta_1 x + \beta_2 x^2 + \cdots + \beta_p x^p \\ \mathbb{E}(z_i e_i) &= 0 \\ z_i &= (1, x_i, \dots, x_i^p)' \\ g(x) &= \frac{d}{dx} m(x) \end{aligned}$$

with iid observations  $(y_i, x_i)$ ,  $i = 1, \dots, n$ . The order of the polynomial  $p$  is known.

- How should we interpret the function  $m(x)$  given the projection assumption? How should we interpret  $g(x)$ ? (Briefly)
- Describe an estimator  $\hat{g}(x)$  of  $g(x)$ .

- (c) Find the asymptotic distribution of  $\sqrt{n}(\hat{g}(x) - g(x))$  as  $n \rightarrow \infty$ .
- (d) Show how to construct an asymptotic 95% confidence interval for  $g(x)$  (for a single  $x$ ).
- (e) Assume  $p = 2$ . Describe how to estimate  $g(x)$  imposing the constraint that  $m(x)$  is concave.
- (f) Assume  $p = 2$ . Describe how to estimate  $g(x)$  imposing the constraint that  $m(u)$  is increasing on the region  $u \in [x_L, x_U]$ .

**Exercise 8.21** Take the linear model with restrictions

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= \mathbf{0} \\ \mathbf{R}' \boldsymbol{\beta} &= \mathbf{c} \end{aligned}$$

with  $n$  observations. Consider three estimators for  $\boldsymbol{\beta}$

- $\hat{\boldsymbol{\beta}}$ , the unconstrained least squares estimator
- $\tilde{\boldsymbol{\beta}}$ , the constrained least squares estimator
- $\bar{\boldsymbol{\beta}}$ , the constrained efficient minimum distance estimator

For each estimator, define the residuals  $\hat{e}_i = y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}$ ,  $\tilde{e}_i = y_i - \mathbf{x}_i' \tilde{\boldsymbol{\beta}}$ ,  $\bar{e}_i = y_i - \mathbf{x}_i' \bar{\boldsymbol{\beta}}$ , and variance estimators  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{e}_i^2$ ,  $\tilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \tilde{e}_i^2$ , and  $\bar{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \bar{e}_i^2$ .

- (a) As  $\bar{\boldsymbol{\beta}}$  is the most efficient estimator and  $\hat{\boldsymbol{\beta}}$  the least, do you expect that  $\bar{\sigma}^2 < \tilde{\sigma}^2 < \hat{\sigma}^2$ , in large samples?
- (b) Consider the statistic

$$T_n = \hat{\sigma}^{-2} \sum_{i=1}^n (\hat{e}_i - \tilde{e}_i)^2$$

Find the asymptotic distribution for  $T_n$  when  $\mathbf{R}' \boldsymbol{\beta} = \mathbf{c}$  is true.

- (c) Does the result of the previous question simplify when the error  $e_i$  is homoskedastic?

**Exercise 8.22** Take the linear model

$$\begin{aligned} y_i &= x_{1i} \beta_1 + x_{2i} \beta_2 + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= \mathbf{0} \end{aligned}$$

with  $n$  observations. Consider the restriction

$$\frac{\beta_1}{\beta_2} = 2$$

- (a) Find an explicit expression for the constrained least-squares (CLS) estimator  $\tilde{\boldsymbol{\beta}} = (\tilde{\beta}_1, \tilde{\beta}_2)$  of  $\boldsymbol{\beta} = (\beta_1, \beta_2)$  under the restriction. Your answer should be specific to the restriction, it should not be a generic formula for an abstract general restriction.
- (b) Derive the asymptotic distribution of  $\tilde{\beta}_1$  under the assumption that the restriction is true.

# Chapter 9

## Hypothesis Testing

In Chapter 5 we briefly introduced hypothesis testing in the context of the normal regression model. In this chapter we explore hypothesis testing in greater detail, with a particular emphasis on asymptotic inference.

### 9.1 Hypotheses

In Chapter 8 we discussed estimation subject to restrictions, including linear restrictions (8.1), nonlinear restrictions (8.47), and inequality restrictions (8.53). In this chapter we discuss **tests** of such restrictions.

Hypothesis tests attempt to assess whether there is evidence to contradict a proposed parametric restriction. Let

$$\boldsymbol{\theta} = \boldsymbol{r}(\boldsymbol{\beta})$$

be a  $q \times 1$  parameter of interest where  $\boldsymbol{r} : \mathbb{R}^k \rightarrow \Theta \subset \mathbb{R}^q$  is some transformation. For example,  $\boldsymbol{\theta}$  may be a single coefficient, e.g.  $\boldsymbol{\theta} = \beta_j$ , the difference between two coefficients, e.g.  $\boldsymbol{\theta} = \beta_j - \beta_\ell$ , or the ratio of two coefficients, e.g.  $\boldsymbol{\theta} = \beta_j / \beta_\ell$ .

A point hypothesis concerning  $\boldsymbol{\theta}$  is a proposed restriction such as

$$\boldsymbol{\theta} = \boldsymbol{\theta}_0 \tag{9.1}$$

where  $\boldsymbol{\theta}_0$  is a hypothesized (known) value.

More generally, letting  $\boldsymbol{\beta} \in \boldsymbol{B} \subset \mathbb{R}^k$  be the parameter space, a hypothesis is a restriction  $\boldsymbol{\beta} \in \boldsymbol{B}_0$  where  $\boldsymbol{B}_0$  is a proper subset of  $\boldsymbol{B}$ . This specializes to (9.1) by setting  $\boldsymbol{B}_0 = \{\boldsymbol{\beta} \in \boldsymbol{B} : \boldsymbol{r}(\boldsymbol{\beta}) = \boldsymbol{\theta}_0\}$ .

In this chapter we will focus exclusively on point hypotheses of the form (9.1) as they are the most common and relatively simple to handle.

The hypothesis to be tested is called the null hypothesis.

**Definition 9.1.1** The *null hypothesis*, written  $\mathbb{H}_0$ , is the restriction  $\boldsymbol{\theta} = \boldsymbol{\theta}_0$  or  $\boldsymbol{\beta} \in \boldsymbol{B}_0$ .

We often write the null hypothesis as  $\mathbb{H}_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$  or  $\mathbb{H}_0 : \boldsymbol{r}(\boldsymbol{\beta}) = \boldsymbol{\theta}_0$ .

The complement of the null hypothesis (the collection of parameter values which do not satisfy the null hypothesis) is called the alternative hypothesis.

**Definition 9.1.2** The *alternative hypothesis*, written  $\mathbb{H}_1$ , is the set  $\{\boldsymbol{\theta} \in \Theta : \boldsymbol{\theta} \neq \boldsymbol{\theta}_0\}$  or  $\{\boldsymbol{\beta} \in \boldsymbol{B} : \boldsymbol{\beta} \notin \boldsymbol{B}_0\}$ .

We often write the alternative hypothesis as  $\mathbb{H}_1 : \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$  or  $\mathbb{H}_1 : \boldsymbol{r}(\boldsymbol{\beta}) \neq \boldsymbol{\theta}_0$ . For simplicity, we often refer to the hypotheses as “the null” and “the alternative”.

In hypothesis testing, we assume that there is a true (but unknown) value of  $\boldsymbol{\theta}$  and this value either satisfies  $\mathbb{H}_0$  or does not satisfy  $\mathbb{H}_0$ . The goal of hypothesis testing is to assess whether or not  $\mathbb{H}_0$  is true, by asking if  $\mathbb{H}_0$  is consistent with the observed data.

To be specific, take our example of wage determination and consider the question: Does union membership affect wages? We can turn this into a hypothesis test by specifying the null as the restriction that a coefficient on union membership is zero in a wage regression. Consider, for example, the estimates reported in Table 4.1. The coefficient for “Male Union Member” is 0.095 (a wage premium of 9.5%) and the coefficient for “Female Union Member” is 0.022 (a wage premium of 2.2%). These are estimates, not the true values. The question is: Are the true coefficients zero? To answer this question, the testing method asks the question: Are the observed estimates compatible with the hypothesis, in the sense that the deviation from the hypothesis can be reasonably explained by stochastic variation? Or are the observed estimates incompatible with the hypothesis, in the sense that that the observed estimates would be highly unlikely if the hypothesis were true?

## 9.2 Acceptance and Rejection

A hypothesis test either accepts the null hypothesis or rejects the null hypothesis in favor of the alternative hypothesis. We can describe these two decisions as “Accept  $\mathbb{H}_0$ ” and “Reject  $\mathbb{H}_0$ ”. In the example given in the previous section, the decision would be either to accept the hypothesis that union membership does not affect wages, or to reject the hypothesis in favor of the alternative that union membership does affect wages.

The decision is based on the data, and so is a mapping from the sample space to the decision set. This splits the sample space into two regions  $S_0$  and  $S_1$  such that if the observed sample falls into  $S_0$  we accept  $\mathbb{H}_0$ , while if the sample falls into  $S_1$  we reject  $\mathbb{H}_0$ . The set  $S_0$  is called the **acceptance region** and the set  $S_1$  the **rejection** or **critical region**.

It is convenient to express this mapping as a real-valued function called a **test statistic**

$$T = T((y_1, \boldsymbol{x}_1), \dots, (y_n, \boldsymbol{x}_n))$$

relative to a **critical value**  $c$ . The hypothesis test then consists of the decision rule

1. Accept  $\mathbb{H}_0$  if  $T \leq c$ .
2. Reject  $\mathbb{H}_0$  if  $T > c$ .

A test statistic  $T$  should be designed so that small values are likely when  $\mathbb{H}_0$  is true and large values are likely when  $\mathbb{H}_1$  is true. There is a well developed statistical theory concerning the design of optimal tests. We will not review that theory here, but instead refer the reader to Lehmann and Romano (2005). In this chapter we will summarize the main approaches to the design of test statistics.

The most commonly used test statistic is the absolute value of the t-statistic

$$T = |T(\theta_0)| \tag{9.2}$$

where

$$T(\theta) = \frac{\hat{\theta} - \theta}{s(\hat{\theta})} \tag{9.3}$$

is the t-statistic from (7.43),  $\hat{\theta}$  is a point estimate, and  $s(\hat{\theta})$  its standard error.  $T$  is an appropriate statistic when testing hypotheses on individual coefficients or real-valued parameters  $\theta = h(\boldsymbol{\beta})$  and  $\theta_0$  is the hypothesized value. Quite typically,  $\theta_0 = 0$ , as interest focuses on whether or not a coefficient equals zero, but this is not the only possibility. For example, interest may focus on whether an elasticity  $\theta$  equals 1, in which case we may wish to test  $\mathbb{H}_0 : \theta = 1$ .

### 9.3 Type I Error

A false rejection of the null hypothesis  $\mathbb{H}_0$  (rejecting  $\mathbb{H}_0$  when  $\mathbb{H}_0$  is true) is called a **Type I error**. The probability of a Type I error is

$$\Pr(\text{Reject } \mathbb{H}_0 \mid \mathbb{H}_0 \text{ true}) = \Pr(T > c \mid \mathbb{H}_0 \text{ true}). \quad (9.4)$$

The finite sample **size** of the test is defined as the supremum of (9.4) across all data distributions which satisfy  $\mathbb{H}_0$ . A primary goal of test construction is to limit the incidence of Type I error by bounding the size of the test.

For the reasons discussed in Chapter 7, in typical econometric models the exact sampling distributions of estimators and test statistics are unknown and hence we cannot explicitly calculate (9.4). Instead, we typically rely on asymptotic approximations. Suppose that the test statistic has an asymptotic distribution under  $\mathbb{H}_0$ . That is, when  $\mathbb{H}_0$  is true

$$T \xrightarrow{d} \xi \quad (9.5)$$

as  $n \rightarrow \infty$  for some continuously-distributed random variable  $\xi$ . This is not a substantive restriction, as most conventional econometric tests satisfy (9.5). Let  $G(u) = \Pr(\xi \leq u)$  denote the distribution of  $\xi$ . We call  $\xi$  (or  $G$ ) the **asymptotic null distribution**.

It is generally desirable to design test statistics  $T$  whose asymptotic null distribution  $G$  is known and does not depend on unknown parameters. In this case we say that the statistic  $T$  is **asymptotically pivotal**.

For example, if the test statistic equals the absolute t-statistic from (9.2), then we know from Theorem 7.12.1 that if  $\theta = \theta_0$  (that is, the null hypothesis holds), then  $T \xrightarrow{d} |Z|$  as  $n \rightarrow \infty$  where  $Z \sim N(0, 1)$ . This means that  $G(u) = \Pr(|Z| \leq u) = 2\Phi(u) - 1$ , the distribution of the absolute value of the standard normal as shown in (7.44). This distribution does not depend on unknowns and is pivotal.

We define the **asymptotic size** of the test as the asymptotic probability of a Type I error:

$$\begin{aligned} \lim_{n \rightarrow \infty} \Pr(T > c \mid \mathbb{H}_0 \text{ true}) &= \Pr(\xi > c) \\ &= 1 - G(c). \end{aligned}$$

We see that the asymptotic size of the test is a simple function of the asymptotic null distribution  $G$  and the critical value  $c$ . For example, the asymptotic size of a test based on the absolute t-statistic with critical value  $c$  is  $2(1 - \Phi(c))$ .

In the dominant approach to hypothesis testing, the researcher pre-selects a **significance level**  $\alpha \in (0, 1)$  and then selects  $c$  so that the (asymptotic) size is no larger than  $\alpha$ . When the asymptotic null distribution  $G$  is pivotal, we can accomplish this by setting  $c$  equal to the  $1 - \alpha$  quantile of the distribution  $G$ . (If the distribution  $G$  is not pivotal, more complicated methods must be used, pointing out the great convenience of using asymptotically pivotal test statistics.) We call  $c$  the **asymptotic critical value** because it has been selected from the asymptotic null distribution. For example, since  $2(1 - \Phi(1.96)) = 0.05$ , it follows that the 5% asymptotic critical value for the absolute t-statistic is  $c = 1.96$ . Calculation of normal critical values is done numerically in statistical software. For example, in MATLAB the command is `norminv(1- $\alpha$ /2)`.

### 9.4 t tests

As we mentioned earlier, the most common test of the one-dimensional hypothesis

$$\mathbb{H}_0 : \theta = \theta_0 \quad (9.6)$$

against the alternative

$$\mathbb{H}_1 : \theta \neq \theta_0 \quad (9.7)$$

is the absolute value of the t-statistic (9.3). We now formally state its asymptotic null distribution, which is a simple application of Theorem 7.12.1.

**Theorem 9.4.1** *Under Assumptions 7.1.2, 7.10.1, and  $\mathbb{H}_0 : \theta = \theta_0$ ,*

$$T(\theta_0) \xrightarrow{d} Z.$$

*For  $c$  satisfying  $\alpha = 2(1 - \Phi(c))$ ,*

$$\Pr(|T(\theta_0)| > c \mid \mathbb{H}_0) \longrightarrow \alpha,$$

*and the test “Reject  $\mathbb{H}_0$  if  $|T(\theta_0)| > c$ ” has asymptotic size  $\alpha$ .*

The theorem shows that asymptotic critical values can be taken from the normal distribution. As in our discussion of asymptotic confidence intervals (Section 7.13), the critical value could alternatively be taken from the student  $t$  distribution, which would be the exact test in the normal regression model (Section 5.14). Indeed,  $t$  critical values are the default in packages such as Stata. Since the critical values from the student  $t$  distribution are (slightly) larger than those from the normal distribution, using student  $t$  critical values decreases the rejection probability of the test. In practical applications the difference is typically unimportant unless the sample size is quite small (in which case the asymptotic approximation should be questioned as well).

The alternative hypothesis  $\theta \neq \theta_0$  is sometimes called a “two-sided” alternative. In contrast, sometimes we are interested in testing for one-sided alternatives such as

$$\mathbb{H}_1 : \theta > \theta_0 \quad (9.8)$$

or

$$\mathbb{H}_1 : \theta < \theta_0. \quad (9.9)$$

Tests of  $\theta = \theta_0$  against  $\theta > \theta_0$  or  $\theta < \theta_0$  are based on the signed t-statistic  $T = T(\theta_0)$ . The hypothesis  $\theta = \theta_0$  is rejected in favor of  $\theta > \theta_0$  if  $T > c$  where  $c$  satisfies  $\alpha = 1 - \Phi(c)$ . Negative values of  $T$  are not taken as evidence against  $\mathbb{H}_0$ , as point estimates  $\hat{\theta}$  less than  $\theta_0$  do not point to  $\theta > \theta_0$ . Since the critical values are taken from the single tail of the normal distribution, they are smaller than for two-sided tests. Specifically, the asymptotic 5% critical value is  $c = 1.645$ . Thus, we reject  $\theta = \theta_0$  in favor of  $\theta > \theta_0$  if  $T > 1.645$ .

Conversely, tests of  $\theta = \theta_0$  against  $\theta < \theta_0$  reject  $\mathbb{H}_0$  for negative t-statistics, e.g. if  $T \leq -c$ . For this alternative large positive values of  $T$  are not evidence against  $\mathbb{H}_0$ . An asymptotic 5% test rejects if  $T < -1.645$ .

There seems to be an ambiguity. Should we use the two-sided critical value 1.96 or the one-sided critical value 1.645? The answer is that we should use one-sided tests and critical values only when the parameter space is known to satisfy a one-sided restriction such as  $\theta \geq \theta_0$ . This is when the test of  $\theta = \theta_0$  against  $\theta > \theta_0$  makes sense. If the restriction  $\theta \geq \theta_0$  is not known *a priori*, then imposing this restriction to test  $\theta = \theta_0$  against  $\theta > \theta_0$  does not make sense. Since linear regression coefficients typically do not have *a priori* sign restrictions, the standard convention is to use two-sided critical values.

This may seem contrary to the way testing is presented in statistical textbooks, which often focus on one-sided alternative hypotheses. The latter focus is primarily for pedagogy, as the one-sided theoretical problem is cleaner and easier to understand.

## 9.5 Type II Error and Power

A false acceptance of the null hypothesis  $\mathbb{H}_0$  (accepting  $\mathbb{H}_0$  when  $\mathbb{H}_1$  is true) is called a **Type II error**. The rejection probability under the alternative hypothesis is called the **power** of the test, and equals 1 minus the probability of a Type II error:

$$\pi(\boldsymbol{\theta}) = \Pr(\text{Reject } \mathbb{H}_0 \mid \mathbb{H}_1 \text{ true}) = \Pr(T > c \mid \mathbb{H}_1 \text{ true}).$$

We call  $\pi(\boldsymbol{\theta})$  the **power function** and is written as a function of  $\boldsymbol{\theta}$  to indicate its dependence on the true value of the parameter  $\boldsymbol{\theta}$ .

In the dominant approach to hypothesis testing, the goal of test construction is to have high power subject to the constraint that the size of the test is lower than the pre-specified significance level. Generally, the power of a test depends on the true value of the parameter  $\boldsymbol{\theta}$ , and for a well behaved test the power is increasing both as  $\boldsymbol{\theta}$  moves away from the null hypothesis  $\boldsymbol{\theta}_0$  and as the sample size  $n$  increases.

Given the two possible states of the world ( $\mathbb{H}_0$  or  $\mathbb{H}_1$ ) and the two possible decisions (Accept  $\mathbb{H}_0$  or Reject  $\mathbb{H}_0$ ), there are four possible pairings of states and decisions as is depicted in the following chart.

Hypothesis Testing Decisions		
	Accept $\mathbb{H}_0$	Reject $\mathbb{H}_0$
$\mathbb{H}_0$ true	Correct Decision	Type I Error
$\mathbb{H}_1$ true	Type II Error	Correct Decision

Given a test statistic  $T$ , increasing the critical value  $c$  increases the acceptance region  $S_0$  while decreasing the rejection region  $S_1$ . This decreases the likelihood of a Type I error (decreases the size) but increases the likelihood of a Type II error (decreases the power). Thus the choice of  $c$  involves a trade-off between size and the power. This is why the significance level  $\alpha$  of the test cannot be set arbitrarily small. (Otherwise the test will not have meaningful power.)

It is important to consider the power of a test when interpreting hypothesis tests, as an overly narrow focus on size can lead to poor decisions. For example, it is easy to design a test which has perfect size yet has trivial power. Specifically, for any hypothesis we can use the following test: Generate a random variable  $U \sim U[0, 1]$  and reject  $\mathbb{H}_0$  if  $U < \alpha$ . This test has exact size of  $\alpha$ . Yet the test also has power precisely equal to  $\alpha$ . When the power of a test equals the size, we say that the test has **trivial power**. Nothing is learned from such a test.

## 9.6 Statistical Significance

Testing requires a pre-selected choice of significance level  $\alpha$ , yet there is no objective scientific basis for choice of  $\alpha$ . Nevertheless the common practice is to set  $\alpha = 0.05$  (5%). Alternative values are  $\alpha = 0.10$  (10%) and  $\alpha = 0.01$  (1%). These choices are somewhat the by-product of traditional tables of critical values and statistical software.

The informal reasoning behind the choice of a 5% critical value is to ensure that Type I errors should be relatively unlikely – that the decision “Reject  $\mathbb{H}_0$ ” has scientific strength – yet the test retains power against reasonable alternatives. The decision “Reject  $\mathbb{H}_0$ ” means that the evidence is inconsistent with the null hypothesis, in the sense that it is relatively unlikely (1 in 20) that data generated by the null hypothesis would yield the observed test result.

In contrast, the decision “Accept  $\mathbb{H}_0$ ” is not a strong statement. It does not mean that the evidence supports  $\mathbb{H}_0$ , only that there is insufficient evidence to reject  $\mathbb{H}_0$ . Because of this, it is more accurate to use the label “Do not Reject  $\mathbb{H}_0$ ” instead of “Accept  $\mathbb{H}_0$ ”.

When a test rejects  $\mathbb{H}_0$  at the 5% significance level it is common to say that the statistic is **statistically significant** and if the test accepts  $\mathbb{H}_0$  it is common to say that the statistic is **not statistically significant** or that it is **statistically insignificant**. It is helpful to remember that this is simply a compact way of saying “Using the statistic  $T$ , the hypothesis  $\mathbb{H}_0$  can [cannot] be rejected at the asymptotic 5% level.” Furthermore, when the null hypothesis  $\mathbb{H}_0 : \theta = 0$  is rejected it is common to say that the coefficient  $\theta$  is statistically significant, because the test has rejected the hypothesis that the coefficient is equal to zero.

Let us return to the example about the union wage premium as measured in Table 4.1. The absolute t-statistic for the coefficient on “Male Union Member” is  $0.095/0.020 = 4.7$ , which is greater than the 5% asymptotic critical value of 1.96. Therefore we reject the hypothesis that union membership does not affect wages for men. In this case, we can say that union membership is statistically significant for men. However, the absolute t-statistic for the coefficient on “Female Union Member” is  $0.023/0.020 = 1.2$ , which is less than 1.96 and therefore we do not reject the hypothesis that union membership does not affect wages for women. In this case we find that membership for women is not statistically significant.

When a test accepts a null hypothesis (when a test is not statistically significant) a common misinterpretation is that this is evidence that the null hypothesis is true. This is incorrect. Failure to reject is by itself not evidence. Without an analysis of power, we do not know the likelihood of making a Type II error, and thus are uncertain. In our wage example, it would be a mistake to write that “the regression finds that female union membership has no effect on wages”. This is an incorrect and most unfortunate interpretation. The test has failed to reject the hypothesis that the coefficient is zero, but that does not mean that the coefficient is actually zero.

When a test rejects a null hypothesis (when a test is statistically significant) it is strong evidence against the hypothesis (since if the hypothesis were true then rejection is an unlikely event). Rejection should be taken as evidence against the null hypothesis. However, we can never conclude that the null hypothesis is indeed false, as we cannot exclude the possibility that we are making a Type I error.

Perhaps more importantly, there is an important distinction between statistical and economic significance. If we correctly reject the hypothesis  $\mathbb{H}_0 : \theta = 0$  it means that the true value of  $\theta$  is non-zero. This includes the possibility that  $\theta$  may be non-zero but close to zero in magnitude. This only makes sense if we interpret the parameters in the context of their relevant models. In our wage regression example, we might consider wage effects of 1% magnitude or less as being “close to zero”. In a log wage regression this corresponds to a dummy variable with a coefficient less than 0.01. If the standard error is sufficiently small (less than 0.005) then a coefficient estimate of 0.01 will be statistically significant, but not economically significant. This occurs frequently in applications with very large sample sizes where standard errors can be quite small.

The solution is to focus whenever possible on confidence intervals and the economic meaning of the coefficients. For example, if the coefficient estimate is 0.005 with a standard error of 0.002 then a 95% confidence interval would be  $[0.001, 0.009]$  indicating that the true effect is likely between 0% and 1%, and hence is slightly positive but small. This is much more informative than the misleading statement “the effect is statistically positive”.

## 9.7 P-Values

Continuing with the wage regression estimates reported in Table 4.1, consider another question: Does marriage status affect wages? To test the hypothesis that marriage status has no effect on wages, we examine the t-statistics for the coefficients on “Married Male” and “Married Female” in Table 4.1, which are  $0.211/0.010 = 22$  and  $0.016/0.010 = 1.7$ , respectively. The first exceeds the asymptotic 5% critical value of 1.96, so we reject the hypothesis for men, though not for women. But the statistic for men is exceptionally high, and that for women is only slightly below the critical value. Suppose in contrast that the t-statistic had been 2.0, which is more than the critical



value. This would lead to the decision “Reject  $\mathbb{H}_0$ ” rather than “Accept  $\mathbb{H}_0$ ”. Should we really be making a different decision if the  $t$ -statistic is 1.7 rather than 2.0? The difference in values is small, shouldn’t the difference in the decision be also small? Thinking through these examples it seems unsatisfactory to simply report “Accept  $\mathbb{H}_0$ ” or “Reject  $\mathbb{H}_0$ ”. These two decisions do not summarize the evidence. Instead, the magnitude of the statistic  $T$  suggests a “degree of evidence” against  $\mathbb{H}_0$ . How can we take this into account?

The answer is to report what is known as the **asymptotic p-value**

$$p = 1 - G(T).$$

Since the distribution function  $G$  is monotonically increasing, the  $p$ -value is a monotonically decreasing function of  $T$  and is an equivalent test statistic. Instead of rejecting  $\mathbb{H}_0$  at the significance level  $\alpha$  if  $T > c$ , we can reject  $\mathbb{H}_0$  if  $p < \alpha$ . Thus it is sufficient to report  $p$ , and let the reader decide. In practice, the  $p$ -value is calculated numerically. For example, in MATLAB the command is `2*(1-normalcdf(abs(t)))`.

It is instructive to interpret  $p$  as the **marginal significance level**: the largest value of  $\alpha$  for which the test  $T$  “rejects” the null hypothesis. That is,  $p = 0.11$  means that  $T$  rejects  $\mathbb{H}_0$  for all significance levels greater than 0.11, but fails to reject  $\mathbb{H}_0$  for significance levels less than 0.11.

Furthermore, the asymptotic  $p$ -value has a very convenient asymptotic null distribution. Since  $T \xrightarrow{d} \xi$  under  $\mathbb{H}_0$ , then  $p = 1 - G(T) \xrightarrow{d} 1 - G(\xi)$ , which has the distribution

$$\begin{aligned} \Pr(1 - G(\xi) \leq u) &= \Pr(1 - u \leq G(\xi)) \\ &= 1 - \Pr(\xi \leq G^{-1}(1 - u)) \\ &= 1 - G(G^{-1}(1 - u)) \\ &= 1 - (1 - u) \\ &= u, \end{aligned}$$

which is the uniform distribution on  $[0, 1]$ . (This calculation assumes that  $G(u)$  is strictly increasing which is true for conventional asymptotic distributions such as the normal.) Thus  $p \xrightarrow{d} U[0, 1]$ . This means that the “unusualness” of  $p$  is easier to interpret than the “unusualness” of  $T$ .

An important caveat is that the  $p$ -value  $p$  should not be interpreted as the probability that either hypothesis is true. A common mis-interpretation is that  $p$  is the probability “that the null hypothesis is true.” This is incorrect. Rather,  $p$  is the marginal significance level – a measure of the strength of information against the null hypothesis.

For a  $t$ -statistic, the  $p$ -value can be calculated either using the normal distribution or the student  $t$  distribution, the latter presented in Section 5.14.  $p$ -values calculated using the student  $t$  will be slightly larger, though the difference is small when the sample size is large.

Returning to our empirical example, for the test that the coefficient on “Married Male” is zero, the  $p$ -value is 0.000. This means that it would be nearly impossible to observe a  $t$ -statistic as large as 22 when the true value of the coefficient is zero. When presented with such evidence we can say that we “strongly reject” the null hypothesis, that the test is “highly significant”, or that “the test rejects at any conventional critical value”. In contrast, the  $p$ -value for the coefficient on “Married Female” is 0.094. In this context it is typical to say that the test is “close to significant”, meaning that the  $p$ -value is larger than 0.05, but not too much larger.

A related (but somewhat inferior) empirical practice is to append asterisks (\*) to coefficient estimates or test statistics to indicate the level of significance. A common practice is to append a single asterisk (\*) for an estimate or test statistic which exceeds the 10% critical value (i.e., is significant at the 10% level), append a double asterisk (\*\*) for a test which exceeds the 5% critical value, or append a triple asterisk (\*\*\*) for a test which exceeds the 1% critical value. Such a practice can be better than a table of raw test statistics as the asterisks permit a quick interpretation of significance. On the other hand, asterisks are inferior to  $p$ -values, which are also easy and quick to

interpret. The goal is essentially the same; it seems wiser to report p-values whenever possible and avoid the use of asterisks.

Our recommendation is that the best empirical practice is to compute and report the asymptotic p-value  $p$  rather than simply the test statistic  $T$ , the binary decision Accept/Reject, or appending asterisks. The p-value is a simple statistic, easy to interpret, and contains more information than the other choices.

We now summarize the main features of hypothesis testing.

1. Select a significance level  $\alpha$ .
2. Select a test statistic  $T$  with asymptotic distribution  $T \xrightarrow{d} \xi$  under  $\mathbb{H}_0$ .
3. Set the asymptotic critical value  $c$  so that  $1 - G(c) = \alpha$ , where  $G$  is the distribution function of  $\xi$ .
4. Calculate the asymptotic p-value  $p = 1 - G(T)$ .
5. Reject  $\mathbb{H}_0$  if  $T > c$ , or equivalently  $p < \alpha$ .
6. Accept  $\mathbb{H}_0$  if  $T \leq c$ , or equivalently  $p \geq \alpha$ .
7. Report  $p$  to summarize the evidence concerning  $\mathbb{H}_0$  versus  $\mathbb{H}_1$ .

## 9.8 t-ratios and the Abuse of Testing

In Section 4.18, we argued that a good applied practice is to report coefficient estimates  $\hat{\theta}$  and standard errors  $s(\hat{\theta})$  for all coefficients of interest in estimated models. With  $\hat{\theta}$  and  $s(\hat{\theta})$  the reader can easily construct confidence intervals  $[\hat{\theta} \pm 2s(\hat{\theta})]$  and t-statistics  $(\hat{\theta} - \theta_0)/s(\hat{\theta})$  for hypotheses of interest.

Some applied papers (especially older ones) report t-ratios  $T = \hat{\theta}/s(\hat{\theta})$  instead of standard errors. This is poor econometric practice. While the same information is being reported (you can back out standard errors by division, e.g.  $s(\hat{\theta}) = \hat{\theta}/T$ ), standard errors are generally more helpful to readers than t-ratios. Standard errors help the reader focus on the estimation precision and confidence intervals, while t-ratios focus attention on statistical significance. While statistical significance is important, it is less important that the parameter estimates themselves and their confidence intervals. The focus should be on the meaning of the parameter estimates, their magnitudes, and their interpretation, not on listing which variables have significant (e.g. non-zero) coefficients. In many modern applications, sample sizes are very large so standard errors can be very small. Consequently t-ratios can be large even if the coefficient estimates are economically small. In such contexts it may not be interesting to announce “The coefficient is non-zero!” Instead, what is interesting to announce is that “The coefficient estimate is economically interesting!”

In particular, some applied papers report coefficient estimates and t-ratios, and limit their discussion of the results to describing which variables are “significant” (meaning that their t-ratios exceed 2) and the signs of the coefficient estimates. This is very poor empirical work, and should be studiously avoided. It is also a recipe for banishment of your work to lower tier economics journals.

Fundamentally, the common t-ratio is a test for the hypothesis that a coefficient equals zero. This should be reported and discussed when this is an interesting economic hypothesis of interest. But if this is not the case, it is distracting.

One problem is that standard packages, such as Stata, by default report t-statistics and p-values for every estimated coefficient. While this can be useful (as a user doesn’t need to explicitly ask to test an desired coefficient) it can be misleading as it may unintentionally suggest that the entire list of t-statistics and p-values are important. Instead, a user should focus on tests of scientifically motivated hypotheses.

In general, when a coefficient  $\theta$  is of interest, it is constructive to focus on the point estimate, its standard error, and its confidence interval. The point estimate gives our “best guess” for the value. The standard error is a measure of precision. The confidence interval gives us the range of values consistent with the data. If the standard error is large then the point estimate is not a good summary about  $\theta$ . The endpoints of the confidence interval describe the bounds on the likely possibilities. If the confidence interval embraces too broad a set of values for  $\theta$ , then the dataset is not sufficiently informative to render useful inferences about  $\theta$ . On the other hand if the confidence interval is tight, then the data have produced an accurate estimate, and the focus should be on the value and interpretation of this estimate. In contrast, the statement “the t-ratio is highly significant” has little interpretive value.

The above discussion requires that the researcher knows what the coefficient  $\theta$  means (in terms of the economic problem) and can interpret values and magnitudes, not just signs. This is critical for good applied econometric practice.

For example, consider the question about the effect of marriage status on mean log wages. We had found that the effect is “highly significant” for men and “close to significant” for women. Now, let’s construct asymptotic 95% confidence intervals for the coefficients. The one for men is [0.19, 0.23] and that for women is [−0.00, 0.03]. This shows that average wages for married men are about 19-23% higher than for unmarried men, which is substantial, while the difference for women is about 0-3%, which is small. These *magnitudes* are more informative than the results of the hypothesis tests.

## 9.9 Wald Tests

The t-test is appropriate when the null hypothesis is a real-valued restriction. More generally, there may be multiple restrictions on the coefficient vector  $\beta$ . Suppose that we have  $q > 1$  restrictions which can be written in the form (9.1). It is natural to estimate  $\theta = r(\beta)$  by the plug-in estimate  $\hat{\theta} = r(\hat{\beta})$ . To test  $\mathbb{H}_0 : \theta = \theta_0$  against  $\mathbb{H}_1 : \theta \neq \theta_0$  one approach is to measure the magnitude of the discrepancy  $\hat{\theta} - \theta_0$ . As this is a vector, there is more than one measure of its length. One simple measure is the weighted quadratic form known as the **Wald statistic**. This is (7.47) evaluated at the null hypothesis

$$W = W(\theta_0) = (\hat{\theta} - \theta_0)' \hat{\mathbf{V}}_{\hat{\theta}}^{-1} (\hat{\theta} - \theta_0) \quad (9.10)$$

where  $\hat{\mathbf{V}}_{\hat{\theta}} = \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\hat{\beta}} \hat{\mathbf{R}}$  is an estimate of  $\mathbf{V}_{\hat{\theta}}$  and  $\hat{\mathbf{R}} = \frac{\partial}{\partial \beta} r(\hat{\beta})'$ . Notice that we can write  $W$  alternatively as

$$W = n (\hat{\theta} - \theta_0)' \hat{\mathbf{V}}_{\theta}^{-1} (\hat{\theta} - \theta_0)$$

using the asymptotic variance estimate  $\hat{\mathbf{V}}_{\theta}$ , or we can write it directly as a function of  $\hat{\beta}$  as

$$W = (r(\hat{\beta}) - \theta_0)' (\hat{\mathbf{R}}' \hat{\mathbf{V}}_{\hat{\beta}} \hat{\mathbf{R}})^{-1} (r(\hat{\beta}) - \theta_0). \quad (9.11)$$

Also, when  $r(\beta) = \mathbf{R}'\beta$  is a linear function of  $\beta$ , then the Wald statistic simplifies to

$$W = (\mathbf{R}'\hat{\beta} - \theta_0)' (\mathbf{R}' \hat{\mathbf{V}}_{\hat{\beta}} \mathbf{R})^{-1} (\mathbf{R}'\hat{\beta} - \theta_0).$$

The Wald statistic  $W$  is a weighted Euclidean measure of the length of the vector  $\hat{\theta} - \theta_0$ . When  $q = 1$  then  $W = T^2$ , the square of the t-statistic, so hypothesis tests based on  $W$  and  $|T|$  are equivalent. The Wald statistic (9.10) is a generalization of the t-statistic to the case of multiple restrictions. As the Wald statistic is symmetric in the argument  $\hat{\theta} - \theta_0$  it treats positive and negative alternatives symmetrically. Thus the inherent alternative is always two-sided.

As shown in Theorem 7.16.1, when  $\beta$  satisfies  $r(\beta) = \theta_0$  then  $W \xrightarrow{d} \chi_q^2$ , a chi-square random variable with  $q$  degrees of freedom. Let  $G_q(u)$  denote the  $\chi_q^2$  distribution function. For a given significance level  $\alpha$ , the asymptotic critical value  $c$  satisfies  $\alpha = 1 - G_q(c)$ . For example, the 5% critical values for  $q = 1$ ,  $q = 2$ , and  $q = 3$  are 3.84, 5.99, and 7.82, respectively, and in general the level  $\alpha$  critical value can be calculated in MATLAB as `chi2inv(1- $\alpha$ ,q)`. An asymptotic test rejects  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  if  $W > c$ . As with t-tests, it is conventional to describe a Wald test as “significant” if  $W$  exceeds the 5% asymptotic critical value.

**Theorem 9.9.1** *Under Assumptions 7.1.2 and 7.10.1, and  $\mathbb{H}_0 : \theta = \theta_0$ , then*

$$W \xrightarrow{d} \chi_q^2,$$

*and for  $c$  satisfying  $\alpha = 1 - G_q(c)$ ,*

$$\Pr(W > c \mid \mathbb{H}_0) \longrightarrow \alpha$$

*so the test “Reject  $\mathbb{H}_0$  if  $W > c$ ” has asymptotic size  $\alpha$ .*

Notice that the asymptotic distribution in Theorem 9.9.1 depends solely on  $q$ , the number of restrictions being tested. It does not depend on  $k$ , the number of parameters estimated.

The asymptotic p-value for  $W$  is  $p = 1 - G_q(W)$ , and this is particularly useful when testing multiple restrictions. For example, if you write that a Wald test on eight restrictions ( $q = 8$ ) has the value  $W = 11.2$ , it is difficult for a reader to assess the magnitude of this statistic unless they have quick access to a statistical table or software. Instead, if you write that the p-value is  $p = 0.19$  (as is the case for  $W = 11.2$  and  $q = 8$ ) then it is simple for a reader to interpret its magnitude as “insignificant”. To calculate the asymptotic p-value for a Wald statistic in MATLAB, use the command `1-chi2cdf(w,q)`.

Some packages (including Stata) and papers report  $F$  versions of Wald statistics. That is, for any Wald statistic  $W$  which tests a  $q$ -dimensional restriction, the  $F$  version of the test is

$$F = W/q.$$

When  $F$  is reported, it is conventional to use  $F_{q,n-k}$  critical values and p-values rather than  $\chi_q^2$  values. The connection between Wald and F statistics is demonstrated in Section 9.14 we show that when Wald statistics are calculated using a homoskedastic covariance matrix, then  $F = W/q$  is identical to the F statistic of (5.23). While there is no formal justification to using the  $F_{q,n-k}$  distribution for non-homoskedastic covariance matrices, the  $F_{q,n-k}$  distribution provides continuity with the exact distribution theory under normality and is a bit more conservative than the  $\chi_q^2$  distribution. (Furthermore, the difference is small when  $n - k$  is moderately large.)

To implement a test of zero restrictions in Stata, an easy method is to use the command “test X1 X2” where X1 and X2 are the names of the variables whose coefficients are hypothesized to equal zero. This command should be executed after executing a regression command. The  $F$  version of the Wald statistic is reported, using the covariance matrix calculated using the method specified in the regression command. A p-value is reported, calculated using the  $F_{q,n-k}$  distribution.

To illustrate, consider the empirical results presented in Table 4.1. The hypothesis “Union membership does not affect wages” is the joint restriction that both coefficients on “Male Union Member” and “Female Union Member” are zero. We calculate the Wald statistic for this joint hypothesis and find  $W = 23$  (or  $F = 12.5$ ) with a p-value of  $p = 0.000$ . Thus we reject the null hypothesis in favor of the alternative that at least one of the coefficients is non-zero. This does not mean that both coefficients are non-zero, just that one of the two is non-zero. Therefore examining both the joint Wald statistic and the individual t-statistics is useful for interpretation.

As a second example from the same regression, take the hypothesis that married status has no effect on mean wages for women. This is the joint restriction that the coefficients on “Married Female” and “Formerly Married Female” are zero. The Wald statistic for this hypothesis is  $W = 6.4$  ( $F = 3.2$ ) with a p-value of 0.04. Such a p-value is typically called “marginally significant”, in the sense that it is slightly smaller than 0.05.

### Abraham Wald

The Hungarian mathematician/statistician/econometrician Abraham Wald (1902-1950) developed an optimality property for the Wald test in terms of weighted average power. He also developed the field of sequential testing and the design of experiments.

## 9.10 Homoskedastic Wald Tests

If the error is known to be homoskedastic, then it is appropriate to use the homoskedastic Wald statistic (7.49) which replaces  $\widehat{\mathbf{V}}_{\widehat{\boldsymbol{\theta}}}$  with the homoskedastic estimate  $\widehat{\mathbf{V}}_{\widehat{\boldsymbol{\theta}}}^0$ . This statistic equals

$$\begin{aligned} W^0 &= (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)' (\widehat{\mathbf{V}}_{\widehat{\boldsymbol{\theta}}}^0)^{-1} (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \\ &= (\mathbf{r}(\widehat{\boldsymbol{\beta}}) - \boldsymbol{\theta}_0)' (\mathbf{R}' (\mathbf{X}' \mathbf{X})^{-1} \widehat{\mathbf{R}})^{-1} (\mathbf{r}(\widehat{\boldsymbol{\beta}}) - \boldsymbol{\theta}_0) / s^2. \end{aligned} \quad (9.12)$$

In the case of linear hypotheses  $\mathbb{H}_0 : \mathbf{R}' \boldsymbol{\beta} = \boldsymbol{\theta}_0$  we can write this as

$$W^0 = (\mathbf{R}' \widehat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0)' (\mathbf{R}' (\mathbf{X}' \mathbf{X})^{-1} \mathbf{R})^{-1} (\mathbf{R}' \widehat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0) / s^2. \quad (9.13)$$

We call (9.12) or (9.13) a **homoskedastic Wald statistic** as it is an appropriate test when the errors are conditionally homoskedastic.

As for  $W$ , when  $q = 1$  then  $W^0 = T^2$ , the square of the t-statistic where the latter is computed with a homoskedastic standard error.

**Theorem 9.10.1** Under Assumptions 7.1.2 and 7.10.1,  $\mathbb{E}(e_i^2 | \mathbf{x}_i) = \sigma^2$ , and  $\mathbb{H}_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$ , then

$$W^0 \xrightarrow{d} \chi_q^2,$$

and for  $c$  satisfying  $\alpha = 1 - G_q(c)$ ,

$$\Pr(W^0 > c | \mathbb{H}_0) \longrightarrow \alpha$$

so the test “Reject  $\mathbb{H}_0$  if  $W^0 > c$ ” has asymptotic size  $\alpha$ .

## 9.11 Criterion-Based Tests

The Wald statistic is based on the length of the vector  $\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0$ : the discrepancy between the estimate  $\hat{\boldsymbol{\theta}} = \mathbf{r}(\hat{\boldsymbol{\beta}})$  and the hypothesized value  $\boldsymbol{\theta}_0$ . An alternative class of tests is based on the discrepancy between the criterion function minimized with and without the restriction.

Criterion-based testing applies when we have a criterion function, say  $J(\boldsymbol{\beta})$  with  $\boldsymbol{\beta} \in \mathbf{B}$ , which is minimized for estimation, and the goal is to test  $\mathbb{H}_0 : \boldsymbol{\beta} \in \mathbf{B}_0$  versus  $\mathbb{H}_1 : \boldsymbol{\beta} \notin \mathbf{B}_0$  where  $\mathbf{B}_0 \subset \mathbf{B}$ . Minimizing the criterion function over  $\mathbf{B}$  and  $\mathbf{B}_0$  we obtain the unrestricted and restricted estimators

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= \operatorname{argmin}_{\boldsymbol{\beta} \in \mathbf{B}} J(\boldsymbol{\beta}) \\ \tilde{\boldsymbol{\beta}} &= \operatorname{argmin}_{\boldsymbol{\beta} \in \mathbf{B}_0} J(\boldsymbol{\beta}).\end{aligned}$$

The **criterion-based statistic** for  $\mathbb{H}_0$  versus  $\mathbb{H}_1$  is proportional to

$$\begin{aligned}J &= \min_{\boldsymbol{\beta} \in \mathbf{B}_0} J(\boldsymbol{\beta}) - \min_{\boldsymbol{\beta} \in \mathbf{B}} J(\boldsymbol{\beta}) \\ &= J(\tilde{\boldsymbol{\beta}}) - J(\hat{\boldsymbol{\beta}}).\end{aligned}$$

The criterion-based statistic  $J$  is sometimes called a **distance** statistic, a **minimum-distance** statistic, or a **likelihood-ratio-like** statistic.

Since  $\mathbf{B}_0$  is a subset of  $\mathbf{B}$ ,  $J(\tilde{\boldsymbol{\beta}}) \geq J(\hat{\boldsymbol{\beta}})$  and thus  $J \geq 0$ . The statistic  $J$  measures the cost (on the criterion) of imposing the null restriction  $\boldsymbol{\beta} \in \mathbf{B}_0$ .

## 9.12 Minimum Distance Tests

The minimum distance test is a criterion-based test where  $J(\boldsymbol{\beta})$  is the minimum distance criterion (8.20)

$$J(\boldsymbol{\beta}) = n \left( \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \right)' \widehat{\mathbf{W}} \left( \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \right) \quad (9.14)$$

with  $\hat{\boldsymbol{\beta}}$  the unrestricted (LS) estimator. The restricted estimator  $\tilde{\boldsymbol{\beta}}_{\text{md}}$  minimizes (9.14) subject to  $\boldsymbol{\beta} \in \mathbf{B}_0$ . Observing that  $J(\hat{\boldsymbol{\beta}}) = 0$ , the minimum distance statistic simplifies to

$$J = J(\tilde{\boldsymbol{\beta}}_{\text{md}}) = n \left( \hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}_{\text{md}} \right)' \widehat{\mathbf{W}} \left( \hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}_{\text{md}} \right). \quad (9.15)$$

The efficient minimum distance estimator  $\tilde{\boldsymbol{\beta}}_{\text{emd}}$  is obtained by setting  $\widehat{\mathbf{W}} = \widehat{\mathbf{V}}_{\boldsymbol{\beta}}^{-1}$  in (9.14) and (9.15). The efficient minimum distance statistic for  $\mathbb{H}_0 : \boldsymbol{\beta} \in \mathbf{B}_0$  is therefore

$$J^* = n \left( \hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}_{\text{emd}} \right)' \widehat{\mathbf{V}}_{\boldsymbol{\beta}}^{-1} \left( \hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}_{\text{emd}} \right). \quad (9.16)$$

Consider the class of linear hypotheses  $\mathbb{H}_0 : \mathbf{R}'\boldsymbol{\beta} = \boldsymbol{\theta}_0$ . In this case we know from (8.28) that the efficient minimum distance estimator  $\tilde{\boldsymbol{\beta}}_{\text{emd}}$  subject to the constraint  $\mathbf{R}'\boldsymbol{\beta} = \boldsymbol{\theta}_0$  is

$$\tilde{\boldsymbol{\beta}}_{\text{emd}} = \hat{\boldsymbol{\beta}} - \widehat{\mathbf{V}}_{\boldsymbol{\beta}} \mathbf{R} \left( \mathbf{R}' \widehat{\mathbf{V}}_{\boldsymbol{\beta}} \mathbf{R} \right)^{-1} \left( \mathbf{R}' \hat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0 \right)$$

and thus

$$\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}_{\text{emd}} = \widehat{\mathbf{V}}_{\boldsymbol{\beta}} \mathbf{R} \left( \mathbf{R}' \widehat{\mathbf{V}}_{\boldsymbol{\beta}} \mathbf{R} \right)^{-1} \left( \mathbf{R}' \hat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0 \right).$$

Substituting into (9.16) we find

$$\begin{aligned}
 J^* &= n \left( \mathbf{R}' \hat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0 \right)' \left( \mathbf{R}' \hat{\mathbf{V}}_{\boldsymbol{\beta}} \mathbf{R} \right)^{-1} \mathbf{R}' \hat{\mathbf{V}}_{\boldsymbol{\beta}} \hat{\mathbf{V}}_{\boldsymbol{\beta}}^{-1} \hat{\mathbf{V}}_{\boldsymbol{\beta}} \mathbf{R} \left( \mathbf{R}' \hat{\mathbf{V}}_{\boldsymbol{\beta}} \mathbf{R} \right)^{-1} \left( \mathbf{R}' \hat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0 \right) \\
 &= n \left( \mathbf{R}' \hat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0 \right)' \left( \mathbf{R}' \hat{\mathbf{V}}_{\boldsymbol{\beta}} \mathbf{R} \right)^{-1} \left( \mathbf{R}' \hat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0 \right) \\
 &= W,
 \end{aligned} \tag{9.17}$$

which is the Wald statistic (9.10).

Thus for linear hypotheses  $\mathbb{H}_0 : \mathbf{R}' \boldsymbol{\beta} = \boldsymbol{\theta}_0$ , the efficient minimum distance statistic  $J^*$  is identical to the Wald statistic (9.10). For non-linear hypotheses, however, the Wald and minimum distance statistics are different.

Newey and West (1987) established the asymptotic null distribution of  $J^*$  for linear and non-linear hypotheses.

**Theorem 9.12.1** *Under Assumptions 7.1.2 and 7.10.1, and  $\mathbb{H}_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$ , then  $J^* \xrightarrow{d} \chi_q^2$ .*

Testing using the minimum distance statistic  $J^*$  is similar to testing using the Wald statistic  $W$ . Critical values and p-values are computed using the  $\chi_q^2$  distribution.  $\mathbb{H}_0$  is rejected in favor of  $\mathbb{H}_1$  if  $J^*$  exceeds the level  $\alpha$  critical value, which can be calculated in MATLAB as `chi2inv(1- $\alpha$ ,q)`. The asymptotic p-value is  $p = 1 - G_q(J^*)$ . In MATLAB, use the command `1-chi2cdf(J,q)`.

### 9.13 Minimum Distance Tests Under Homoskedasticity

If we set  $\widehat{\mathbf{W}} = \hat{\mathbf{Q}}_{xx}/s^2$  in (9.14) we obtain the criterion (8.22)

$$J^0(\boldsymbol{\beta}) = n \left( \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \right)' \hat{\mathbf{Q}}_{xx} \left( \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \right) / s^2.$$

A minimum distance statistic for  $\mathbb{H}_0 : \boldsymbol{\beta} \in \mathbf{B}_0$  is

$$J^0 = \min_{\boldsymbol{\beta} \in \mathbf{B}_0} J^0(\boldsymbol{\beta}).$$

Equation (8.23) showed that

$$SSE(\boldsymbol{\beta}) = n\hat{\sigma}^2 + s^2 J^0(\boldsymbol{\beta})$$

and so the minimizers of  $SSE(\boldsymbol{\beta})$  and  $J^0(\boldsymbol{\beta})$  are identical. Thus the constrained minimizer of  $J^0(\boldsymbol{\beta})$  is constrained least-squares

$$\tilde{\boldsymbol{\beta}}_{\text{cls}} = \underset{\boldsymbol{\beta} \in \mathbf{B}_0}{\operatorname{argmin}} J^0(\boldsymbol{\beta}) = \underset{\boldsymbol{\beta} \in \mathbf{B}_0}{\operatorname{argmin}} SSE(\boldsymbol{\beta}) \tag{9.18}$$

and therefore

$$\begin{aligned}
 J_n^0 &= J_n^0(\tilde{\boldsymbol{\beta}}_{\text{cls}}) \\
 &= n \left( \hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}_{\text{cls}} \right)' \hat{\mathbf{Q}}_{xx} \left( \hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}_{\text{cls}} \right) / s^2.
 \end{aligned}$$

In the special case of linear hypotheses  $\mathbb{H}_0 : \mathbf{R}' \boldsymbol{\beta} = \boldsymbol{\theta}_0$ , the constrained least-squares estimator subject to  $\mathbf{R}' \boldsymbol{\beta} = \boldsymbol{\theta}_0$  has the solution (8.10)

$$\tilde{\boldsymbol{\beta}}_{\text{cls}} = \hat{\boldsymbol{\beta}} - \hat{\mathbf{Q}}_{xx}^{-1} \mathbf{R} \left( \mathbf{R}' \hat{\mathbf{Q}}_{xx}^{-1} \mathbf{R} \right)^{-1} \left( \mathbf{R}' \hat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0 \right)$$

and solving we find

$$J^0 = n \left( \mathbf{R}'\hat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0 \right)' \left( \mathbf{R}'\hat{\mathbf{Q}}_{xx}^{-1}\mathbf{R} \right)^{-1} \left( \mathbf{R}'\hat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0 \right) / s^2 = W^0. \quad (9.19)$$

This is the homoskedastic Wald statistic (9.13). Thus for testing linear hypotheses, homoskedastic minimum distance and Wald statistics agree.

For nonlinear hypotheses they disagree, but have the same null asymptotic distribution.

**Theorem 9.13.1** *Under Assumptions 7.1.2 and 7.10.1,  $\mathbb{E}(e_i^2 | \mathbf{x}_i) = \sigma^2$ , and  $\mathbb{H}_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$ , then  $J^0 \xrightarrow{d} \chi_q^2$ .*

## 9.14 F Tests

In Section 5.15 we introduced the  $F$  test for exclusion restrictions in the normal regression model. More generally, the  $F$  statistic for testing  $\mathbb{H}_0 : \boldsymbol{\beta} \in \mathbf{B}_0$  is

$$F = \frac{(\tilde{\sigma}^2 - \hat{\sigma}^2) / q}{\hat{\sigma}^2 / (n - k)} \quad (9.20)$$

where

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \left( y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}} \right)^2$$

and  $\hat{\boldsymbol{\beta}}$  are the unconstrained estimators of  $\boldsymbol{\beta}$  and  $\sigma^2$ ,

$$\tilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \left( y_i - \mathbf{x}_i' \tilde{\boldsymbol{\beta}}_{\text{cls}} \right)^2$$

and  $\tilde{\boldsymbol{\beta}}_{\text{cls}}$  are the constrained least-squares estimators from (9.18),  $q$  is the number of restrictions, and  $k$  is the number of unconstrained coefficients.

We can alternatively write

$$F = \frac{SSE(\tilde{\boldsymbol{\beta}}_{\text{cls}}) - SSE(\hat{\boldsymbol{\beta}})}{qs^2} \quad (9.21)$$

where

$$SSE(\boldsymbol{\beta}) = \sum_{i=1}^n \left( y_i - \mathbf{x}_i' \boldsymbol{\beta} \right)^2$$

is the sum-of-squared errors. Thus  $F$  is a criterion-based statistic. Using (8.23) we can also write  $F$  as

$$F = J^0 / q,$$

so the  $F$  statistic is identical to the homoskedastic minimum distance statistic divided by the number of restrictions  $q$ .

As we discussed in the previous section, in the special case of linear hypotheses  $\mathbb{H}_0 : \mathbf{R}'\boldsymbol{\beta} = \boldsymbol{\theta}_0$ ,  $J^0 = W^0$ . It follows that in this case  $F = W^0 / q$ . Thus for linear restrictions the  $F$  statistic equals the homoskedastic Wald statistic divided by  $q$ . It follows that they are equivalent tests for  $\mathbb{H}_0$  against  $\mathbb{H}_1$ .



**Theorem 9.14.1** For tests of linear hypotheses  $\mathbb{H}_0 : \mathbf{R}'\boldsymbol{\beta} = \boldsymbol{\theta}_0$ ,

$$F = W^0/q$$

the  $F$  statistic equals the homoskedastic Wald statistic divided by the degrees of freedom. Thus under 7.1.2 and 7.10.1,  $\mathbb{E}(e_i^2 | \mathbf{x}_i) = \sigma^2$ , and  $\mathbb{H}_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$ , then

$$F \xrightarrow{d} \chi_q^2/q.$$

When using an  $F$  statistic, it is conventional to use the  $F_{q,n-k}$  distribution for critical values and p-values. Critical values are given in MATLAB by `finv(1- $\alpha$ , $q$ , $n-k$ )`, and p-values by `1-fcdf(F, $q$ , $n-k$ )`. Alternatively, the  $\chi_q^2/q$  distribution can be used, using `chi2inv(1- $\alpha$ , $q$ )/ $q$`  and `1-chi2cdf(F* $q$ , $q$ )`, respectively. Using the  $F_{q,n-k}$  distribution is a prudent small sample adjustment which yields exact answers if the errors are normal, and otherwise slightly increasing the critical values and p-values relative to the asymptotic approximation. Once again, if the sample size is small enough that the choice makes a difference, then probably we shouldn't be trusting the asymptotic approximation anyway!

An elegant feature about (9.20) or (9.21) is that they are directly computable from the standard output from two simple OLS regressions, as the sum of squared errors (or regression variance) is a typical printed output from statistical packages, and is often reported in applied tables. Thus  $F$  can be calculated by hand from standard reported statistics even if you don't have the original data (or if you are sitting in a seminar and listening to a presentation!).

If you are presented with an  $F$  statistic (or a Wald statistic, as you can just divide by  $q$ ) but don't have access to critical values, a useful rule of thumb is to know that for large  $n$ , the 5% asymptotic critical value is decreasing as  $q$  increases, and is less than 2 for  $q \geq 7$ .

A word of warning: In many statistical packages, when an OLS regression is estimated an “ $F$ -statistic” is automatically reported, even though no hypothesis test was requested. What the package is reporting is an  $F$  statistic of the hypothesis that all slope coefficients<sup>1</sup> are zero. This was a popular statistic in the early days of econometric reporting when sample sizes were very small and researchers wanted to know if there was “any explanatory power” to their regression. This is rarely an issue today, as sample sizes are typically sufficiently large that this  $F$  statistic is nearly always highly significant. While there are special cases where this  $F$  statistic is useful, these cases are not typical. As a general rule, there is no reason to report this  $F$  statistic.

## 9.15 Hausman Tests

Hausman (1978) introduced a general idea about how to test a hypothesis  $\mathbb{H}_0$ . If you have two estimators, one which is efficient under  $\mathbb{H}_0$  but inconsistent under  $\mathbb{H}_1$ , and another which is consistent under  $\mathbb{H}_1$ , then construct a test as a quadratic form in the differences of the estimators. In the case of testing a hypothesis  $\mathbb{H}_0 : \mathbf{r}(\boldsymbol{\beta}) = \boldsymbol{\theta}_0$  let  $\hat{\boldsymbol{\beta}}_{\text{ols}}$  denote the unconstrained least-squares estimator and let  $\tilde{\boldsymbol{\beta}}_{\text{emd}}$  denote the efficient minimum distance estimator which imposes  $\mathbf{r}(\boldsymbol{\beta}) = \boldsymbol{\theta}_0$ . Both estimators are consistent under  $\mathbb{H}_0$ , but  $\tilde{\boldsymbol{\beta}}_{\text{emd}}$  is asymptotically efficient. Under  $\mathbb{H}_1$ ,  $\hat{\boldsymbol{\beta}}_{\text{ols}}$  is consistent for  $\boldsymbol{\beta}$  but  $\tilde{\boldsymbol{\beta}}_{\text{emd}}$  is inconsistent. The difference has the asymptotic distribution

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_{\text{ols}} - \tilde{\boldsymbol{\beta}}_{\text{emd}}) \xrightarrow{d} N\left(\mathbf{0}, \mathbf{V}_{\boldsymbol{\beta}} \mathbf{R} (\mathbf{R}' \mathbf{V}_{\boldsymbol{\beta}} \mathbf{R})^{-1} \mathbf{R}' \mathbf{V}_{\boldsymbol{\beta}}\right).$$

<sup>1</sup>All coefficients except the intercept.

Let  $\mathbf{A}^-$  denote the Moore-Penrose generalized inverse. The Hausman statistic for  $\mathbb{H}_0$  is

$$\begin{aligned} H &= \left( \hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{emd}} \right)' \widehat{\text{avar}} \left( \hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{emd}} \right)^- \left( \hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{emd}} \right) \\ &= n \left( \hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{emd}} \right)' \left( \hat{\mathbf{V}}_{\beta} \hat{\mathbf{R}} \left( \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\beta} \hat{\mathbf{R}} \right)^{-1} \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\beta} \right)^- \left( \hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{emd}} \right). \end{aligned}$$

The matrix  $\hat{\mathbf{V}}_{\beta}^{1/2} \hat{\mathbf{R}} \left( \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\beta} \hat{\mathbf{R}} \right)^{-1} \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\beta}^{1/2}$  idempotent so its generalized inverse is itself. (See Section ??.) It follows that

$$\begin{aligned} \left( \hat{\mathbf{V}}_{\beta} \hat{\mathbf{R}} \left( \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\beta} \hat{\mathbf{R}} \right)^{-1} \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\beta} \right)^- &= \hat{\mathbf{V}}_{\beta}^{-1/2} \left( \hat{\mathbf{V}}_{\beta}^{1/2} \hat{\mathbf{R}} \left( \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\beta} \hat{\mathbf{R}} \right)^{-1} \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\beta}^{1/2} \right)^- \hat{\mathbf{V}}_{\beta}^{-1/2} \\ &= \hat{\mathbf{V}}_{\beta}^{-1/2} \hat{\mathbf{V}}_{\beta}^{1/2} \hat{\mathbf{R}} \left( \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\beta} \hat{\mathbf{R}} \right)^{-1} \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\beta}^{1/2} \hat{\mathbf{V}}_{\beta}^{-1/2} \\ &= \hat{\mathbf{R}} \left( \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\beta} \hat{\mathbf{R}} \right)^{-1} \hat{\mathbf{R}}'. \end{aligned}$$

Thus the Hausman statistic is

$$H = n \left( \hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{emd}} \right)' \hat{\mathbf{R}} \left( \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\beta} \hat{\mathbf{R}} \right)^{-1} \hat{\mathbf{R}}' \left( \hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{emd}} \right),$$

In the context of linear restrictions,  $\hat{\mathbf{R}} = \mathbf{R}$  and  $\mathbf{R}'\tilde{\beta} = \theta_0$  so the statistic takes the form

$$H = n \left( \mathbf{R}'\hat{\beta}_{\text{ols}} - \theta_0 \right)' \hat{\mathbf{R}} \left( \mathbf{R}'\hat{\mathbf{V}}_{\beta} \mathbf{R} \right)^{-1} \left( \mathbf{R}'\hat{\beta}_{\text{ols}} - \theta_0 \right),$$

which is precisely the Wald statistic. With nonlinear restrictions then can differ.

In either case we see that the asymptotic null distribution of the Hausman statistic  $H$  is  $\chi_q^2$ , so the appropriate test is to reject  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  if  $H > c$  where  $c$  is a critical value taken from the  $\chi_q^2$  distribution.

**Theorem 9.15.1** *For general hypotheses the Hausman test statistic is*

$$H = n \left( \hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{emd}} \right)' \hat{\mathbf{R}} \left( \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\beta} \hat{\mathbf{R}} \right)^{-1} \hat{\mathbf{R}}' \left( \hat{\beta}_{\text{ols}} - \tilde{\beta}_{\text{emd}} \right),$$

*and has the asymptotic distribution under  $\mathbb{H}_0 : \mathbf{r}(\beta) = \theta_0$ ,*

$$H \xrightarrow{d} \chi_q^2.$$

### Jerry Hausman

Jerry Hausman (1946- ) of the United States is a leading micro-econometrician, best known for his influential contributions on specification testing and panel data.

## 9.16 Score Tests

Score tests are traditionally derived in likelihood analysis, but can more generally be constructed from first-order conditions evaluated at restricted estimates. We focus on the likelihood derivation.

Given the log likelihood function  $\log L(\boldsymbol{\beta}, \sigma^2)$ , a restriction  $\mathbb{H}_0 : \mathbf{r}(\boldsymbol{\beta}) = \mathbf{0}$ , and restricted estimators  $\tilde{\boldsymbol{\beta}}$  and  $\tilde{\sigma}^2$ , the **score statistic** for  $\mathbb{H}_0$  is defined as

$$S = \left( \frac{\partial}{\partial \boldsymbol{\beta}} \log L(\tilde{\boldsymbol{\beta}}, \tilde{\sigma}^2) \right)' \left( -\frac{\partial^2}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} \log L(\tilde{\boldsymbol{\beta}}, \tilde{\sigma}^2) \right)^{-1} \left( \frac{\partial}{\partial \boldsymbol{\beta}} \log L(\tilde{\boldsymbol{\beta}}, \tilde{\sigma}^2) \right).$$

The idea is that if the restriction is true, then the restricted estimators should be close to the maximum of the log-likelihood where the derivative should be small. However if the restriction is false then the restricted estimators should be distant from the maximum and the derivative should be large. Hence small values of  $S$  are expected under  $\mathbb{H}_0$  and large values under  $\mathbb{H}_1$ . Tests of  $\mathbb{H}_0$  thus reject for large values of  $S$ .

We explore the score statistic in the context of the normal regression model and linear hypotheses  $\mathbf{r}(\boldsymbol{\beta}) = \mathbf{R}'\boldsymbol{\beta}$ . Recall that in the normal regression log-likelihood function is

$$\log L(\boldsymbol{\beta}, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i'\boldsymbol{\beta})^2.$$

The constrained MLE under linear hypotheses is constrained least squares

$$\begin{aligned} \tilde{\boldsymbol{\beta}} &= \hat{\boldsymbol{\beta}} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R} \left[ \mathbf{R}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R} \right]^{-1} (\mathbf{R}'\hat{\boldsymbol{\beta}} - \mathbf{c}) \\ \tilde{e}_i &= y_i - \mathbf{x}_i'\tilde{\boldsymbol{\beta}} \\ \tilde{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n \tilde{e}_i^2. \end{aligned}$$

We can calculate that the derivative and Hessian are

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\beta}} \log L(\tilde{\boldsymbol{\beta}}, \tilde{\sigma}^2) &= \frac{1}{\tilde{\sigma}^2} \sum_{i=1}^n \mathbf{x}_i (y_i - \mathbf{x}_i'\tilde{\boldsymbol{\beta}}) = \frac{1}{\tilde{\sigma}^2} \mathbf{X}'\tilde{\mathbf{e}} \\ -\frac{\partial^2}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} \log L(\tilde{\boldsymbol{\beta}}, \tilde{\sigma}^2) &= \frac{1}{\tilde{\sigma}^2} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' = \frac{1}{\tilde{\sigma}^2} \mathbf{X}'\mathbf{X}. \end{aligned}$$

Since  $\tilde{\mathbf{e}} = \mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}$  we can further calculate that

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\beta}} \log L(\tilde{\boldsymbol{\beta}}, \tilde{\sigma}^2) &= \frac{1}{\tilde{\sigma}^2} (\mathbf{X}'\mathbf{X}) \left( (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X}\tilde{\boldsymbol{\beta}} \right) \\ &= \frac{1}{\tilde{\sigma}^2} (\mathbf{X}'\mathbf{X}) (\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}) \\ &= \frac{1}{\tilde{\sigma}^2} \mathbf{R} \left[ \mathbf{R}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R} \right]^{-1} (\mathbf{R}'\hat{\boldsymbol{\beta}} - \mathbf{c}). \end{aligned}$$

Together we find that

$$S = (\mathbf{R}'\hat{\boldsymbol{\beta}} - \mathbf{c})' \left( \mathbf{R}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R} \right)^{-1} (\mathbf{R}'\hat{\boldsymbol{\beta}} - \mathbf{c}) / \tilde{\sigma}^2$$

This is identical to the homoskedastic Wald statistic, with  $s^2$  replaced by  $\tilde{\sigma}^2$ . We can also write  $S$  as a monotonic transformation of the  $F$  statistic, since

$$S = n \frac{(\tilde{\sigma}^2 - \hat{\sigma}^2)}{\hat{\sigma}^2} = n \left( 1 - \frac{\hat{\sigma}^2}{\tilde{\sigma}^2} \right) = n \left( 1 - \frac{1}{1 + \frac{q}{n-k} F} \right).$$

The test “Reject  $\mathbb{H}_0$  for large values of  $S$ ” is identical to the test “Reject  $\mathbb{H}_0$  for large values of  $F$ ”, so they are identical tests. Since for the normal regression model the exact distribution of  $F$  is known, it is better to use the  $F$  statistic with  $F$  p-values.

In more complicated settings a potential advantage of score tests is that they are calculated using the restricted parameter estimates  $\tilde{\beta}$  rather than the unrestricted estimates  $\hat{\beta}$ . Thus when  $\tilde{\beta}$  is relatively easy to calculate there can be a preference for score statistics. This is not a concern for linear restrictions.

More generally, score and score-like statistics can be constructed from first-order conditions evaluated at restricted parameter estimates. Also, when test statistics are constructed using covariance matrix estimators which are calculated using restricted parameter estimates (e.g. restricted residuals) then these are often described as score tests.

An example of the latter is the Wald-type statistic

$$W = \left( r(\hat{\beta}) - \theta_0 \right)' \left( \hat{R}' \tilde{V}_{\hat{\beta}} \hat{R} \right)^{-1} \left( r(\hat{\beta}) - \theta_0 \right)$$

where the covariance matrix estimate  $\tilde{V}_{\hat{\beta}}$  is calculated using the restricted residuals  $\tilde{e}_i = y_i - \mathbf{x}_i' \tilde{\beta}$ . This may be done when  $\beta$  and  $\theta$  are high-dimensional, so there is worry that the estimator  $\hat{V}_{\hat{\beta}}$  is imprecise.

## 9.17 Problems with Tests of Nonlinear Hypotheses

While the t and Wald tests work well when the hypothesis is a linear restriction on  $\beta$ , they can work quite poorly when the restrictions are nonlinear. This can be seen by a simple example introduced by Lafontaine and White (1986). Take the model

$$\begin{aligned} y_i &= \beta + e_i \\ e_i &\sim N(0, \sigma^2) \end{aligned}$$

and consider the hypothesis

$$\mathbb{H}_0 : \beta = 1.$$

Let  $\hat{\beta}$  and  $\hat{\sigma}^2$  be the sample mean and variance of  $y_i$ . The standard Wald test for  $\mathbb{H}_0$  is

$$W = n \frac{(\hat{\beta} - 1)^2}{\hat{\sigma}^2}.$$

Now notice that  $\mathbb{H}_0$  is equivalent to the hypothesis

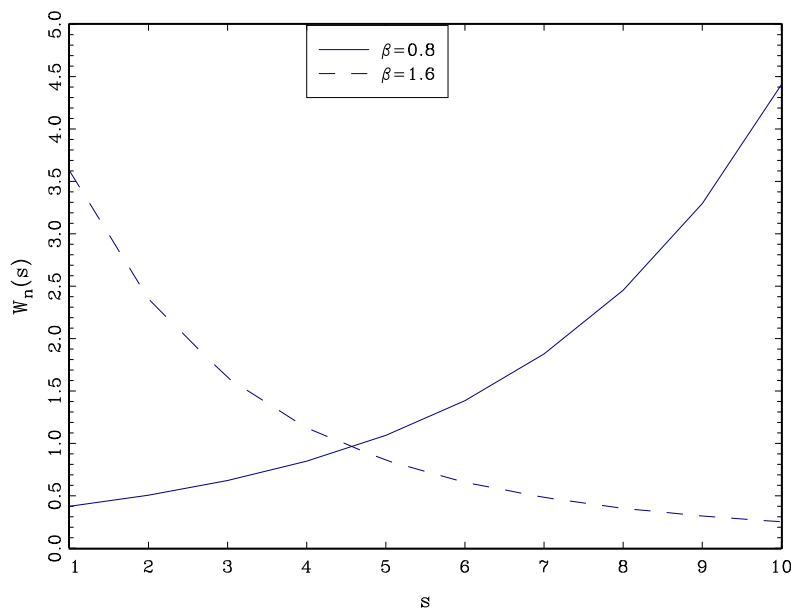
$$\mathbb{H}_0(s) : \beta^s = 1$$

for any positive integer  $s$ . Letting  $r(\beta) = \beta^s$ , and noting  $\mathbf{R} = s\beta^{s-1}$ , we find that the standard Wald test for  $\mathbb{H}_0(s)$  is

$$W(s) = n \frac{(\hat{\beta}^s - 1)^2}{\hat{\sigma}^2 s^2 \hat{\beta}^{2s-2}}.$$

While the hypothesis  $\beta^s = 1$  is unaffected by the choice of  $s$ , the statistic  $W(s)$  varies with  $s$ . This is an unfortunate feature of the Wald statistic.

To demonstrate this effect, we have plotted in Figure 9.1 the Wald statistic  $W(s)$  as a function of  $s$ , setting  $n/\hat{\sigma}^2 = 10$ . The increasing solid line is for the case  $\hat{\beta} = 0.8$ . The decreasing dashed line is for the case  $\hat{\beta} = 1.6$ . It is easy to see that in each case there are values of  $s$  for which the test statistic is significant relative to asymptotic critical values, while there are other values of  $s$

Figure 9.1: Wald Statistic as a function of  $s$ 

for which the test statistic is insignificant. This is distressing since the choice of  $s$  is arbitrary and irrelevant to the actual hypothesis.

Our first-order asymptotic theory is not useful to help pick  $s$ , as  $W(s) \xrightarrow{d} \chi_1^2$  under  $\mathbb{H}_0$  for any  $s$ . This is a context where **Monte Carlo simulation** can be quite useful as a tool to study and compare the exact distributions of statistical procedures in finite samples. The method uses random simulation to create artificial datasets, to which we apply the statistical tools of interest. This produces random draws from the statistic's sampling distribution. Through repetition, features of this distribution can be calculated.

In the present context of the Wald statistic, one feature of importance is the Type I error of the test using the asymptotic 5% critical value 3.84 – the probability of a false rejection,  $\Pr(W(s) > 3.84 \mid \beta = 1)$ . Given the simplicity of the model, this probability depends only on  $s$ ,  $n$ , and  $\sigma^2$ . In Table 9.1 we report the results of a Monte Carlo simulation where we vary these three parameters. The value of  $s$  is varied from 1 to 10,  $n$  is varied among 20, 100 and 500, and  $\sigma$  is varied among 1 and 3. The Table reports the simulation estimate of the Type I error probability from 50,000 random samples. Each row of the table corresponds to a different value of  $s$  – and thus corresponds to a particular choice of test statistic. The second through seventh columns contain the Type I error probabilities for different combinations of  $n$  and  $\sigma$ . These probabilities are calculated as the percentage of the 50,000 simulated Wald statistics  $W(s)$  which are larger than 3.84. The null hypothesis  $\beta^s = 1$  is true, so these probabilities are Type I error.

To interpret the table, remember that the ideal Type I error probability is 5% (.05) with deviations indicating distortion. Type I error rates between 3% and 8% are considered reasonable. Error rates above 10% are considered excessive. Rates above 20% are unacceptable. When comparing statistical procedures, we compare the rates row by row, looking for tests for which rejection rates are close to 5% and rarely fall outside of the 3%-8% range. For this particular example the only test which meets this criterion is the conventional  $W = W(1)$  test. Any other choice of  $s$  leads to a test with unacceptable Type I error probabilities.

Table 9.1  
Type I Error Probability of Asymptotic 5%  $W(s)$  Test

$s$	$\sigma = 1$			$\sigma = 3$		
	$n = 20$	$n = 100$	$n = 500$	$n = 20$	$n = 100$	$n = 500$
1	.06	.05	.05	.07	.05	.05
2	.08	.06	.05	.15	.08	.06
3	.10	.06	.05	.21	.12	.07
4	.13	.07	.06	.25	.15	.08
5	.15	.08	.06	.28	.18	.10
6	.17	.09	.06	.30	.20	.11
7	.19	.10	.06	.31	.22	.13
8	.20	.12	.07	.33	.24	.14
9	.22	.13	.07	.34	.25	.15
10	.23	.14	.08	.35	.26	.16

Note: Rejection frequencies from 50,000 simulated random samples

In Table 9.1 you can also see the impact of variation in sample size. In each case, the Type I error probability improves towards 5% as the sample size  $n$  increases. There is, however, no magic choice of  $n$  for which all tests perform uniformly well. Test performance deteriorates as  $s$  increases, which is not surprising given the dependence of  $W(s)$  on  $s$  as shown in Figure 9.1.

In this example it is not surprising that the choice  $s = 1$  yields the best test statistic. Other choices are arbitrary and would not be used in practice. While this is clear in this particular example, in other examples natural choices are not always obvious and the best choices may in fact appear counter-intuitive at first.

This point can be illustrated through another example which is similar to one developed in Gregory and Veall (1985). Take the model

$$\begin{aligned} y_i &= \beta_0 + x_{1i}\beta_1 + x_{2i}\beta_2 + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= \mathbf{0} \end{aligned} \tag{9.22}$$

and the hypothesis

$$\mathbb{H}_0 : \frac{\beta_1}{\beta_2} = \theta_0$$

where  $\theta_0$  is a known constant. Equivalently, define  $\theta = \beta_1/\beta_2$ , so the hypothesis can be stated as  $\mathbb{H}_0 : \theta = \theta_0$ .

Let  $\hat{\boldsymbol{\beta}} = (\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2)$  be the least-squares estimates of (9.22), let  $\hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}$  be an estimate of the covariance matrix for  $\hat{\boldsymbol{\beta}}$  and set  $\hat{\theta} = \hat{\beta}_1/\hat{\beta}_2$ . Define

$$\hat{\mathbf{R}}_1 = \begin{pmatrix} 0 \\ \frac{1}{\hat{\beta}_2} \\ -\frac{\hat{\beta}_1}{\hat{\beta}_2^2} \end{pmatrix}$$

so that the standard error for  $\hat{\theta}$  is  $s(\hat{\theta}) = \left( \hat{\mathbf{R}}_1' \hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}} \hat{\mathbf{R}}_1 \right)^{1/2}$ . In this case a t-statistic for  $\mathbb{H}_0$  is

$$T_1 = \frac{\left( \frac{\hat{\beta}_1}{\hat{\beta}_2} - \theta_0 \right)}{s(\hat{\theta})}.$$

An alternative statistic can be constructed through reformulating the null hypothesis as

$$\mathbb{H}_0 : \beta_1 - \theta_0 \beta_2 = 0.$$

A t-statistic based on this formulation of the hypothesis is

$$T_2 = \frac{\hat{\beta}_1 - \theta_0 \hat{\beta}_2}{\left(\mathbf{R}_2' \hat{\mathbf{V}}_{\hat{\beta}} \mathbf{R}_2\right)^{1/2}}.$$

where

$$\mathbf{R}_2 = \begin{pmatrix} 0 \\ 1 \\ -\theta_0 \end{pmatrix}.$$

To compare  $T_1$  and  $T_2$  we perform another simple Monte Carlo simulation. We let  $x_{1i}$  and  $x_{2i}$  be mutually independent  $N(0, 1)$  variables,  $e_i$  be an independent  $N(0, \sigma^2)$  draw with  $\sigma = 3$ , and normalize  $\beta_0 = 0$  and  $\beta_1 = 1$ . This leaves  $\beta_2$  as a free parameter, along with sample size  $n$ . We vary  $\beta_2$  among .1, .25, .50, .75, and 1.0 and  $n$  among 100 and 500.

Table 9.2  
Type I Error Probability of Asymptotic 5% t-tests

	$n = 100$				$n = 500$			
	$\Pr(T < -1.645)$		$\Pr(T > 1.645)$		$\Pr(T < -1.645)$		$\Pr(T > 1.645)$	
$\beta_2$	$T_1$	$T_2$	$T_1$	$T_2$	$T_1$	$T_2$	$T_1$	$T_2$
.10	.47	.06	.00	.06	.28	.05	.00	.05
.25	.26	.06	.00	.06	.15	.05	.00	.05
.50	.15	.06	.00	.06	.10	.05	.00	.05
.75	.12	.06	.00	.06	.09	.05	.00	.05
1.00	.10	.06	.00	.06	.07	.05	.02	.05

The one-sided Type I error probabilities  $\Pr(T < -1.645)$  and  $\Pr(T > 1.645)$  are calculated from 50,000 simulated samples. The results are presented in Table 9.2. Ideally, the entries in the table should be 0.05. However, the rejection rates for the  $T_1$  statistic diverge greatly from this value, especially for small values of  $\beta_2$ . The left tail probabilities  $\Pr(T_1 < -1.645)$  greatly exceed 5%, while the right tail probabilities  $\Pr(T_1 > 1.645)$  are close to zero in most cases. In contrast, the rejection rates for the linear  $T_2$  statistic are invariant to the value of  $\beta_2$ , and are close to the ideal 5% rate for both sample sizes. The implication of Table 8.2 is that the two t-ratios have dramatically different sampling behavior.

The common message from both examples is that Wald statistics are sensitive to the algebraic formulation of the null hypothesis.

A simple solution is to use the minimum distance statistic  $J$ , which equals  $W$  with  $r = 1$  in the first example, and  $|T_2|$  in the second example. The minimum distance statistic is invariant to the algebraic formulation of the null hypothesis, so is immune to this problem. Whenever possible, the Wald statistic should not be used to test nonlinear hypotheses.

## 9.18 Monte Carlo Simulation

In Section 9.17 we introduced the method of Monte Carlo simulation to illustrate the small sample problems with tests of nonlinear hypotheses. In this section we describe the method in more detail.

Recall, our data consist of observations  $(y_i, \mathbf{x}_i)$  which are random draws from a population distribution  $F$ . Let  $\boldsymbol{\theta}$  be a parameter and let  $T = T((y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n), \boldsymbol{\theta})$  be a statistic of interest, for example an estimator  $\hat{\theta}$  or a t-statistic  $(\hat{\theta} - \theta)/s(\hat{\theta})$ . The exact distribution of  $T$  is

$$G(u, F) = \Pr(T \leq u \mid F).$$

While the asymptotic distribution of  $T$  might be known, the exact (finite sample) distribution  $G$  is generally unknown.

Monte Carlo simulation uses numerical simulation to compute  $G(u, F)$  for selected choices of  $F$ . This is useful to investigate the performance of the statistic  $T$  in reasonable situations and sample sizes. The basic idea is that for any given  $F$ , the distribution function  $G(u, F)$  can be calculated numerically through simulation. The name Monte Carlo derives from the famous Mediterranean gambling resort where games of chance are played.

The method of Monte Carlo is quite simple to describe. The researcher chooses  $F$  (the distribution of the data) and the sample size  $n$ . A “true” value of  $\theta$  is implied by this choice, or equivalently the value  $\theta$  is selected directly by the researcher which implies restrictions on  $F$ .

Then the following experiment is conducted by computer simulation:

1.  $n$  independent random pairs  $(y_i^*, \mathbf{x}_i^*)$ ,  $i = 1, \dots, n$ , are drawn from the distribution  $F$  using the computer’s random number generator.
2. The statistic  $T = T((y_1^*, \mathbf{x}_1^*), \dots, (y_n^*, \mathbf{x}_n^*), \theta)$  is calculated on this pseudo data.

For step 1, computer packages have built-in random number procedures including  $U[0, 1]$  and  $N(0, 1)$ . From these most random variables can be constructed. (For example, a chi-square can be generated by sums of squares of normals.)

For step 2, it is important that the statistic be evaluated at the “true” value of  $\theta$  corresponding to the choice of  $F$ .

The above experiment creates one random draw from the distribution  $G(u, F)$ . This is one observation from an unknown distribution. Clearly, from one observation very little can be said. So the researcher repeats the experiment  $B$  times, where  $B$  is a large number. Typically, we set  $B = 1000$  or  $B = 5000$ . We will discuss this choice later.

Notationally, let the  $b^{th}$  experiment result in the draw  $T_b$ ,  $b = 1, \dots, B$ . These results are stored. After all  $B$  experiments have been calculated, these results constitute a random sample of size  $B$  from the distribution of  $G(u, F) = \Pr(T_b \leq u) = \Pr(T \leq u \mid F)$ .

From a random sample, we can estimate any feature of interest using (typically) a method of moments estimator. We now describe some specific examples.

Suppose we are interested in the bias, mean-squared error (MSE), and/or variance of the distribution of  $\hat{\theta} - \theta$ . We then set  $T = \hat{\theta} - \theta$ , run the above experiment, and calculate

$$\begin{aligned}\widehat{Bias}(\hat{\theta}) &= \frac{1}{B} \sum_{b=1}^B T_b = \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b - \theta \\ \widehat{MSE}(\hat{\theta}) &= \frac{1}{B} \sum_{b=1}^B (T_b)^2 = \frac{1}{B} \sum_{b=1}^B (\hat{\theta}_b - \theta)^2 \\ \widehat{var}(\hat{\theta}) &= \widehat{MSE}(\hat{\theta}) - \left( \widehat{Bias}(\hat{\theta}) \right)^2\end{aligned}$$

Suppose we are interested in the Type I error associated with an asymptotic 5% two-sided t-test. We would then set  $T = \left| \hat{\theta} - \theta \right| / s(\hat{\theta})$  and calculate

$$\hat{P} = \frac{1}{B} \sum_{b=1}^B 1(T_b \geq 1.96), \quad (9.23)$$

the percentage of the simulated t-ratios which exceed the asymptotic 5% critical value.

Suppose we are interested in the 5% and 95% quantile of  $T = \hat{\theta}$  or  $T = (\hat{\theta} - \theta) / s(\hat{\theta})$ . We then compute the 5% and 95% sample quantiles of the sample  $\{T_b\}$ . The  $\alpha$  sample quantile is a number



$q_\alpha$  such that  $100\alpha\%$  of the sample are less than  $q_\alpha$ . A simple way to compute sample quantiles is to sort the sample  $\{T_b\}$  from low to high. Then  $q_\alpha$  is the  $N^{th}$  number in this ordered sequence, where  $N = (B + 1)\alpha$ . It is therefore convenient to pick  $B$  so that  $N$  is an integer. For example, if we set  $B = 999$ , then the 5% sample quantile is  $50^{th}$  sorted value and the 95% sample quantile is the  $950^{th}$  sorted value.

The typical purpose of a Monte Carlo simulation is to investigate the performance of a statistical procedure (estimator or test) in realistic settings. Generally, the performance will depend on  $n$  and  $F$ . In many cases, an estimator or test may perform wonderfully for some values, and poorly for others. It is therefore useful to conduct a variety of experiments, for a selection of choices of  $n$  and  $F$ .

As discussed above, the researcher must select the number of experiments,  $B$ . Often this is called the number of **replications**. Quite simply, a larger  $B$  results in more precise estimates of the features of interest of  $G$ , but requires more computational time. In practice, therefore, the choice of  $B$  is often guided by the computational demands of the statistical procedure. Since the results of a Monte Carlo experiment are estimates computed from a random sample of size  $B$ , it is straightforward to calculate standard errors for any quantity of interest. If the standard error is too large to make a reliable inference, then  $B$  will have to be increased.

In particular, it is simple to make inferences about rejection probabilities from statistical tests, such as the percentage estimate reported in (9.23). The random variable  $1(T_b \geq 1.96)$  is iid Bernoulli, equalling 1 with probability  $p = \mathbb{E}(1(T_b \geq 1.96))$ . The average (9.23) is therefore an unbiased estimator of  $p$  with standard error  $s(\hat{p}) = \sqrt{p(1-p)/B}$ . As  $p$  is unknown, this may be approximated by replacing  $p$  with  $\hat{p}$  or with an hypothesized value. For example, if we are assessing an asymptotic 5% test, then we can set  $s(\hat{p}) = \sqrt{(.05)(.95)/B} \simeq .22/\sqrt{B}$ . Hence, standard errors for  $B = 100, 1000$ , and  $5000$ , are, respectively,  $s(\hat{p}) = .022, .007$ , and  $.003$ .

Most papers in econometric methods, and some empirical papers, include the results of Monte Carlo simulations to illustrate the performance of their methods. When extending existing results, it is good practice to start by replicating existing (published) results. This is not exactly possible in the case of simulation results, as they are inherently random. For example suppose a paper investigates a statistical test, and reports a simulated rejection probability of 0.07 based on a simulation with  $B = 100$  replications. Suppose you attempt to replicate this result, and find a rejection probability of 0.03 (again using  $B = 100$  simulation replications). Should you conclude that you have failed in your attempt? Absolutely not! Under the hypothesis that both simulations are identical, you have two independent estimates,  $\hat{p}_1 = 0.07$  and  $\hat{p}_2 = 0.03$ , of a common probability  $p$ . The asymptotic (as  $B \rightarrow \infty$ ) distribution of their difference is  $\sqrt{B}(\hat{p}_1 - \hat{p}_2) \xrightarrow{d} N(0, 2p(1-p))$ , so a standard error for  $\hat{p}_1 - \hat{p}_2 = 0.04$  is  $\hat{s} = \sqrt{2p(1-p)/B} \simeq 0.03$ , using the estimate  $p = (\hat{p}_1 + \hat{p}_2)/2$ . Since the t-ratio  $0.04/0.03 = 1.3$  is not statistically significant, it is incorrect to reject the null hypothesis that the two simulations are identical. The difference between the results  $\hat{p}_1 = 0.07$  and  $\hat{p}_2 = 0.03$  is consistent with random variation.

What should be done? The first mistake was to copy the previous paper's choice of  $B = 100$ . Instead, suppose you set  $B = 5000$ . Suppose you now obtain  $\hat{p}_2 = 0.04$ . Then  $\hat{p}_1 - \hat{p}_2 = 0.03$  and a standard error is  $\hat{s} = \sqrt{p(1-p)(1/100 + 1/5000)} \simeq 0.02$ . Still we cannot reject the hypothesis that the two simulations are different. Even though the estimates (0.07 and 0.04) appear to be quite different, the difficulty is that the original simulation used a very small number of replications ( $B = 100$ ) so the reported estimate is quite imprecise. In this case, it is appropriate to conclude that your results "replicate" the previous study, as there is no statistical evidence to reject the hypothesis that they are equivalent.

Most journals have policies requiring authors to make available their data sets and computer programs required for empirical results. They do not have similar policies regarding simulations. Nevertheless, it is good professional practice to make your simulations available. The best practice is to post your simulation code on your webpage. This invites others to build on and use your results, leading to possible collaboration, citation, and/or advancement.

## 9.19 Confidence Intervals by Test Inversion

There is a close relationship between hypothesis tests and confidence intervals. We observed in Section 7.13 that the standard 95% asymptotic confidence interval for a parameter  $\theta$  is

$$\begin{aligned}\widehat{C} &= \left[ \widehat{\theta} - 1.96 \cdot s(\widehat{\theta}), \quad \widehat{\theta} + 1.96 \cdot s(\widehat{\theta}) \right] \\ &= \{ \theta : |T(\theta)| \leq 1.96 \}.\end{aligned}\tag{9.24}$$

That is, we can describe  $\widehat{C}$  as “The point estimate plus or minus 2 standard errors” or “The set of parameter values not rejected by a two-sided t-test.” The second definition, known as **test statistic inversion** is a general method for finding confidence intervals, and typically produces confidence intervals with excellent properties.

Given a test statistic  $T(\theta)$  and critical value  $c$ , the acceptance region “Accept if  $T(\theta) \leq c$ ” is identical to the confidence interval  $\widehat{C} = \{ \theta : T(\theta) \leq c \}$ . Since the regions are identical, the probability of coverage  $\Pr(\theta \in \widehat{C})$  equals the probability of correct acceptance  $\Pr(\text{Accept}|\theta)$  which is exactly 1 minus the Type I error probability. Thus inverting a test with good Type I error probabilities yields a confidence interval with good coverage probabilities.

Now suppose that the parameter of interest  $\theta = r(\beta)$  is a nonlinear function of the coefficient vector  $\beta$ . In this case the standard confidence interval for  $\theta$  is the set  $\widehat{C}$  as in (9.24) where  $\widehat{\theta} = r(\widehat{\beta})$  is the point estimate and  $s(\widehat{\theta}) = \sqrt{\widehat{R}' \widehat{V}_{\widehat{\beta}} \widehat{R}}$  is the delta method standard error. This confidence interval is inverting the t-test based on the nonlinear hypothesis  $r(\beta) = \theta$ . The trouble is that in Section 9.17 we learned that there is no unique t-statistic for tests of nonlinear hypotheses and that the choice of parameterization matters greatly.

For example, if  $\theta = \beta_1/\beta_2$  then the coverage probability of the standard interval (9.24) is 1 minus the probability of the Type I error, which as shown in Table 8.2 can be far from the nominal 5%.

In this example a good solution is the same as discussed in Section 9.17 – to rewrite the hypothesis as a linear restriction. The hypothesis  $\theta = \beta_1/\beta_2$  is the same as  $\theta\beta_2 = \beta_1$ . The t-statistic for this restriction is

$$T(\theta) = \frac{\widehat{\beta}_1 - \widehat{\beta}_2\theta}{\left(\mathbf{R}' \widehat{V}_{\widehat{\beta}} \mathbf{R}\right)^{1/2}}$$

where

$$\mathbf{R} = \begin{pmatrix} 1 \\ -\theta \end{pmatrix}$$

and  $\widehat{V}_{\widehat{\beta}}$  is the covariance matrix for  $(\widehat{\beta}_1 \ \widehat{\beta}_2)$ . A 95% confidence interval for  $\theta = \beta_1/\beta_2$  is the set of values of  $\theta$  such that  $|T(\theta)| \leq 1.96$ . Since  $\theta$  appears in both the numerator and denominator,  $T(\theta)$  is a non-linear function of  $\theta$  so the easiest method to find the confidence set is by grid search over  $\theta$ .

For example, in the wage equation

$$\log(\text{Wage}) = \beta_1 \text{Experience} + \beta_2 \text{Experience}^2/100 + \dots$$

the highest expected wage occurs at  $\text{Experience} = -50\beta_1/\beta_2$ . From Table 4.1 we have the point estimate  $\widehat{\theta} = 29.8$  and we can calculate the standard error  $s(\widehat{\theta}) = 0.022$  for a 95% confidence interval  $[29.8, 29.9]$ . However, if we instead invert the linear form of the test we can numerically find the interval  $[29.1, 30.6]$  which is much larger. From the evidence presented in Section 9.17 we know the first interval can be quite inaccurate and the second interval is greatly preferred.

## 9.20 Multiple Tests and Bonferroni Corrections

In most applications, economists examine a large number of estimates, test statistics, and p-values. What does it mean (or does it mean anything) if one statistic appears to be “significant” after examining a large number of statistics? This is known as the problem of **multiple testing** or **multiple comparisons**.

To be specific, suppose we examine a set of  $k$  coefficients, standard errors and t-ratios, and consider the “significance” of each statistic. Based on conventional reasoning, for each coefficient we would reject the hypothesis that the coefficient is zero with asymptotic size  $\alpha$  if the absolute t-statistic exceeds the  $1 - \alpha$  critical value of the normal distribution, or equivalently if the p-value for the t-statistic is smaller than  $\alpha$ . If we observe that one of the  $k$  statistics is “significant” based on this criteria, that means that one of the p-values is smaller than  $\alpha$ , or equivalently, that the smallest p-value is smaller than  $\alpha$ . We can then rephrase the question: Under the joint hypothesis that a set of  $k$  hypotheses are all true, what is the probability that the smallest p-value is smaller than  $\alpha$ ? In general, we cannot provide a precise answer to this question, but the Bonferroni correction bounds this probability by  $\alpha k$ . The Bonferroni method furthermore suggests that if we want the familywise error probability (the probability that one of the tests falsely rejects) is bounded below  $\alpha$ , then an appropriate rule is to reject only if the smallest p-value is smaller than  $\alpha/k$ . Equivalently, the Bonferroni familywise p-value is  $k \min_{j \leq k} p_j$ .

Formally, suppose we have  $k$  hypotheses  $\mathbb{H}_j$ ,  $j = 1, \dots, k$ . For each we have a test and associated p-value  $p_j$  with the property that when  $\mathbb{H}_j$  is true  $\lim_{n \rightarrow \infty} \Pr(p_j < \alpha) = \alpha$ . We then observe that among the  $k$  tests, one of the  $k$  will appear “significant” if  $\min_{j \leq k} p_j < \alpha$ . This event can be written as

$$\left\{ \min_{j \leq k} p_j < \alpha \right\} = \bigcup_{j=1}^k \{p_j < \alpha\}.$$

Boole’s inequality states that for any  $k$  events  $A_j$ ,  $\Pr\left(\bigcup_{j=1}^k A_j\right) \leq \sum_{j=1}^k \Pr(A_j)$ . Thus

$$\Pr\left(\min_{j \leq k} p_j < \alpha\right) \leq \sum_{j=1}^k \Pr(p_j < \alpha) \longrightarrow k\alpha$$

as stated. This demonstrates that the familywise rejection probability is at most  $k$  times the individual rejection probability.

Furthermore,

$$\Pr\left(\min_{j \leq k} p_j < \frac{\alpha}{k}\right) \leq \sum_{j=1}^k \Pr\left(p_j < \frac{\alpha}{k}\right) \longrightarrow \alpha.$$

This demonstrates that the family rejection probability can be controlled (bounded below  $\alpha$ ) if each individual test is subjected to the stricter standard that a p-value must be smaller than  $\alpha/k$  to be labeled as “significant.”

To illustrate, suppose we have two coefficient estimates, with individual p-values 0.04 and 0.15. Based on a conventional 5% level, the standard individual tests would suggest that the first coefficient estimate is “significant” but not the second. A Bonferroni 5% test, however, does not reject as it would require that the smallest p-value be smaller than 0.025, which is not the case in this example. Alternatively, the Bonferroni familywise p-value is 0.08, which is not significant at the 5% level.

In contrast, if the two p-values are 0.01 and 0.15, then the Bonferroni familywise p-value is 0.02, which is significant at the 5% level.

## 9.21 Power and Test Consistency

The **power** of a test is the probability of rejecting  $\mathbb{H}_0$  when  $\mathbb{H}_1$  is true.

For simplicity suppose that  $y_i$  is i.i.d.  $N(\theta, \sigma^2)$  with  $\sigma^2$  known, consider the t-statistic  $T(\theta) = \sqrt{n}(\bar{y} - \theta)/\sigma$ , and tests of  $\mathbb{H}_0 : \theta = 0$  against  $\mathbb{H}_1 : \theta > 0$ . We reject  $\mathbb{H}_0$  if  $T = T(0) > c$ . Note that

$$T = T(\theta) + \sqrt{n}\theta/\sigma$$

and  $T(\theta)$  has an exact  $N(0, 1)$  distribution. This is because  $T(\theta)$  is centered at the true mean  $\theta$ , while the test statistic  $T(0)$  is centered at the (false) hypothesized mean of 0.

The power of the test is

$$\Pr(T > c \mid \theta) = \Pr(Z + \sqrt{n}\theta/\sigma > c) = 1 - \Phi(c - \sqrt{n}\theta/\sigma).$$

This function is monotonically increasing in  $\mu$  and  $n$ , and decreasing in  $\sigma$  and  $c$ .

Notice that for any  $c$  and  $\theta \neq 0$ , the power increases to 1 as  $n \rightarrow \infty$ . This means that for  $\theta \in \mathbb{H}_1$ , the test will reject  $\mathbb{H}_0$  with probability approaching 1 as the sample size gets large. We call this property **test consistency**.

**Definition 9.21.1** A test of  $\mathbb{H}_0 : \theta \in \Theta_0$  is **consistent against fixed alternatives** if for all  $\theta \in \Theta_1$ ,  $\Pr(\text{Reject } \mathbb{H}_0 \mid \theta) \rightarrow 1$  as  $n \rightarrow \infty$ .

For tests of the form “Reject  $\mathbb{H}_0$  if  $T > c$ ”, a sufficient condition for test consistency is that the  $T$  diverges to positive infinity with probability one for all  $\theta \in \Theta_1$ .

**Definition 9.21.2**  $T \xrightarrow{p} \infty$  as  $n \rightarrow \infty$  if for all  $M < \infty$ ,  $\Pr(T \leq M) \rightarrow 0$  as  $n \rightarrow \infty$ . Similarly,  $T \xrightarrow{p} -\infty$  as  $n \rightarrow \infty$  if for all  $M < \infty$ ,  $\Pr(T \geq -M) \rightarrow 0$  as  $n \rightarrow \infty$ .

In general, t-tests and Wald tests are consistent against fixed alternatives. Take a t-statistic for a test of  $\mathbb{H}_0 : \theta = \theta_0$

$$T = \frac{\hat{\theta} - \theta_0}{s(\hat{\theta})}$$

where  $\theta_0$  is a known value and  $s(\hat{\theta}) = \sqrt{n^{-1}\hat{V}_\theta}$ . Note that

$$T = \frac{\hat{\theta} - \theta}{s(\hat{\theta})} + \frac{\sqrt{n}(\theta - \theta_0)}{\sqrt{\hat{V}_\theta}}.$$

The first term on the right-hand-side converges in distribution to  $N(0, 1)$ . The second term on the right-hand-side equals zero if  $\theta = \theta_0$ , converges in probability to  $+\infty$  if  $\theta > \theta_0$ , and converges in probability to  $-\infty$  if  $\theta < \theta_0$ . Thus the two-sided t-test is consistent against  $\mathbb{H}_1 : \theta \neq \theta_0$ , and one-sided t-tests are consistent against the alternatives for which they are designed.

**Theorem 9.21.1** Under Assumptions 7.1.2 and 7.10.1, for  $\theta = \mathbf{r}(\beta) \neq \theta_0$  and  $q = 1$ , then  $|T| \xrightarrow{p} \infty$ , so for any  $c < \infty$  the test “Reject  $\mathbb{H}_0$  if  $|T| > c$ ” is consistent against fixed alternatives.

The Wald statistic for  $\mathbb{H}_0 : \boldsymbol{\theta} = \mathbf{r}(\boldsymbol{\beta}) = \boldsymbol{\theta}_0$  against  $\mathbb{H}_1 : \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$  is

$$W = n \left( \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \right)' \hat{\mathbf{V}}_{\boldsymbol{\theta}}^{-1} \left( \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \right).$$

Under  $\mathbb{H}_1$ ,  $\hat{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$ . Thus  $\left( \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \right)' \hat{\mathbf{V}}_{\boldsymbol{\theta}}^{-1} \left( \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \right) \xrightarrow{p} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)' \mathbf{V}_{\boldsymbol{\theta}}^{-1} (\boldsymbol{\theta} - \boldsymbol{\theta}_0) > 0$ . Hence under  $\mathbb{H}_1$ ,  $W \xrightarrow{p} \infty$ . Again, this implies that Wald tests are consistent tests.

**Theorem 9.21.2** *Under Assumptions 7.1.2 and 7.10.1, for  $\boldsymbol{\theta} = \mathbf{r}(\boldsymbol{\beta}) \neq \boldsymbol{\theta}_0$ , then  $W \xrightarrow{p} \infty$ , so for any  $c < \infty$  the test “Reject  $\mathbb{H}_0$  if  $W > c$ ” is consistent against fixed alternatives.*

## 9.22 Asymptotic Local Power

Consistency is a good property for a test, but does not give a useful approximation to the power of a test. To approximate the power function we need a distributional approximation.

The standard asymptotic method for power analysis uses what are called **local alternatives**. This is similar to our analysis of restriction estimation under misspecification (Section 8.13). The technique is to index the parameter by sample size so that the asymptotic distribution of the statistic is continuous in a localizing parameter. In this section we consider t-tests on real-valued parameters and in the next section consider Wald tests. Specifically, we consider parameter vectors  $\boldsymbol{\beta}_n$  which are indexed by sample size  $n$  and satisfy the real-valued relationship

$$\theta_n = r(\boldsymbol{\beta}_n) = \theta_0 + n^{-1/2}h \quad (9.25)$$

where the scalar  $h$  is called a **localizing parameter**. We index  $\boldsymbol{\beta}_n$  and  $\theta_n$  by sample size to indicate their dependence on  $n$ . The way to think of (9.25) is that the true value of the parameters are  $\boldsymbol{\beta}_n$  and  $\theta_n$ . The parameter  $\theta_n$  is close to the hypothesized value  $\theta_0$ , with deviation  $n^{-1/2}h$ .

The specification (9.25) states that for any fixed  $h$ ,  $\theta_n$  approaches  $\theta_0$  as  $n$  gets large. Thus  $\theta_n$  is “close” or “local” to  $\theta_0$ . The concept of a localizing sequence (9.25) might seem odd since in the actual world the sample size cannot mechanically affect the value of the parameter. Thus (9.25) should not be interpreted literally. Instead, it should be interpreted as a technical device which allows the asymptotic distribution of the test statistic to be continuous in the alternative hypothesis.

To evaluate the asymptotic distribution of the test statistic we start by examining the scaled estimate centered at the hypothesized value  $\theta_0$ . Breaking it into a term centered at the true value  $\theta_n$  and a remainder we find

$$\begin{aligned} \sqrt{n} \left( \hat{\theta} - \theta_0 \right) &= \sqrt{n} \left( \hat{\theta} - \theta_n \right) + \sqrt{n} (\theta_n - \theta_0) \\ &= \sqrt{n} \left( \hat{\theta} - \theta_n \right) + h \end{aligned}$$

where the second equality is (9.25). The first term is asymptotically normal:

$$\sqrt{n} \left( \hat{\theta} - \theta_n \right) \xrightarrow{d} \sqrt{V_{\theta}} Z.$$

where  $Z \sim N(0, 1)$ . Therefore

$$\sqrt{n} \left( \hat{\theta} - \theta_0 \right) \xrightarrow{d} \sqrt{V_{\theta}} Z + h$$

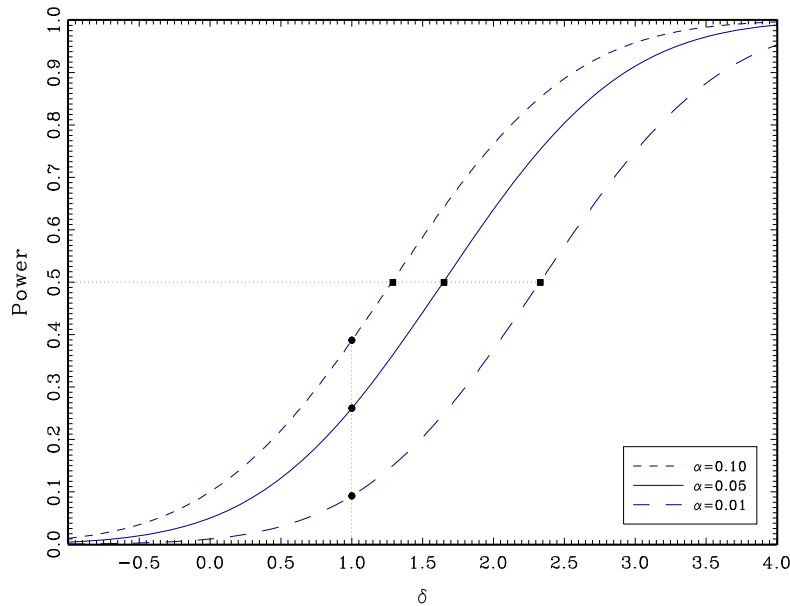


Figure 9.2: Asymptotic Local Power Function of One-Sided t Test

or  $N(h, V_\theta)$ . This is a continuous asymptotic distribution, and depends continuously on the localizing parameter  $h$ .

Applied to the t statistic we find

$$\begin{aligned}
 T &= \frac{\hat{\theta} - \theta_0}{s(\hat{\theta})} \\
 &\xrightarrow{d} \frac{\sqrt{V_\theta}Z + h}{\sqrt{V_\theta}} \\
 &\sim Z + \delta
 \end{aligned} \tag{9.26}$$

where  $\delta = h/\sqrt{V_\theta}$ . This generalizes Theorem 9.4.1 (which assumes  $\mathbb{H}_0$  is true) to allow for local alternatives of the form (9.25).

Consider a t-test of  $\mathbb{H}_0$  against the one-sided alternative  $\mathbb{H}_1 : \theta > \theta_0$  which rejects  $\mathbb{H}_0$  for  $T > c$  where  $\Phi(c) = 1 - \alpha$ . The **asymptotic local power** of this test is the limit (as the sample size diverges) of the rejection probability under the local alternative (9.25)

$$\begin{aligned}
 \lim_{n \rightarrow \infty} \Pr(\text{Reject } \mathbb{H}_0) &= \lim_{n \rightarrow \infty} \Pr(T > c) \\
 &= \Pr(Z + \delta > c) \\
 &= 1 - \Phi(c - \delta) \\
 &= \Phi(\delta - c) \\
 &\stackrel{\text{def}}{=} \pi(\delta).
 \end{aligned}$$

We call  $\pi(\delta)$  the **asymptotic local power function**.

In Figure 9.2 we plot the local power function  $\pi(\delta)$  as a function of  $\delta \in [-1, 4]$  for tests of asymptotic size  $\alpha = 0.10$ ,  $\alpha = 0.05$ , and  $\alpha = 0.01$ .  $\delta = 0$  corresponds to the null hypothesis so  $\pi(\delta) = \alpha$ . The power functions are monotonically increasing in  $\delta$ . Note that the power is lower than  $\alpha$  for  $\delta < 0$  due to the one-sided nature of the test.

We can see that the three power functions are ranked by  $\alpha$  so that the test with  $\alpha = 0.10$  has higher power than the test with  $\alpha = 0.01$ . This is the inherent trade-off between size and power. Decreasing size induces a decrease in power, and conversely.

The coefficient  $\delta$  can be interpreted as the parameter deviation measured as a multiple of the standard error  $s(\hat{\theta})$ . To see this, recall that  $s(\hat{\theta}) = n^{-1/2}\sqrt{\hat{V}_\theta} \simeq n^{-1/2}\sqrt{V_\theta}$  and then note that

$$\delta = \frac{h}{\sqrt{V_\theta}} \simeq \frac{n^{-1/2}h}{s(\hat{\theta})} = \frac{\theta_n - \theta_0}{s(\hat{\theta})}.$$

Thus  $\delta$  approximately equals the deviation  $\theta_n - \theta_0$  expressed as multiples of the standard error  $s(\hat{\theta})$ . Thus as we examine Figure 9.2, we can interpret the power function at  $\delta = 1$  (e.g. 26% for a 5% size test) as the power when the parameter  $\theta_n$  is one standard error above the hypothesized value. For example, from Table 4.1 the standard error for the coefficient on “Married Female” is 0.010. Thus in this example,  $\delta = 1$  corresponds to  $\theta_n = 0.010$  or an 1.0% wage premium for married females. Our calculations show that the asymptotic power of a one-sided 5% test against this alternative is about 26%.

The difference between power functions can be measured either vertically or horizontally. For example, in Figure 9.2 there is a vertical dotted line at  $\delta = 1$ , showing that the asymptotic local power function  $\pi(\delta)$  equals 39% for  $\alpha = 0.10$ , equals 26% for  $\alpha = 0.05$  and equals 9% for  $\alpha = 0.01$ . This is the difference in power across tests of differing size, holding fixed the parameter in the alternative.

A horizontal comparison can also be illuminating. To illustrate, in Figure 9.2 there is a horizontal dotted line at 50% power. 50% power is a useful benchmark, as it is the point where the test has equal odds of rejection and acceptance. The dotted line crosses the three power curves at  $\delta = 1.29$  ( $\alpha = 0.10$ ),  $\delta = 1.65$  ( $\alpha = 0.05$ ), and  $\delta = 2.33$  ( $\alpha = 0.01$ ). This means that the parameter  $\theta$  must be at least 1.65 standard errors above the hypothesized value for a one-sided 5% test to have 50% (approximate) power.

The ratio of these values (e.g.  $1.65/1.29 = 1.28$  for the asymptotic 5% versus 10% tests) measures the relative parameter magnitude needed to achieve the same power. (Thus, for a 5% size test to achieve 50% power, the parameter must be 28% larger than for a 10% size test.) Even more interesting, the square of this ratio (e.g.  $(1.65/1.29)^2 = 1.64$ ) can be interpreted as the increase in sample size needed to achieve the same power under fixed parameters. That is, to achieve 50% power, a 5% size test needs 64% more observations than a 10% size test. This interpretation follows by the following informal argument. By definition and (9.25)  $\delta = h/\sqrt{V_\theta} = \sqrt{n}(\theta_n - \theta_0)/\sqrt{V_\theta}$ . Thus holding  $\theta$  and  $V_\theta$  fixed,  $\delta^2$  is proportional to  $n$ .

The analysis of a two-sided t test is similar. (9.26) implies that

$$T = \left| \frac{\hat{\theta} - \theta_0}{s(\hat{\theta})} \right| \xrightarrow{d} |Z + \delta|$$

and thus the local power of a two-sided t test is

$$\begin{aligned} \lim_{n \rightarrow \infty} \Pr(\text{Reject } \mathbb{H}_0) &= \lim_{n \rightarrow \infty} \Pr(T > c) \\ &= \Pr(|Z + \delta| > c) \\ &= \Phi(\delta - c) - \Phi(-\delta - c) \end{aligned}$$

which is monotonically increasing in  $|\delta|$ .

**Theorem 9.22.1** Under Assumptions 7.1.2 and 7.10.1, and  $\theta_n = r(\beta_n) = r_0 + n^{-1/2}h$ , then

$$T(\theta_0) = \frac{\hat{\theta} - \theta_0}{s(\hat{\theta})} \xrightarrow{d} Z + \delta$$

where  $Z \sim N(0, 1)$  and  $\delta = h/\sqrt{V_{\theta}}$ . For  $c$  such that  $\Phi(c) = 1 - \alpha$ ,

$$\Pr(T(\theta_0) > c) \longrightarrow \Phi(\delta - c).$$

Furthermore, for  $c$  such that  $\Phi(c) = 1 - \alpha/2$ ,

$$\Pr(|T(\theta_0)| > c) \longrightarrow \Phi(\delta - c) - \Phi(-\delta - c).$$

### 9.23 Asymptotic Local Power, Vector Case

In this section we extend the local power analysis of the previous section to the case of vector-valued alternatives. We generalize (9.25) to allow  $\theta_n$  to be vector-valued. The local parameterization takes the form

$$\theta_n = r(\beta_n) = \theta_0 + n^{-1/2}h \quad (9.27)$$

where  $h$  is  $q \times 1$ .

Under (9.27),

$$\begin{aligned} \sqrt{n}(\hat{\theta} - \theta_0) &= \sqrt{n}(\hat{\theta} - \theta_n) + h \\ &\xrightarrow{d} Z_h \sim N(h, V_{\theta}), \end{aligned}$$

a normal random vector with mean  $h$  and variance matrix  $V_{\theta}$ .

Applied to the Wald statistic we find

$$\begin{aligned} W &= n(\hat{\theta} - \theta_0)' \hat{V}_{\theta}^{-1} (\hat{\theta} - \theta_0) \\ &\xrightarrow{d} Z_h' V_{\theta}^{-1} Z_h \sim \chi_q^2(\lambda) \end{aligned} \quad (9.28)$$

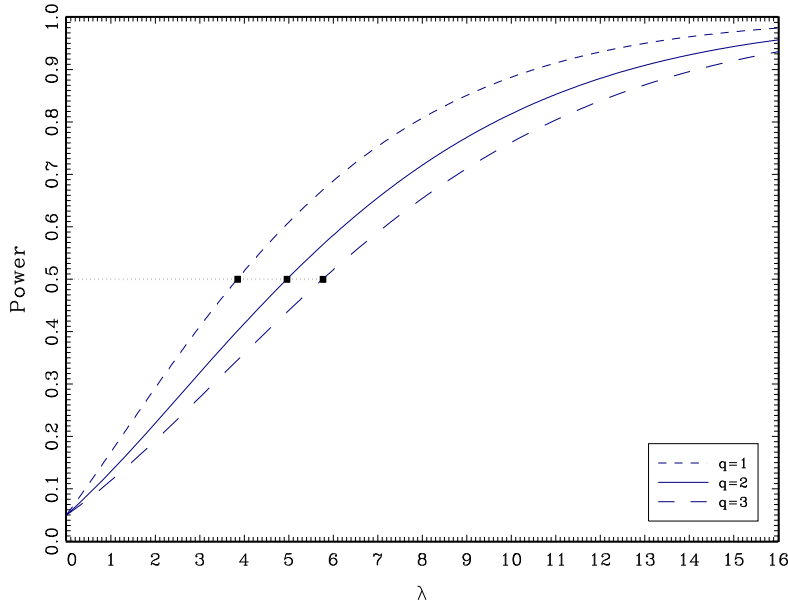
where  $\lambda = h' V_{\theta}^{-1} h$ .  $\chi_q^2(\lambda)$  is a non-central chi-square random variable with non-centrality parameter  $\lambda$ . (See Section 5.3 and Theorem 5.3.3.)

The convergence (9.28) shows that under the local alternatives (9.27),  $W \xrightarrow{d} \chi_q^2(\lambda)$ . This generalizes the null asymptotic distribution which obtains as the special case  $\lambda = 0$ . We can use this result to obtain a continuous asymptotic approximation to the power function. For any significance level  $\alpha > 0$  set the asymptotic critical value  $c$  so that  $\Pr(\chi_q^2 > c) = \alpha$ . Then as  $n \rightarrow \infty$ ,

$$\Pr(W > c) \longrightarrow \Pr(\chi_q^2(\lambda) > c) \stackrel{\text{def}}{=} \pi(\lambda).$$

The asymptotic local power function  $\pi(\lambda)$  depends only on  $\alpha$ ,  $q$ , and  $\lambda$ .



Figure 9.3: Asymptotic Local Power Function, Varying  $q$ 

**Theorem 9.23.1** Under Assumptions 7.1.2 and 7.10.1, and  $\theta_n = r(\beta_n) = \theta_0 + n^{-1/2}h$ , then

$$W \xrightarrow{d} \chi_q^2(\lambda)$$

where  $\lambda = h'V_\theta^{-1}h$ . Furthermore, for  $c$  such that  $\Pr(\chi_q^2 > c) = \alpha$ ,

$$\Pr(W > c) \longrightarrow \Pr(\chi_q^2(\lambda) > c).$$

Figure 9.3 plots  $\pi(\lambda)$  as a function of  $\lambda$  for  $q = 1$ ,  $q = 2$ , and  $q = 3$ , and  $\alpha = 0.05$ . The asymptotic power functions are monotonically increasing in  $\lambda$  and asymptote to one.

Figure 9.3 also shows the power loss for fixed non-centrality parameter  $\lambda$  as the dimensionality of the test increases. The power curves shift to the right as  $q$  increases, resulting in a decrease in power. This is illustrated by the dotted line at 50% power. The dotted line crosses the three power curves at  $\lambda = 3.85$  ( $q = 1$ ),  $\lambda = 4.96$  ( $q = 2$ ), and  $\lambda = 5.77$  ( $q = 3$ ). The ratio of these  $\lambda$  values correspond to the relative sample sizes needed to obtain the same power. Thus increasing the dimension of the test from  $q = 1$  to  $q = 2$  requires a 28% increase in sample size, or an increase from  $q = 1$  to  $q = 3$  requires a 50% increase in sample size, to obtain a test with 50% power.

## 9.24 Technical Proofs\*

**Proof of Theorem 9.12.1.** The conditions of Theorem 8.14.1 hold, since  $\mathbb{H}_0$  implies Assumption 8.6.1. From (8.58) with  $\widehat{W} = \widehat{V}_\beta$ , we see that

$$\begin{aligned} \sqrt{n}(\widehat{\beta} - \widetilde{\beta}_{\text{emd}}) &= \widehat{V}_\beta \widehat{R} (R_n^{*'} \widehat{V}_\beta \widehat{R})^{-1} R_n^{*'} \sqrt{n}(\widehat{\beta} - \beta) \\ &\xrightarrow{d} V_\beta R (R' V_\beta R)^{-1} R' N(0, V_\beta) \\ &= V_\beta R Z. \end{aligned}$$

where  $\mathbf{Z} \sim N(\mathbf{0}, (\mathbf{R}' \mathbf{V}_\beta \mathbf{R})^{-1})$ . Thus

$$\begin{aligned} J^* &= n \left( \hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}_{\text{emd}} \right)' \hat{\mathbf{V}}_\beta^{-1} \left( \hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}_{\text{emd}} \right) \\ &\xrightarrow{d} \mathbf{Z}' \mathbf{R}' \mathbf{V}_\beta \mathbf{V}_\beta^{-1} \mathbf{V}_\beta \mathbf{R} \mathbf{Z} \\ &= \mathbf{Z}' (\mathbf{R}' \mathbf{V}_\beta \mathbf{R}) \mathbf{Z} \\ &= \chi_q^2. \end{aligned}$$

■

## Exercises

**Exercise 9.1** Prove that if an additional regressor  $\mathbf{X}_{k+1}$  is added to  $\mathbf{X}$ , Theil's adjusted  $\bar{R}^2$  increases if and only if  $|T_{k+1}| > 1$ , where  $T_{k+1} = \hat{\beta}_{k+1}/s(\hat{\beta}_{k+1})$  is the t-ratio for  $\hat{\beta}_{k+1}$  and

$$s(\hat{\beta}_{k+1}) = (s^2[(\mathbf{X}'\mathbf{X})^{-1}]_{k+1,k+1})^{1/2}$$

is the homoskedasticity-formula standard error.

**Exercise 9.2** You have two independent samples  $(\mathbf{y}_1, \mathbf{X}_1)$  and  $(\mathbf{y}_2, \mathbf{X}_2)$  which satisfy  $\mathbf{y}_1 = \mathbf{X}_1\beta_1 + \mathbf{e}_1$  and  $\mathbf{y}_2 = \mathbf{X}_2\beta_2 + \mathbf{e}_2$ , where  $\mathbb{E}(\mathbf{x}_{1i}e_{1i}) = \mathbf{0}$  and  $\mathbb{E}(\mathbf{x}_{2i}e_{2i}) = \mathbf{0}$ , and both  $\mathbf{X}_1$  and  $\mathbf{X}_2$  have  $k$  columns. Let  $\hat{\beta}_1$  and  $\hat{\beta}_2$  be the OLS estimates of  $\beta_1$  and  $\beta_2$ . For simplicity, you may assume that both samples have the same number of observations  $n$ .

- (a) Find the asymptotic distribution of  $\sqrt{n} \left( (\hat{\beta}_2 - \hat{\beta}_1) - (\beta_2 - \beta_1) \right)$  as  $n \rightarrow \infty$ .
- (b) Find an appropriate test statistic for  $\mathbb{H}_0 : \beta_2 = \beta_1$ .
- (c) Find the asymptotic distribution of this statistic under  $\mathbb{H}_0$ .

**Exercise 9.3** Let  $T$  be a t-statistic for  $\mathbb{H}_0 : \theta = 0$  versus  $H_1 : \theta \neq 0$ . Since  $|T| \rightarrow_d |Z|$  under  $H_0$ , someone suggests the test “Reject  $\mathbb{H}_0$  if  $|T| < c_1$  or  $|T| > c_2$ , where  $c_1$  is the  $\alpha/2$  quantile of  $|Z|$  and  $c_2$  is the  $1 - \alpha/2$  quantile of  $|Z|$ .”

- (a) Show that the asymptotic size of the test is  $\alpha$ .
- (b) Is this a good test of  $\mathbb{H}_0$  versus  $\mathbb{H}_1$ ? Why or why not?

**Exercise 9.4** Let  $W$  be a Wald statistic for  $\mathbb{H}_0 : \boldsymbol{\theta} = \mathbf{0}$  versus  $\mathbb{H}_1 : \boldsymbol{\theta} \neq \mathbf{0}$ , where  $\boldsymbol{\theta}$  is  $q \times 1$ . Since  $W \rightarrow_d \chi_q^2$  under  $H_0$ , someone suggests the test “Reject  $\mathbb{H}_0$  if  $W < c_1$  or  $W > c_2$ , where  $c_1$  is the  $\alpha/2$  quantile of  $\chi_q^2$  and  $c_2$  is the  $1 - \alpha/2$  quantile of  $\chi_q^2$ .”

- (a) Show that the asymptotic size of the test is  $\alpha$ .
- (b) Is this a good test of  $\mathbb{H}_0$  versus  $\mathbb{H}_1$ ? Why or why not?

**Exercise 9.5** Take the linear model

$$\begin{aligned} y_i &= \mathbf{x}'_{1i}\beta_1 + \mathbf{x}'_{2i}\beta_2 + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= 0 \end{aligned}$$

where both  $\mathbf{x}_{1i}$  and  $\mathbf{x}_{2i}$  are  $q \times 1$ . Show how to test the hypotheses  $\mathbb{H}_0 : \beta_1 = \beta_2$  against  $\mathbb{H}_1 : \beta_1 \neq \beta_2$ .

**Exercise 9.6** Suppose a researcher wants to know which of a set of 20 regressors has an effect on a variable *testscore*. He regresses *testscore* on the 20 regressors and reports the results. One of the 20 regressors (*studytime*) has a large t-ratio (about 2.5), while other t-ratios are insignificant (smaller than 2 in absolute value). He argues that the data show that *studytime* is the key predictor for *testscore*. Do you agree with this conclusion? Is there a deficiency in his reasoning?

**Exercise 9.7** Take the model

$$\begin{aligned} y_i &= x_i\beta_1 + x_i^2\beta_2 + e_i \\ \mathbb{E}(e_i | x_i) &= 0 \end{aligned}$$

where  $y_i$  is wages (dollars per hour) and  $x_i$  is age. Describe how you would test the hypothesis that the expected wage for a 40-year-old worker is \$20 an hour.

**Exercise 9.8** You want to test  $\mathbb{H}_0 : \beta_2 = 0$  against  $\mathbb{H}_1 : \beta_2 \neq 0$  in the model

$$y_i = \mathbf{x}'_{1i}\beta_1 + \mathbf{x}'_{2i}\beta_2 + e_i$$

$$\mathbb{E}(\mathbf{x}_i e_i) = 0$$

You read a paper which estimates model

$$y_i = \mathbf{x}'_{1i}\hat{\gamma}_1 + (\mathbf{x}_{2i} - \mathbf{x}_{1i})'\hat{\gamma}_2 + \hat{e}_i$$

and reports a test of  $\mathbb{H}_0 : \gamma_2 = 0$  against  $\mathbb{H}_1 : \gamma_2 \neq 0$ . Is this related to the test you wanted to conduct?

**Exercise 9.9** Suppose a researcher uses one dataset to test a specific hypothesis  $\mathbb{H}_0$  against  $\mathbb{H}_1$ , and finds that he can reject  $\mathbb{H}_0$ . A second researcher gathers a similar but independent dataset, uses similar methods and finds that she cannot reject  $\mathbb{H}_0$ . How should we (as interested professionals) interpret these mixed results?

**Exercise 9.10** In Exercise 7.8, you showed that  $\sqrt{n}(\hat{\sigma}^2 - \sigma^2) \rightarrow_d N(0, V)$  as  $n \rightarrow \infty$  for some  $V$ . Let  $\hat{V}$  be an estimate of  $V$ .

- (a) Using this result, construct a t-statistic for  $\mathbb{H}_0 : \sigma^2 = 1$  against  $\mathbb{H}_1 : \sigma^2 \neq 1$ .
- (b) Using the Delta Method, find the asymptotic distribution of  $\sqrt{n}(\hat{\sigma} - \sigma)$ .
- (c) Use the previous result to construct a t-statistic for  $\mathbb{H}_0 : \sigma = 1$  against  $\mathbb{H}_1 : \sigma \neq 1$ .
- (d) Are the null hypotheses in (a) and (c) the same or are they different? Are the tests in (a) and (c) the same or are they different? If they are different, describe a context in which the two tests would give contradictory results.

**Exercise 9.11** Consider a regression such as Table 4.1 where both *experience* and its square are included. A researcher wants to test the hypothesis that *experience* does not affect mean wages, and does this by computing the t-statistic for *experience*. Is this the correct approach? If not, what is the appropriate testing method?

**Exercise 9.12** A researcher estimates a regression and computes a test of  $\mathbb{H}_0$  against  $\mathbb{H}_1$  and finds a p-value of  $p = 0.08$ , or “not significant”. She says “I need more data. If I had a larger sample the test will have more power and then the test will reject.” Is this interpretation correct?

**Exercise 9.13** A common view is that “If the sample size is large enough, any hypothesis will be rejected.” What does this mean? Interpret and comment.

**Exercise 9.14** Take the model

$$y_i = \mathbf{x}'_i \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(\mathbf{x}_i e_i) = 0$$

with parameter of interest  $\theta = \mathbf{R}'\boldsymbol{\beta}$  with  $\mathbf{R}$   $k \times 1$ . Let  $\hat{\boldsymbol{\beta}}$  be the least-squares estimate and  $\hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}$  its variance estimate.

- (a) Write down  $\hat{C}$ , the 95% asymptotic confidence interval for  $\theta$ , in terms of  $\hat{\boldsymbol{\beta}}$ ,  $\hat{\mathbf{V}}_{\hat{\boldsymbol{\beta}}}$ ,  $\mathbf{R}$ , and  $z = 1.96$  (the 97.5% quantile of  $N(0, 1)$ ).
- (b) Show that the decision “Reject  $\mathbb{H}_0$  if  $\theta_0 \notin \hat{C}$ ” is an asymptotic 5% test of  $\mathbb{H}_0 : \theta = \theta_0$ .

**Exercise 9.15** You are at a seminar where a colleague presents a simulation study of a test of a hypothesis  $\mathbb{H}_0$  with nominal size 5%. Based on  $B = 100$  simulation replications under  $\mathbb{H}_0$  the estimated size is 7%. Your colleague says: “Unfortunately the test over-rejects.”

- Do you agree or disagree with your colleague? Explain. Hint: Use an asymptotic (large  $B$ ) approximation.
- Suppose the number of simulation replications were  $B = 1000$  yet the estimated size is still 7%. Does your answer change?

**Exercise 9.16** You have  $n$  iid observations  $(y_i, x_{1i}, x_{2i})$ , and consider two alternative regression models

$$\begin{aligned} y_i &= \mathbf{x}'_{1i}\boldsymbol{\beta}_1 + e_{1i} \\ \mathbb{E}(\mathbf{x}_{1i}e_{1i}) &= 0 \end{aligned} \tag{9.29}$$

$$\begin{aligned} y_i &= \mathbf{x}'_{2i}\boldsymbol{\beta}_2 + e_{2i} \\ \mathbb{E}(\mathbf{x}_{2i}e_{2i}) &= 0 \end{aligned} \tag{9.30}$$

where  $\mathbf{x}_{1i}$  and  $\mathbf{x}_{2i}$  have at least some different regressors. (For example, (9.29) is a wage regression on geographic variables and (2) is a wage regression on personal appearance measurements.) You want to know if model (9.29) or model (9.30) fits the data better. Define  $\sigma_1^2 = E(e_{1i}^2)$  and  $\sigma_2^2 = E(e_{2i}^2)$ . You decide that the model with the smaller variance fit (e.g., model (9.29) fits better if  $\sigma_1^2 < \sigma_2^2$ .) You decide to test for this by testing the hypothesis of equal fit  $\mathbb{H}_0 : \sigma_1^2 = \sigma_2^2$  against the alternative of unequal fit  $\mathbb{H}_1 : \sigma_1^2 \neq \sigma_2^2$ . For simplicity, suppose that  $e_{1i}$  and  $e_{2i}$  are observed.

- Construct an estimate  $\hat{\theta}$  of  $\theta = \sigma_1^2 - \sigma_2^2$ .
- Find the asymptotic distribution of  $\sqrt{n}(\hat{\theta} - \theta)$  as  $n \rightarrow \infty$ .
- Find an estimator of the asymptotic variance of  $\hat{\theta}$ .
- Propose a test of asymptotic size  $\alpha$  of  $\mathbb{H}_0$  against  $\mathbb{H}_1$ .
- Suppose the test accepts  $\mathbb{H}_0$ . Briefly, what is your interpretation?

**Exercise 9.17** You have two regressors  $x_1$  and  $x_2$ , and estimate a regression with all quadratic terms

$$y_i = \alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{1i}^2 + \beta_4 x_{2i}^2 + \beta_5 x_{1i} x_{2i} + e_i$$

One of your advisors asks: Can we exclude the variable  $x_{2i}$  from this regression?

How do you translate this question into a statistical test? When answering these questions, be specific, not general.

- What is the relevant null and alternative hypotheses?
- What is an appropriate test statistic? Be specific.
- What is the appropriate asymptotic distribution for the statistic? Be specific.
- What is the rule for acceptance/rejection of the null hypothesis?

**Exercise 9.18** The observed data is  $\{y_i, \mathbf{x}_i, \mathbf{z}_i\} \in \mathbb{R} \times \mathbb{R}^k \times \mathbb{R}^\ell$ ,  $k > 1$  and  $\ell > 1$ ,  $i = 1, \dots, n$ . An econometrician first estimates

$$y_i = \mathbf{x}_i' \hat{\boldsymbol{\beta}} + \hat{e}_i$$

by least squares. The econometrician next regresses the residual  $\hat{e}_i$  on  $\mathbf{z}_i$ , which can be written as

$$\hat{e}_i = \mathbf{z}_i' \tilde{\boldsymbol{\gamma}} + \tilde{u}_i.$$

- Define the population parameter  $\boldsymbol{\gamma}$  being estimated in this second regression.
- Find the probability limit for  $\tilde{\boldsymbol{\gamma}}$ .
- Suppose the econometrician constructs a Wald statistic  $W_n$  for  $\mathbb{H}_0 : \boldsymbol{\gamma} = \mathbf{0}$  from the second regression, ignoring the regression. Write down the formula for  $W_n$ .
- Assuming  $\mathbb{E}(\mathbf{z}_i \mathbf{x}_i') = \mathbf{0}$ , find the asymptotic distribution for  $W_n$  under  $\mathbb{H}_0 : \boldsymbol{\gamma} = \mathbf{0}$ .
- If  $\mathbb{E}(\mathbf{z}_i \mathbf{x}_i') \neq \mathbf{0}$  will your answer to (d) change?

**Exercise 9.19** An economist estimates  $y_i = x_{1i}\beta_1 + x_{2i}\beta_2 + e_i$  by least-squares and tests the hypothesis  $\mathbb{H}_0 : \beta_2 = 0$  against  $\mathbb{H}_1 : \beta_2 \neq 0$ . She obtains a Wald statistic  $W_n = 0.34$ . The sample size is  $n = 500$ .

- What is the correct degrees of freedom for the  $\chi^2$  distribution to evaluate the significance of the Wald statistic?
- The Wald statistic  $W_n$  is very small. Indeed, is it less than the 1% quantile of the appropriate  $\chi^2$  distribution? If so, should you reject  $\mathbb{H}_0$ ? Explain your reasoning.

**Exercise 9.20** You are reading a paper, and it reports the results from two nested OLS regressions:

$$\begin{aligned} y_i &= \mathbf{x}_{1i}' \tilde{\boldsymbol{\beta}}_1 + \tilde{e}_i \\ y_i &= \mathbf{x}_{1i}' \hat{\boldsymbol{\beta}}_1 + \mathbf{x}_{2i}' \hat{\boldsymbol{\beta}}_2 + \hat{e}_i \end{aligned}$$

Some summary statistics are reported:

Short Regression	Long Regression
$R^2 = .20$	$R^2 = .26$
$\sum_{i=1}^n \tilde{e}_i^2 = 106$	$\sum_{i=1}^n \hat{e}_i^2 = 100$
# of coefficients=5	# of coefficients=8
$n = 50$	$n = 50$

You are curious if the estimate  $\hat{\boldsymbol{\beta}}_2$  is statistically different from the zero vector. Is there a way to determine an answer from this information? Do you have to make any assumptions (beyond the standard regularity conditions) to justify your answer?

**Exercise 9.21** Take the model

$$\begin{aligned} y_i &= x_{1i}\beta_1 + x_{2i}\beta_2 + x_{3i}\beta_3 + x_{4i}\beta_4 + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= \mathbf{0} \end{aligned}$$

Describe how you would test

$$\mathbb{H}_0 : \frac{\beta_1}{\beta_2} = \frac{\beta_3}{\beta_4}$$

against

$$\mathbb{H}_1 : \frac{\beta_1}{\beta_2} \neq \frac{\beta_3}{\beta_4}.$$

**Exercise 9.22** You have a random sample from the model

$$y_i = x_i\beta_1 + x_i^2\beta_2 + e_i$$

$$\mathbb{E}(e_i | x_i) = 0$$

where  $y_i$  is wages (dollars per hour) and  $x_i$  is age. Describe how you would test the hypothesis that the expected wage for a 40-year-old worker is \$20 an hour.

**Exercise 9.23** Let  $T_n$  be a test statistic such that under  $\mathbb{H}_0$ ,  $T_n \rightarrow_d \chi_3^2$ . Since  $P(\chi_3^2 > 7.815) = 0.05$ , an asymptotic 5% test of  $\mathbb{H}_0$  rejects when  $T_n > 7.815$ . An econometrician is interested in the Type I error of this test when  $n = 100$  and the data structure is well specified. She performs the following Monte Carlo experiment.

- $B = 200$  samples of size  $n = 100$  are generated from a distribution satisfying  $\mathbb{H}_0$ .
- On each sample, the test statistic  $T_{nb}$  is calculated.
- She calculates  $\hat{p} = \frac{1}{B} \sum_{b=1}^B 1(T_{nb} > 7.815) = 0.070$
- The econometrician concludes that the test  $T_n$  is oversized in this context – it rejects too frequently under  $\mathbb{H}_0$ .

Is her conclusion correct, incorrect, or incomplete? Be specific in your answer.

**Exercise 9.24** Do a Monte Carlo simulation. Take the model

$$y_i = \alpha + x_i\beta + e_i$$

$$\mathbb{E}(x_ie_i) = 0$$

where the parameter of interest is  $\theta = \exp(\beta)$ . Your data generating process (DGP) for the simulation is:  $x_i$  is  $U[0, 1]$ ,  $e_i$  is independent of  $x_i$  and  $N(0, 1)$ ,  $n = 50$ . Set  $\alpha = 0$  and  $\beta = 1$ . Generate  $B = 1000$  independent samples with  $\alpha$ . On each, estimate the regression by least-squares, calculate the covariance matrix using a standard (heteroskedasticity-robust) formula, and similarly estimate  $\theta$  and its standard error. For each replication, store  $\hat{\beta}$ ,  $\hat{\theta}$ ,  $t_\beta = (\hat{\beta} - \beta) / s(\hat{\beta})$ , and  $t_\theta = (\hat{\theta} - \theta) / s(\hat{\theta})$

- Does the value of  $\alpha$  matter? Explain why the described statistics are **invariant** to  $\alpha$  and thus setting  $\alpha = 0$  is irrelevant.
- From the 1000 replications estimate  $\mathbb{E}(\hat{\beta})$  and  $\mathbb{E}(\hat{\theta})$ . Discuss if you see evidence if either estimator is biased or unbiased.
- From the 1000 replications estimate  $\Pr(t_\beta > 1.645)$  and  $\Pr(t_\theta > 1.645)$ . What does asymptotic theory predict these probabilities should be in large samples? What do your simulation results indicate?

**Exercise 9.25** The data set `invest` on the textbook website contains data on 565 U.S. firms extracted from Compustat for the year 1987. (This is one year from a panel data set used by B. E. Hansen (1999). The original data was compiled by Hall and Hall (1993).) The variables are

- $I$  Investment to Capital Ratio (multiplied by 100).
- $Q$  Total Market Value to Asset Ratio (Tobin's Q).

- $C$       Cash Flow to Asset Ratio.
- $D$       Long Term Debt to Asset Ratio.

The flow variables are annual sums for 1987. The stock variables are beginning of year.

- (a) Estimate a linear regression of  $I_i$  on the other variables. Calculate appropriate standard errors.
- (b) Calculate asymptotic confidence intervals for the coefficients.
- (c) This regression is related to Tobin's  $q$  theory of investment, which suggests that investment should be predicted solely by  $Q_i$ . Thus the coefficient on  $Q_i$  should be positive and the others should be zero. Test the joint hypothesis that the coefficients on  $C_i$  and  $D_i$  are zero. Test the hypothesis that the coefficient on  $Q_i$  is zero. Are the results consistent with the predictions of the theory?
- (d) Now try a non-linear (quadratic) specification. Regress  $I_i$  on  $Q_i, C_i, D_i, Q_i^2, C_i^2, D_i^2, Q_i C_i, Q_i D_i, C_i D_i$ . Test the joint hypothesis that the six interaction and quadratic coefficients are zero.

**Exercise 9.26** In a paper in 1963, Marc Nerlove analyzed a cost function for 145 American electric companies. His data set `Nerlove1963` is on the textbook website. The variables are

- $C$       Total cost
- $Q$       Output
- $PL$       Unit price of labor
- $PK$       Unit price of capital
- $PF$       Unit price of labor

Nerlov was interested in estimating a *cost function*:  $C = f(Q, PL, PF, PK)$ .

- (a) First estimate an unrestricted Cobb-Douglass specification

$$\log C_i = \beta_1 + \beta_2 \log Q_i + \beta_3 \log PL_i + \beta_4 \log PK_i + \beta_5 \log PF_i + e_i. \quad (9.31)$$

Report parameter estimates and standard errors.

- (b) What is the economic meaning of the restriction  $\mathbb{H}_0 : \beta_3 + \beta_4 + \beta_5 = 1$ ?
- (c) Estimate (9.31) by constrained least-squares imposing  $\beta_3 + \beta_4 + \beta_5 = 1$ . Report your parameter estimates and standard errors.
- (d) Estimate (9.31) by efficient minimum distance imposing  $\beta_3 + \beta_4 + \beta_5 = 1$ . Report your parameter estimates and standard errors.
- (e) Test  $\mathbb{H}_0 : \beta_3 + \beta_4 + \beta_5 = 1$  using a Wald statistic.
- (f) Test  $\mathbb{H}_0 : \beta_3 + \beta_4 + \beta_5 = 1$  using a minimum distance statistic.

**Exercise 9.27** In Section 8.12 we report estimates from Mankiw, Romer and Weil (1992). We reported estimation both by unrestricted least-squares and by constrained estimation, imposing the constraint that three coefficients ( $2^{nd}$ ,  $3^{rd}$  and  $4^{th}$  coefficients) sum to zero, as implied by the Solow growth theory. Using the same dataset `MRW1992` estimate the unrestricted model and test the hypothesis that the three coefficients sum to zero.



**Exercise 9.28** Using the CPS dataset and the subsample of non-hispanic blacks (race code = 2), test the hypothesis that marriage status does not affect mean wages.

- (a) Take the regression reported in Table 4.1. Which variables will need to be omitted to estimate a regression for the subsample of blacks?
- (b) Express the hypothesis “marriage status does not affect mean wages” as a restriction on the coefficients. How many restrictions is this?
- (c) Find the Wald (or F) statistic for this hypothesis. What is the appropriate distribution for the test statistic? Calculate the p-value of the test.
- (d) What do you conclude?

**Exercise 9.29** Using the CPS dataset and the subsample of non-hispanic blacks (race code = 2) and whites (race code = 1), test the hypothesis that the returns to education is common across groups.

- (a) Allow the return to education to vary across the four groups (white male, white female, black male, black female) by interacting dummy variables with *education*. Estimate an appropriate version of the regression reported in Table 4.1.
- (b) Find the Wald (or F) statistic for this hypothesis. What is the appropriate distribution for the test statistic? Calculate the p-value of the test.
- (c) What do you conclude?

# Chapter 10

## Multivariate Regression

### 10.1 Introduction

**Multivariate regression** is a system of regression equations. Multivariate regression is used as reduced form models for instrumental variable estimation (explored in Chapter 11), vector autoregressions (explored in Chapter 15), demand systems (demand for multiple goods), and other contexts.

Multivariate regression is also called by the name **systems of regression equations**. Closely related is the method of **Seemingly Unrelated Regressions** (SUR) which we introduce in Section 10.7.

Most of the tools of single equation regression generalize naturally to multivariate regression. A major difference is a new set of notation to handle matrix estimates.

### 10.2 Regression Systems

A system of linear regressions takes the form

$$y_{ji} = \mathbf{x}'_{ji}\boldsymbol{\beta}_j + e_{ji} \quad (10.1)$$

for variables  $j = 1, \dots, m$  and observations  $i = 1, \dots, n$ , where the regressor vectors  $\mathbf{x}_{ji}$  are  $k_j \times 1$  and  $e_{ji}$  is an error. The coefficient vectors  $\boldsymbol{\beta}_j$  are  $k_j \times 1$ . The total number of coefficients are  $\bar{k} = \sum_{j=1}^n k_j$ . The regression system specializes to univariate regression when  $m = 1$ .

It is typical to treat the observations as independent across observations  $i$  but correlated across variables  $j$ . As an example, the observations  $y_{ji}$  could be expenditures by household  $i$  on good  $j$ . The standard assumptions are that households are mutually independent, but expenditures by an individual household are correlated across goods.

To describe the dependence between the dependent variables, we can define the  $m \times 1$  error vector  $\mathbf{e}_i = (e_{1i}, \dots, e_{mi})'$  and its  $m \times m$  variance matrix

$$\boldsymbol{\Sigma} = \mathbb{E}(\mathbf{e}_i \mathbf{e}_i').$$

The diagonal elements are the variances of the errors  $e_{ji}$ , and the off-diagonals are the covariances across variables. It is typical to allow  $\boldsymbol{\Sigma}$  to be unconstrained.

We can group the  $m$  equations (10.1) into a single equation as follows. Let  $\mathbf{y}_i = (y_{1i}, \dots, y_{mi})'$  be the  $m \times 1$  vector of dependent variables, define the  $\bar{k} \times m$  matrix of regressors

$$\mathbf{X}_i = \begin{pmatrix} \mathbf{x}_{1i} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \mathbf{x}_{2i} & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{x}_{mi} \end{pmatrix},$$

and define the  $\bar{k} \times 1$  stacked coefficient vector

$$\boldsymbol{\beta} = \begin{pmatrix} \boldsymbol{\beta}_1 \\ \vdots \\ \boldsymbol{\beta}_m \end{pmatrix}.$$

Then the  $m$  regression equations can jointly be written as

$$\mathbf{y}_i = \mathbf{X}_i' \boldsymbol{\beta} + \mathbf{e}_i. \quad (10.2)$$

The entire system can be written in matrix notation by stacking the variables. Define

$$\mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_n \end{pmatrix}, \quad \mathbf{e} = \begin{pmatrix} \mathbf{e}_1 \\ \vdots \\ \mathbf{e}_n \end{pmatrix}, \quad \bar{\mathbf{X}} = \begin{pmatrix} \mathbf{X}_1' \\ \vdots \\ \mathbf{X}_n' \end{pmatrix}$$

which are  $mn \times 1$ ,  $mn \times 1$ , and  $mn \times \bar{k}$ , respectively. The system can be written as

$$\mathbf{y} = \bar{\mathbf{X}} \boldsymbol{\beta} + \mathbf{e}.$$

In many (perhaps most) applications the regressor vectors  $\mathbf{x}_{ji}$  are common across the variables  $j$ , so  $\mathbf{x}_{ji} = \mathbf{x}_i$  and  $k_j = k$ . By this we mean that the same variables enter each equation with no exclusion restrictions. Several important simplifications occur in this context. One is that we can write (10.2) using the notation

$$\mathbf{y}_i = \mathbf{B}' \mathbf{x}_i + \mathbf{e}_i$$

where  $\mathbf{B} = (\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \dots, \boldsymbol{\beta}_m)$  is  $k \times m$ . Another is that we can write the system in the  $n \times m$  matrix notation

$$\mathbf{Y} = \mathbf{X} \mathbf{B} + \mathbf{E}$$

where

$$\mathbf{Y} = \begin{pmatrix} \mathbf{y}_1' \\ \vdots \\ \mathbf{y}_n' \end{pmatrix}, \quad \mathbf{E} = \begin{pmatrix} \mathbf{e}_1' \\ \vdots \\ \mathbf{e}_n' \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} \mathbf{x}_1' \\ \vdots \\ \mathbf{x}_n' \end{pmatrix}.$$

Another convenient implication of common regressors is that we have the simplification

$$\mathbf{X}_i = \begin{pmatrix} \mathbf{x}_i & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \mathbf{x}_i & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{x}_i \end{pmatrix} = \mathbf{I}_m \otimes \mathbf{x}_i$$

where  $\otimes$  is the Kronecker product (see Appendix A.16).

### 10.3 Least-Squares Estimator

Consider estimating each equation (10.1) by least-squares. This takes the form

$$\hat{\boldsymbol{\beta}}_j = \left( \sum_{i=1}^n \mathbf{x}_{ji} \mathbf{x}_{ji}' \right)^{-1} \left( \sum_{i=1}^n \mathbf{x}_{ji} y_{ji} \right).$$

The combined estimate of  $\boldsymbol{\beta}$  is the stacked vector

$$\hat{\boldsymbol{\beta}} = \begin{pmatrix} \hat{\boldsymbol{\beta}}_1 \\ \vdots \\ \hat{\boldsymbol{\beta}}_m \end{pmatrix}.$$

It turns that we can write this estimator using the systems notation

$$\hat{\beta} = (\overline{\mathbf{X}}' \overline{\mathbf{X}})^{-1} (\overline{\mathbf{X}}' \mathbf{y}) = \left( \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left( \sum_{i=1}^n \mathbf{X}_i \mathbf{y}_i \right). \quad (10.3)$$

To see this, observe that

$$\begin{aligned} \overline{\mathbf{X}}' \overline{\mathbf{X}} &= \begin{pmatrix} \mathbf{X}_1 & \cdots & \mathbf{X}_n \end{pmatrix} \begin{pmatrix} \mathbf{X}_1' \\ \vdots \\ \mathbf{X}_n' \end{pmatrix} \\ &= \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i' \\ &= \sum_{i=1}^n \begin{pmatrix} \mathbf{x}_{1i} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \mathbf{x}_{2i} & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{x}_{mi} \end{pmatrix} \begin{pmatrix} \mathbf{x}_{1i}' & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \mathbf{x}_{2i}' & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{x}_{mi}' \end{pmatrix} \\ &= \begin{pmatrix} \sum_{i=1}^n \mathbf{x}_{1i} \mathbf{x}_{1i}' & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \sum_{i=1}^n \mathbf{x}_{2i} \mathbf{x}_{2i}' & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \sum_{i=1}^n \mathbf{x}_{mi} \mathbf{x}_{mi}' \end{pmatrix} \end{aligned}$$

and

$$\begin{aligned} \overline{\mathbf{X}}' \mathbf{y} &= \begin{pmatrix} \mathbf{X}_1 & \cdots & \mathbf{X}_n \end{pmatrix} \begin{pmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_n \end{pmatrix} \\ &= \sum_{i=1}^n \mathbf{X}_i \mathbf{y}_i \\ &= \sum_{i=1}^n \begin{pmatrix} \mathbf{x}_{1i} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \mathbf{x}_{2i} & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{x}_{mi} \end{pmatrix} \begin{pmatrix} y_{1i} \\ \vdots \\ y_{mi} \end{pmatrix} \\ &= \begin{pmatrix} \sum_{i=1}^n \mathbf{x}_{1i} y_{1i} \\ \vdots \\ \sum_{i=1}^n \mathbf{x}_{mi} y_{mi} \end{pmatrix}. \end{aligned}$$

Hence

$$\begin{aligned} (\overline{\mathbf{X}}' \overline{\mathbf{X}})^{-1} (\overline{\mathbf{X}}' \mathbf{y}) &= \left( \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left( \sum_{i=1}^n \mathbf{X}_i \mathbf{y}_i \right) \\ &= \begin{pmatrix} (\sum_{i=1}^n \mathbf{x}_{1i} \mathbf{x}_{1i}')^{-1} (\sum_{i=1}^n \mathbf{x}_{1i} y_{1i}) \\ \vdots \\ (\sum_{i=1}^n \mathbf{x}_{mi} \mathbf{x}_{mi}')^{-1} (\sum_{i=1}^n \mathbf{x}_{mi} y_{mi}) \end{pmatrix} \\ &= \hat{\beta} \end{aligned}$$

as claimed.

The  $m \times 1$  residual vector for the  $i^{th}$  observation is

$$\hat{\mathbf{e}}_i = \mathbf{y}_i - \mathbf{X}_i' \hat{\beta}$$

and the least-squares estimate of the  $m \times m$  error variance matrix is

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{e}}_i \hat{\mathbf{e}}_i'. \quad (10.4)$$

In the case of common regressors, observe that

$$\hat{\beta}_j = \left( \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left( \sum_{i=1}^n \mathbf{x}_i y_{ji} \right).$$

We can set

$$\hat{\mathbf{B}} = (\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_m) = (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{Y}). \quad (10.5)$$

In Stata, multivariate regression can be implemented using the `mvreg` command.

## 10.4 Mean and Variance of Systems Least-Squares

We can calculate the finite-sample mean and variance of  $\hat{\beta}$  under the conditional mean assumption

$$\mathbb{E}(\mathbf{e}_i | \mathbf{x}_i) = \mathbf{0} \quad (10.6)$$

where  $\mathbf{x}_i$  is the union of the regressors  $\mathbf{x}_{ji}$ . Equation (10.6) is equivalent to  $\mathbb{E}(y_{ji} | \mathbf{x}_i) = \mathbf{x}_{ji}'\beta_j$ , or that the regression model is correctly specified.

We can center the estimator as

$$\hat{\beta} - \beta = (\overline{\mathbf{X}}'\overline{\mathbf{X}})^{-1} (\overline{\mathbf{X}}'\mathbf{e}) = \left( \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left( \sum_{i=1}^n \mathbf{X}_i \mathbf{e}_i \right).$$

Taking conditional expectations, we find  $\mathbb{E}(\hat{\beta} | \mathbf{X}) = \beta$ . Consequently, systems least-squares is unbiased under correct specification.

To compute the variance of the estimator, define the conditional covariance matrix of the errors of the  $i^{th}$  observation

$$\mathbb{E}(\mathbf{e}_i \mathbf{e}_i' | \mathbf{x}_i) = \Sigma_i$$

which in general is unrestricted. Observe that if the observations are mutually independent, then

$$\begin{aligned} \mathbb{E}(\mathbf{e} \mathbf{e}' | \mathbf{X}) &= \mathbb{E} \left( \begin{pmatrix} \mathbf{e}_1 \mathbf{e}_1 & \mathbf{e}_1 \mathbf{e}_2 & \cdots & \mathbf{e}_1 \mathbf{e}_n \\ \vdots & \ddots & & \vdots \\ \mathbf{e}_n \mathbf{e}_1 & \mathbf{e}_n \mathbf{e}_2 & \cdots & \mathbf{e}_n \mathbf{e}_n \end{pmatrix} \middle| \mathbf{X} \right) \\ &= \begin{pmatrix} \Sigma_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \ddots & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \Sigma_n \end{pmatrix}. \end{aligned}$$

Also, by independence across observations,

$$\text{var} \left( \sum_{i=1}^n \mathbf{X}_i \mathbf{e}_i \middle| \mathbf{X} \right) = \sum_{i=1}^n \text{var}(\mathbf{X}_i \mathbf{e}_i | \mathbf{x}_i) = \sum_{i=1}^n \mathbf{X}_i \Sigma_i \mathbf{X}_i'.$$

It follows that

$$\text{var}(\hat{\beta} | \mathbf{X}) = (\overline{\mathbf{X}}'\overline{\mathbf{X}})^{-1} \left( \sum_{i=1}^n \mathbf{X}_i \Sigma_i \mathbf{X}_i' \right) (\overline{\mathbf{X}}'\overline{\mathbf{X}})^{-1}.$$

When the regressors are common so that  $\mathbf{X}_i = \mathbf{I}_m \otimes \mathbf{x}_i$  then the covariance matrix can be written as

$$\text{var}(\hat{\boldsymbol{\beta}} | \mathbf{X}) = \left( \mathbf{I}_n \otimes (\mathbf{X}'\mathbf{X})^{-1} \right) \left( \sum_{i=1}^n (\boldsymbol{\Sigma}_i \otimes \mathbf{x}_i \mathbf{x}_i') \right) \left( \mathbf{I}_m \otimes (\mathbf{X}'\mathbf{X})^{-1} \right).$$

Alternatively, if the errors are conditionally homoskedastic

$$\mathbb{E}(\mathbf{e}_i \mathbf{e}_i' | \mathbf{x}_i) = \boldsymbol{\Sigma} \quad (10.7)$$

then the covariance matrix takes the form

$$\text{var}(\hat{\boldsymbol{\beta}} | \mathbf{X}) = (\overline{\mathbf{X}}'\overline{\mathbf{X}})^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \boldsymbol{\Sigma} \mathbf{x}_i' \right) (\overline{\mathbf{X}}'\overline{\mathbf{X}})^{-1}.$$

If both simplifications (common regressors and conditional homoskedasticity) hold then we have the considerable simplification

$$\text{var}(\hat{\boldsymbol{\beta}} | \mathbf{X}) = \boldsymbol{\Sigma} \otimes (\mathbf{X}'\mathbf{X})^{-1}.$$

## 10.5 Asymptotic Distribution

For an asymptotic distribution it is sufficient to consider the equation-by-equation projection model in which case

$$\mathbb{E}(\mathbf{x}_{ji} e_{ji}) = \mathbf{0}. \quad (10.8)$$

First, consider consistency. Since  $\hat{\boldsymbol{\beta}}_j$  are the standard least-squares estimators, they are consistent for the projection coefficients  $\boldsymbol{\beta}_j$ .

Second, consider the asymptotic distribution. Again by our single equation theory it is immediate that the  $\hat{\boldsymbol{\beta}}_j$  are asymptotically normally distributed. But our previous theory does not provide a joint distribution of the  $\hat{\boldsymbol{\beta}}_j$  across  $j$ . For this we need a joint theory for the stacked estimates  $\hat{\boldsymbol{\beta}}$ , which we now provide.

Since the vector

$$\mathbf{X}_i \mathbf{e}_i = \begin{pmatrix} \mathbf{x}_{1i} e_{1i} \\ \vdots \\ \mathbf{x}_{mi} e_{mi} \end{pmatrix}$$

is i.i.d. across  $i$  and mean zero under (10.8), the central limit theorem implies

$$\left( \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{X}_i \mathbf{e}_i \right) \xrightarrow{d} \mathbf{N}(\mathbf{0}, \boldsymbol{\Omega})$$

where

$$\boldsymbol{\Omega} = \mathbb{E}(\mathbf{X}_i \mathbf{e}_i \mathbf{e}_i' \mathbf{X}_i') = \mathbb{E}(\mathbf{X}_i \boldsymbol{\Sigma}_i \mathbf{X}_i').$$

The matrix  $\boldsymbol{\Omega}$  is the covariance matrix of the variables  $\mathbf{x}_{ji} e_{ji}$  across equations. Under conditional homoskedasticity (10.7) the matrix  $\boldsymbol{\Omega}$  simplifies to

$$\boldsymbol{\Omega} = \mathbb{E}(\mathbf{X}_i \boldsymbol{\Sigma} \mathbf{X}_i') \quad (10.9)$$

(see Exercise 10.1). When the regressors are common then it simplifies to

$$\boldsymbol{\Omega} = \mathbb{E}(\mathbf{e}_i \mathbf{e}_i' \otimes \mathbf{x}_i \mathbf{x}_i') \quad (10.10)$$

(see Exercise 10.2) and under both conditions (homoskedasticity and common regressors) it simplifies to

$$\boldsymbol{\Omega} = \boldsymbol{\Sigma} \otimes \mathbb{E}(\mathbf{x}_i \mathbf{x}_i') \quad (10.11)$$

(see Exercise 10.3).

Applied to the centered and normalized estimator we obtain the asymptotic distribution.

**Theorem 10.5.1** *Under Assumption 7.1.2,*

$$\sqrt{n} \left( \hat{\beta} - \beta \right) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_\beta)$$

where

$$\mathbf{V}_\beta = \mathbf{Q}^{-1} \mathbf{\Omega} \mathbf{Q}^{-1}$$

$$\mathbf{Q} = \mathbb{E}(\mathbf{X}_i \mathbf{X}_i') = \begin{pmatrix} \mathbb{E}(\mathbf{x}_{1i} \mathbf{x}_{1i}') & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \ddots & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbb{E}(\mathbf{x}_{ni} \mathbf{x}_{ni}') \end{pmatrix}$$

For a proof, see Exercise 10.4.

When the regressors are common then the matrix  $\mathbf{Q}$  simplifies as

$$\mathbf{Q} = \mathbf{I}_m \otimes \mathbb{E}(\mathbf{x}_i \mathbf{x}_i') \quad (10.12)$$

(See Exercise 10.5).

If both the regressors are common and the errors are conditionally homoskedastic (10.7) then we have the simplification

$$\mathbf{V}_\beta = \mathbf{\Sigma} \otimes (\mathbb{E}(\mathbf{x}_i \mathbf{x}_i'))^{-1} \quad (10.13)$$

(see Exercise 10.6).

Sometimes we are interested in parameters  $\theta = r(\beta_1, \dots, \beta_m) = r(\beta)$  which are functions of the coefficients from multiple equations. In this case the least-squares estimate of  $\theta$  is  $\hat{\theta} = r(\hat{\beta})$ . The asymptotic distribution of  $\hat{\theta}$  can be obtained from Theorem 10.5.1 by the delta method.

**Theorem 10.5.2** *Under Assumptions 7.1.2 and 7.10.1,*

$$\sqrt{n} \left( \hat{\theta} - \theta \right) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_\theta)$$

where

$$\mathbf{V}_\theta = \mathbf{R}' \mathbf{V}_\beta \mathbf{R}$$

$$\mathbf{R} = \frac{\partial}{\partial \beta} r(\beta)'$$

For a proof, see Exercise 10.7.

Theorem 10.5.2 is an example where multivariate regression is fundamentally distinct from univariate regression. Only by treating the least-squares estimates as a joint estimator can we obtain a distributional theory for an estimator  $\hat{\theta}$  which is a function of estimates from multiple equations and thereby construct standard errors, confidence intervals, and hypothesis tests.

## 10.6 Covariance Matrix Estimation

From the finite sample and asymptotic theory we can construct appropriate estimators for the variance of  $\hat{\beta}$ . In the general case we have

$$\hat{\mathbf{V}}_{\hat{\beta}} = \left( \overline{\mathbf{X}}' \overline{\mathbf{X}} \right)^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \hat{e}_i \hat{e}_i' \mathbf{x}_i' \right) \left( \overline{\mathbf{X}}' \overline{\mathbf{X}} \right)^{-1}.$$

Under conditional homoskedasticity (10.7) an appropriate estimator is

$$\hat{\mathbf{V}}_{\hat{\beta}}^0 = \left( \overline{\mathbf{X}}' \overline{\mathbf{X}} \right)^{-1} \left( \sum_{i=1}^n \mathbf{x}_i \hat{\Sigma} \mathbf{x}_i' \right) \left( \overline{\mathbf{X}}' \overline{\mathbf{X}} \right)^{-1}.$$

When the regressors are common then these estimators equal

$$\hat{\mathbf{V}}_{\hat{\beta}} = \left( \mathbf{I}_n \otimes (\mathbf{X}' \mathbf{X})^{-1} \right) \left( \sum_{i=1}^n (\hat{e}_i \hat{e}_i' \otimes \mathbf{x}_i \mathbf{x}_i') \right) \left( \mathbf{I}_n \otimes (\mathbf{X}' \mathbf{X})^{-1} \right)$$

and

$$\hat{\mathbf{V}}_{\hat{\beta}}^0 = \hat{\Sigma} \otimes (\mathbf{X}' \mathbf{X})^{-1},$$

respectively.

Covariance matrix estimators for  $\hat{\theta}$  are found as

$$\begin{aligned} \hat{\mathbf{V}}_{\hat{\theta}} &= \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\hat{\beta}} \hat{\mathbf{R}} \\ \hat{\mathbf{V}}_{\hat{\theta}}^0 &= \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\hat{\beta}}^0 \hat{\mathbf{R}} \\ \hat{\mathbf{R}} &= \frac{\partial}{\partial \beta} \mathbf{r}(\hat{\beta})'. \end{aligned}$$

**Theorem 10.6.1** Under Assumption 7.1.2,

$$n \hat{\mathbf{V}}_{\hat{\beta}} \xrightarrow{p} \mathbf{V}_{\beta}$$

and

$$n \hat{\mathbf{V}}_{\hat{\beta}}^0 \xrightarrow{p} \mathbf{V}_{\beta}^0$$

For a proof, see Exercise 10.8.

## 10.7 Seemingly Unrelated Regression

Consider the systems regression model under the conditional mean and conditional homoskedasticity assumptions

$$\begin{aligned} \mathbf{y}_i &= \mathbf{X}_i' \beta + \mathbf{e}_i \\ \mathbb{E}(\mathbf{e}_i \mid \mathbf{x}_i) &= \mathbf{0} \\ \mathbb{E}(\mathbf{e}_i \mathbf{e}_i' \mid \mathbf{x}_i) &= \Sigma \end{aligned} \tag{10.14}$$



Since the errors are correlated across equations we can consider estimation by Generalized Least Squares (GLS). To derive the estimator, premultiply (10.14) by  $\Sigma^{-1/2}$  so that the transformed error vector is i.i.d. with covariance matrix  $\mathbf{I}_m$ . Then apply least-squares and rearrange to find

$$\hat{\beta}_{\text{glS}} = \left( \sum_{i=1}^n \mathbf{X}_i \Sigma^{-1} \mathbf{X}_i' \right)^{-1} \left( \sum_{i=1}^n \mathbf{X}_i \Sigma^{-1} \mathbf{y}_i \right). \quad (10.15)$$

(see Exercise 10.9). Another approach is to take the vector representation

$$\mathbf{y} = \bar{\mathbf{X}}\boldsymbol{\beta} + \mathbf{e}$$

and calculate that the equation error  $\mathbf{e}$  has variance  $\mathbb{E}(\mathbf{e}\mathbf{e}') = \mathbf{I}_n \otimes \Sigma$ . Premultiply the equation by  $\mathbf{I}_n \otimes \Sigma^{-1/2}$  so that the transformed error has variance matrix  $\mathbf{I}_{nm}$  and then apply least-squares to find

$$\hat{\beta}_{\text{glS}} = \left( \bar{\mathbf{X}}' (\mathbf{I}_n \otimes \Sigma^{-1}) \bar{\mathbf{X}} \right)^{-1} \left( \bar{\mathbf{X}}' (\mathbf{I}_n \otimes \Sigma^{-1}) \mathbf{y} \right) \quad (10.16)$$

(see Exercise 10.10).

Expressions (10.15) and (10.16) are algebraically equivalent. To see the equivalence, observe that

$$\begin{aligned} \bar{\mathbf{X}}' (\mathbf{I}_n \otimes \Sigma^{-1}) \bar{\mathbf{X}} &= \begin{pmatrix} \mathbf{X}_1 & \cdots & \mathbf{X}_n \end{pmatrix} \begin{pmatrix} \Sigma^{-1} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \Sigma^{-1} & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \Sigma^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{X}_1' \\ \vdots \\ \mathbf{X}_n' \end{pmatrix} \\ &= \sum_{i=1}^n \mathbf{X}_i \Sigma^{-1} \mathbf{X}_i' \end{aligned}$$

and

$$\begin{aligned} \bar{\mathbf{X}}' (\mathbf{I}_n \otimes \Sigma^{-1}) \mathbf{y} &= \begin{pmatrix} \mathbf{X}_1 & \cdots & \mathbf{X}_n \end{pmatrix} \begin{pmatrix} \Sigma^{-1} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \Sigma^{-1} & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \Sigma^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_n \end{pmatrix} \\ &= \sum_{i=1}^n \mathbf{X}_i \Sigma^{-1} \mathbf{y}_i. \end{aligned}$$

Since  $\Sigma$  is unknown it must be replaced by an estimator. Using  $\hat{\Sigma}$  from (10.4) we obtain a feasible GLS estimator.

$$\begin{aligned} \hat{\beta}_{\text{sur}} &= \left( \sum_{i=1}^n \mathbf{X}_i \hat{\Sigma}^{-1} \mathbf{X}_i' \right)^{-1} \left( \sum_{i=1}^n \mathbf{X}_i \hat{\Sigma}^{-1} \mathbf{y}_i \right) \\ &= \left( \bar{\mathbf{X}}' (\mathbf{I}_n \otimes \hat{\Sigma}^{-1}) \bar{\mathbf{X}} \right)^{-1} \left( \bar{\mathbf{X}}' (\mathbf{I}_n \otimes \hat{\Sigma}^{-1}) \mathbf{y} \right). \end{aligned} \quad (10.17)$$

This is known as the **Seemingly Unrelated Regression (SUR)** estimator.

The estimator  $\hat{\Sigma}$  can be updated by calculating the SUR residuals  $\hat{\mathbf{e}}_i = \mathbf{y}_i - \mathbf{X}_i' \hat{\beta}_{\text{SUR}}$  and the covariance matrix estimate  $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{e}}_i \hat{\mathbf{e}}_i'$ . Substituted into (10.17) we find an iterated SUR estimator, and this can be iterated until convergence.

Under conditional homoskedasticity (10.7) we can derive its asymptotic distribution.

**Theorem 10.7.1** *Under Assumption 7.1.2 and (10.7)*

$$\sqrt{n} \left( \hat{\beta}_{\text{sur}} - \boldsymbol{\beta} \right) \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{V}_{\boldsymbol{\beta}}^*)$$

where

$$\mathbf{V}_{\boldsymbol{\beta}}^* = \left( \mathbb{E}(\mathbf{X}_i \Sigma^{-1} \mathbf{X}_i') \right)^{-1}.$$

For a proof, see Exercise 10.11.

Under these assumptions, SUR is more efficient than least-squares (in particular, under the assumption of conditional homoskedasticity).

**Theorem 10.7.2** *Under Assumption 7.1.2 and (10.7)*

$$\begin{aligned}\mathbf{V}_{\beta}^* &= (\mathbb{E}(\mathbf{X}_i \boldsymbol{\Sigma}^{-1} \mathbf{X}_i'))^{-1} \\ &\leq (\mathbb{E}(\mathbf{X}_i \mathbf{X}_i'))^{-1} \mathbb{E}(\mathbf{X}_i \boldsymbol{\Sigma} \mathbf{X}_i') (\mathbb{E}(\mathbf{X}_i \mathbf{X}_i'))^{-1} \\ &= \mathbf{V}_{\beta}\end{aligned}$$

and thus  $\hat{\beta}_{\text{sur}}$  is asymptotically more efficient than  $\hat{\beta}_{OLS}$ .

For a proof, see Exercise 10.12.

An appropriate estimator of the variance of  $\hat{\beta}_{SUR}$  is

$$\hat{\mathbf{V}}_{\hat{\beta}} = \left( \sum_{i=1}^n \mathbf{X}_i \hat{\boldsymbol{\Sigma}}^{-1} \mathbf{X}_i' \right)^{-1}.$$

**Theorem 10.7.3** *Under Assumption 7.1.2 and (10.7)*

$$n \hat{\mathbf{V}}_{\hat{\beta}} \xrightarrow{p} \mathbf{V}_{\beta}$$

and thus  $\hat{\beta}_{SUR}$  is asymptotically more efficient than  $\hat{\beta}_{OLS}$ .

For a proof, see Exercise 10.13.

In Stata, the seemingly unrelated regressions estimator is implemented using the `sureg` command.

### Arnold Zellner

Arnold Zellner (1927-2000) of the United States was a founding father of the econometrics field. He was a pioneer in Bayesian econometrics. One of his core contributions was the method of Seemingly Unrelated Regressions.

## 10.8 Maximum Likelihood Estimator

Take the linear model under the assumption that the error is independent of the regressors and multivariate normally distributed. Thus

$$\begin{aligned}\mathbf{y}_i &= \mathbf{X}_i' \boldsymbol{\beta} + \mathbf{e}_i \\ \mathbf{e}_i &\sim \mathbf{N}(\mathbf{0}, \boldsymbol{\Sigma}).\end{aligned}$$

In this case we can consider the maximum likelihood estimator (MLE) of the coefficients.

It is convenient to reparameterize the covariance matrix in terms of its inverse, thus  $\mathbf{S} = \boldsymbol{\Sigma}^{-1}$ . With this reparameterization, the conditional density of  $\mathbf{y}_i$  given  $\mathbf{X}_i$  equals

$$f(\mathbf{y}_i|\mathbf{X}_i) = \frac{\det(\mathbf{S})^{1/2}}{(2\pi)^{m/2}} \exp\left(-\frac{1}{2}(\mathbf{y}_i - \mathbf{X}_i'\boldsymbol{\beta})' \mathbf{S} (\mathbf{y}_i - \mathbf{X}_i'\boldsymbol{\beta})\right).$$

The log-likelihood function for the sample is

$$\log L(\boldsymbol{\beta}, \mathbf{S}) = -\frac{nm}{2} \log(2\pi) + \frac{n}{2} \log \det(\mathbf{S}) - \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i - \mathbf{X}_i'\boldsymbol{\beta})' \mathbf{S} (\mathbf{y}_i - \mathbf{X}_i'\boldsymbol{\beta}).$$

The maximum likelihood estimator  $(\hat{\boldsymbol{\beta}}, \hat{\mathbf{S}})$  maximizes the log-likelihood function. The first order conditions are

$$\begin{aligned} 0 &= \frac{\partial}{\partial \boldsymbol{\beta}} \log L(\boldsymbol{\beta}, \mathbf{S}) \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}, \mathbf{S}=\hat{\mathbf{S}}} \\ &= \sum_{i=1}^n \mathbf{X}_i \hat{\mathbf{S}} (\mathbf{y}_i - \mathbf{X}_i'\hat{\boldsymbol{\beta}}) \end{aligned}$$

and

$$\begin{aligned} 0 &= \frac{\partial}{\partial \mathbf{S}} \log L(\boldsymbol{\beta}, \mathbf{S}) \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}, \mathbf{S}=\hat{\mathbf{S}}} \\ &= \frac{n}{2} \hat{\mathbf{S}}^{-1} - \frac{1}{2} \text{tr} \left( \sum_{i=1}^n (\mathbf{y}_i - \mathbf{X}_i'\hat{\boldsymbol{\beta}}) (\mathbf{y}_i - \mathbf{X}_i'\hat{\boldsymbol{\beta}})' \right). \end{aligned}$$

The second equation uses the matrix results  $\frac{\partial}{\partial \mathbf{S}} \log \det(\mathbf{S}) = \mathbf{S}^{-1}$  and  $\frac{\partial}{\partial \mathbf{B}} \text{tr}(\mathbf{AB}) = \mathbf{A}'$  from Appendix A.15.

Solving and making the substitution  $\hat{\boldsymbol{\Sigma}} = \hat{\mathbf{S}}^{-1}$  we obtain

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \left( \sum_{i=1}^n \mathbf{X}_i \hat{\boldsymbol{\Sigma}}^{-1} \mathbf{X}_i' \right)^{-1} \left( \sum_{i=1}^n \mathbf{X}_i \hat{\boldsymbol{\Sigma}}^{-1} \mathbf{y}_i \right) \\ \hat{\boldsymbol{\Sigma}} &= \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i - \mathbf{X}_i'\hat{\boldsymbol{\beta}}) (\mathbf{y}_i - \mathbf{X}_i'\hat{\boldsymbol{\beta}})'. \end{aligned}$$

Notice that each equation refers to the other. Hence these are not closed-form expressions, but can be solved via iteration. The solution is identical to the iterated SUR estimator. Thus the SUR estimator (iterated) is identical to the MLE under normality.

Recall that the SUR estimator simplifies to OLS when the regressors are common across equations. The same occurs for the MLE. Thus when  $\mathbf{X}_i = \mathbf{I}_m \otimes \mathbf{x}_i$  we find that  $\hat{\boldsymbol{\beta}}_{MLE} = \hat{\boldsymbol{\beta}}_{OLS}$  and  $\hat{\boldsymbol{\Sigma}}_{MLE} = \hat{\boldsymbol{\Sigma}}_{OLS}$ .

## 10.9 Reduced Rank Regression

One context where systems estimation is important is when it is desired to impose or test restrictions across equations. Restricted systems are commonly estimated by maximum likelihood under normality. In this section we explore one important special case of restricted multivariate regression known as reduced rank regression. The model was originally proposed by Anderson (1951) and extended by Johansen (1995).

The unrestricted model is

$$\begin{aligned} \mathbf{y}_i &= \mathbf{B}'\mathbf{x}_i + \mathbf{C}'\mathbf{z}_i + e_i \\ \mathbb{E}(\mathbf{e}_i\mathbf{e}_i' | \mathbf{x}_i, \mathbf{z}_i) &= \boldsymbol{\Sigma} \end{aligned} \quad (10.18)$$

where  $\mathbf{B}$  is  $k \times m$ ,  $\mathbf{C}$  is  $\ell \times m$ , and  $\mathbf{x}_i$  and  $\mathbf{z}_i$  are regressors. We separate the regressors  $\mathbf{x}_i$  and  $\mathbf{z}_i$  because the coefficient matrix  $\mathbf{B}$  will be restricted while  $\mathbf{C}$  will be unrestricted.

The matrix  $\mathbf{B}$  is full rank if

$$\text{rank}(\mathbf{B}) = \min(k, m).$$

The reduced rank restriction is that

$$\text{rank}(\mathbf{B}) = r < \min(k, m)$$

for some known  $r$ .

The reduced rank restriction implies that we can write the coefficient matrix  $\mathbf{B}$  in the factored form

$$\mathbf{B} = \mathbf{G}\mathbf{A}' \quad (10.19)$$

where  $\mathbf{A}$  is  $m \times r$  and  $\mathbf{G}$  is  $k \times r$ . This representation is not unique (as we can replace  $\mathbf{G}$  with  $\mathbf{G}\mathbf{Q}$  and  $\mathbf{A}$  with  $\mathbf{A}\mathbf{Q}^{-1'}$  for any invertible  $\mathbf{Q}$  and the same relation holds). Identification therefore requires a normalization of the coefficients. A conventional normalization is

$$\mathbf{G}'\mathbf{D}\mathbf{G} = \mathbf{I}_r$$

for given  $\mathbf{D}$ .

Equivalently, the reduced rank restriction can be imposed by requiring that  $\mathbf{B}$  satisfy the restriction  $\mathbf{B}\mathbf{A}_\perp = \mathbf{G}\mathbf{A}'\mathbf{A}_\perp = \mathbf{0}$  for some  $m \times (m-r)$  coefficient matrix  $\mathbf{A}_\perp$ . Since  $\mathbf{G}$  is full rank this requires that  $\mathbf{A}'\mathbf{A}_\perp = \mathbf{0}$ , hence  $\mathbf{A}_\perp$  is the orthogonal complement to  $\mathbf{A}$ . Note that  $\mathbf{A}_\perp$  is not unique as it can be replaced by  $\mathbf{A}_\perp\mathbf{Q}$  for any  $(m-r) \times (m-r)$  invertible  $\mathbf{Q}$ . Thus if  $\mathbf{A}_\perp$  is to be estimated it requires a normalization.

We discuss methods for estimation of  $\mathbf{G}$ ,  $\mathbf{A}$ ,  $\boldsymbol{\Sigma}$ ,  $\mathbf{C}$ , and  $\mathbf{A}_\perp$ . The standard approach is maximum likelihood under the assumption that  $\mathbf{e}_i \sim N(\mathbf{0}, \boldsymbol{\Sigma})$ . The log-likelihood function for the sample is

$$\begin{aligned} \log L(\mathbf{G}, \mathbf{A}, \mathbf{C}, \boldsymbol{\Sigma}) &= -\frac{nm}{2} \log(2\pi) - \frac{n}{2} \log \det(\boldsymbol{\Sigma}) \\ &\quad - \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i - \mathbf{A}\mathbf{G}'\mathbf{x}_i - \mathbf{C}'\mathbf{z}_i)' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \mathbf{A}\mathbf{G}'\mathbf{x}_i - \mathbf{C}'\mathbf{z}_i). \end{aligned}$$

Anderson (1951) derived the MLE by imposing the constraint  $\mathbf{B}\mathbf{A}_\perp = \mathbf{0}$  via the method of Lagrange multipliers. This turns out to be algebraically cumbersome.

Johansen (1995) instead proposed a concentration method which turns out to be relatively straightforward. The method is as follows. First, treat  $\mathbf{G}$  as if it is known. Then maximize the log-likelihood with respect to the other parameters. Resubstituting these estimates, we obtain the concentrated log-likelihood function with respect to  $\mathbf{G}$ . This can be maximized to find the MLE for  $\mathbf{G}$ . The other parameter estimates are then obtain by substitution. We now describe these steps in detail.

Given  $\mathbf{G}$ , the likelihood is a normal multivariate regression in the variables  $\mathbf{G}'\mathbf{x}_i$  and  $\mathbf{z}_i$ , so the MLE for  $\mathbf{A}$ ,  $\mathbf{C}$  and  $\boldsymbol{\Sigma}$  are least-squares. In particular, using the Frisch-Waugh-Lovell residual regression formula, we can write the estimators for  $\mathbf{A}$  and  $\boldsymbol{\Sigma}$  as

$$\hat{\mathbf{A}}(\mathbf{G}) = \left( \tilde{\mathbf{Y}}' \tilde{\mathbf{X}} \mathbf{G} \right) \left( \mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \mathbf{G} \right)^{-1}$$

and

$$\hat{\boldsymbol{\Sigma}}(\mathbf{G}) = \frac{1}{n} \left( \tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} - \tilde{\mathbf{Y}}' \tilde{\mathbf{X}} \mathbf{G} \left( \mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \mathbf{G} \right)^{-1} \mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{Y}} \right)$$

where

$$\begin{aligned}\tilde{\mathbf{Y}} &= \mathbf{Y} - \mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{Y} \\ \tilde{\mathbf{X}} &= \mathbf{X} - \mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{X}.\end{aligned}$$

Substituting these estimators into the log-likelihood function, we obtain the concentrated likelihood function, which is a function of  $\mathbf{G}$  only

$$\begin{aligned}\log \tilde{L}(\mathbf{G}) &= \log L(\mathbf{G}, \hat{\mathbf{A}}(\mathbf{G}), \hat{\mathbf{C}}(\mathbf{G}), \hat{\mathbf{\Sigma}}(\mathbf{G})) \\ &= \frac{m}{2} (n \log(2\pi) - 1) - \frac{n}{2} \log \det \left( \tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} - \tilde{\mathbf{Y}}' \tilde{\mathbf{X}} \mathbf{G} (\mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \mathbf{G})^{-1} \mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{Y}} \right) \\ &= \frac{m}{2} (n \log(2\pi) - 1) - \frac{n}{2} \log \det \left( \tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} \right) \frac{\det \left( \mathbf{G}' \left( \tilde{\mathbf{X}}' \tilde{\mathbf{X}} - \tilde{\mathbf{X}}' \tilde{\mathbf{Y}} (\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}})^{-1} \mathbf{Y}' \tilde{\mathbf{X}} \right) \mathbf{G} \right)}{\det \left( \mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \mathbf{G} \right)}.\end{aligned}$$

The third equality uses Theorem A.7.1.8. The MLE  $\hat{\mathbf{G}}$  for  $\mathbf{G}$  is the maximizer of  $\log \tilde{L}(\mathbf{G})$ , or equivalently equals

$$\hat{\mathbf{G}} = \underset{\mathbf{G}}{\operatorname{argmin}} \frac{\det \left( \mathbf{G}' \left( \tilde{\mathbf{X}}' \tilde{\mathbf{X}} - \tilde{\mathbf{X}}' \tilde{\mathbf{Y}} (\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}})^{-1} \mathbf{Y}' \tilde{\mathbf{X}} \right) \mathbf{G} \right)}{\det \left( \mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \mathbf{G} \right)} \quad (10.20)$$

$$\begin{aligned}&= \underset{\mathbf{G}}{\operatorname{argmax}} \frac{\det \left( \mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{Y}} (\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}})^{-1} \mathbf{Y}' \tilde{\mathbf{X}} \mathbf{G} \right)}{\det \left( \mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \mathbf{G} \right)} \quad (10.21) \\ &= \{\mathbf{v}_1, \dots, \mathbf{v}_r\}\end{aligned}$$

which are the generalized eigenvectors of  $\tilde{\mathbf{X}}' \tilde{\mathbf{Y}} (\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}})^{-1} \mathbf{Y}' \tilde{\mathbf{X}}$  with respect to  $\tilde{\mathbf{X}}' \tilde{\mathbf{X}}$  corresponding to the  $r$  largest generalized eigenvalues. (Generalized eigenvalues and eigenvectors are discussed in Section A.10.) The estimator satisfies the normalization  $\hat{\mathbf{G}}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \hat{\mathbf{G}} = \mathbf{I}_r$ . Letting  $\mathbf{v}_j^*$  denote the eigenvectors of (10.20), we can also express  $\hat{\mathbf{G}} = \{\mathbf{v}_m^*, \dots, \mathbf{v}_{m-r+1}^*\}$ .

This is computationally straightforward. In MATLAB, for example, the generalized eigenvalues and eigenvectors of a matrix  $\mathbf{A}$  with respect to  $\mathbf{B}$  are found using the command `eig(A,B)`.

Given  $\hat{\mathbf{G}}$ , the MLE  $\hat{\mathbf{A}}$ ,  $\hat{\mathbf{C}}$ ,  $\hat{\mathbf{\Sigma}}$  are found by least-squares regression of  $\mathbf{y}_i$  on  $\hat{\mathbf{G}}' \mathbf{x}_i$  and  $\mathbf{z}_i$ . In particular,  $\hat{\mathbf{A}} = \hat{\mathbf{G}}' \tilde{\mathbf{X}}' \tilde{\mathbf{Y}}$  since  $\hat{\mathbf{G}}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \hat{\mathbf{G}} = \mathbf{I}_r$ .

We now discuss the estimator  $\hat{\mathbf{A}}_\perp$  of  $\mathbf{A}_\perp$ . It turns out that

$$\begin{aligned}\hat{\mathbf{A}}_\perp &= \underset{\mathbf{A}}{\operatorname{argmax}} \frac{\det \left( \mathbf{A}' \left( \tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} - \tilde{\mathbf{Y}}' \tilde{\mathbf{X}} (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}' \tilde{\mathbf{Y}} \right) \mathbf{A} \right)}{\det \left( \mathbf{A}' \tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} \mathbf{A} \right)} \quad (10.22) \\ &= \{\mathbf{w}_1, \dots, \mathbf{w}_{m-r}\}\end{aligned}$$

the eigenvectors of  $\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} - \tilde{\mathbf{Y}}' \tilde{\mathbf{X}} (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}' \tilde{\mathbf{Y}}$  with respect to  $\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}}$  associated with the largest  $m - r$  eigenvalues.

By the dual eigenvalue relation (Theorem A.10.1), the eigenvalue problems in equations (10.20) and (10.22) have the same non-unit eigenvalues  $\lambda_j$ , and the associated eigenvectors  $\mathbf{v}_j^*$  and  $\mathbf{w}_j$  satisfy

the relationship  $\mathbf{w}_j = \lambda_j^{-1/2} \left( \widetilde{\mathbf{Y}}' \widetilde{\mathbf{Y}} \right)^{-1} \widetilde{\mathbf{Y}}' \widetilde{\mathbf{X}} \mathbf{v}_j^*$ . Letting  $\mathbf{\Lambda} = \text{diag}\{\lambda_m, \dots, \lambda_{m-r+1}\}$  this implies

$$\begin{aligned} \{\mathbf{w}_m, \dots, \mathbf{w}_{m-r+1}\} &= \left( \widetilde{\mathbf{Y}}' \widetilde{\mathbf{Y}} \right)^{-1} \widetilde{\mathbf{Y}}' \widetilde{\mathbf{X}} \{\mathbf{v}_m^*, \dots, \mathbf{v}_{m-r+1}^*\} \mathbf{\Lambda} \\ &= \left( \widetilde{\mathbf{Y}}' \widetilde{\mathbf{Y}} \right)^{-1} \widehat{\mathbf{A}} \mathbf{\Lambda}. \end{aligned}$$

The second equality holds since  $\widehat{\mathbf{G}} = \{\mathbf{v}_m^*, \dots, \mathbf{v}_{m-r+1}^*\}$  and  $\widehat{\mathbf{A}} = \widetilde{\mathbf{Y}}' \widetilde{\mathbf{X}} \widehat{\mathbf{G}}$ . Since the eigenvectors  $\mathbf{w}_j$  satisfy the orthogonality property  $\mathbf{w}_j' \widetilde{\mathbf{Y}}' \widetilde{\mathbf{Y}} \mathbf{w}_\ell = 0$  for  $j \neq \ell$ , it follows that

$$0 = \widehat{\mathbf{A}}_\perp' \widetilde{\mathbf{Y}}' \widetilde{\mathbf{Y}} \{\mathbf{w}_m, \dots, \mathbf{w}_{m-r+1}\} = \widehat{\mathbf{A}}_\perp' \widehat{\mathbf{A}} \mathbf{\Lambda}.$$

Since  $\mathbf{\Lambda} > 0$  we conclude that  $\widehat{\mathbf{A}}_\perp' \widehat{\mathbf{A}} = 0$  as desired.

The solution  $\widehat{\mathbf{A}}_\perp$  in (10.22) can be represented several ways. One which is computationally convenient is to observe that

$$\widetilde{\mathbf{Y}}' \widetilde{\mathbf{Y}} - \widetilde{\mathbf{Y}}' \widetilde{\mathbf{X}} \left( \widetilde{\mathbf{X}}' \widetilde{\mathbf{X}} \right)^{-1} \widetilde{\mathbf{Y}}' \widetilde{\mathbf{X}} = \mathbf{Y}' \mathbf{M}_{\mathbf{X}, \mathbf{Z}} \mathbf{Y} = \widetilde{\mathbf{e}}' \widetilde{\mathbf{e}}$$

where  $\mathbf{M}_{\mathbf{X}, \mathbf{Z}} = \mathbf{I}_n - (\mathbf{X}, \mathbf{Z}) \left( (\mathbf{X}, \mathbf{Z})' (\mathbf{X}, \mathbf{Z}) \right)^{-1} (\mathbf{X}, \mathbf{Z})'$  and  $\widetilde{\mathbf{e}} = \mathbf{M}_{\mathbf{X}, \mathbf{Z}} \mathbf{Y}$  is the residual from the unrestricted least-squares regression of  $\mathbf{Y}$  on  $\mathbf{X}$  and  $\mathbf{Z}$ . The first equality follows by the Frisch-Waugh-Lovell theorem. This shows that  $\widehat{\mathbf{A}}_\perp$  are the generalized eigenvectors of  $\widetilde{\mathbf{e}}' \widetilde{\mathbf{e}}$  with respect to  $\widetilde{\mathbf{Y}}' \widetilde{\mathbf{Y}}$  corresponding to the  $m - r$  largest eigenvalues. In MATLAB, for example, these can be computed using the `eig(A,B)` command.

Another representation is to write  $\mathbf{M}_{\mathbf{Z}} = \mathbf{I}_n - \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}'$  so that

$$\widehat{\mathbf{A}}_\perp = \underset{\mathbf{A}}{\text{argmax}} \frac{\det(\mathbf{A}' \mathbf{Y}' \mathbf{M}_{\mathbf{X}, \mathbf{Z}} \mathbf{Y} \mathbf{A})}{\det(\mathbf{A}' \mathbf{Y}' \mathbf{M}_{\mathbf{Z}} \mathbf{Y} \mathbf{A})} = \underset{\mathbf{A}}{\text{argmin}} \frac{\det(\mathbf{A}' \mathbf{Y}' \mathbf{M}_{\mathbf{Z}} \mathbf{Y} \mathbf{A})}{\det(\mathbf{A}' \mathbf{Y}' \mathbf{M}_{\mathbf{X}, \mathbf{Z}} \mathbf{Y} \mathbf{A})}$$

We summarize our findings.

**Theorem 10.9.1** *The MLE for the reduced rank model (10.18) under  $\mathbf{e}_i \sim N(\mathbf{0}, \mathbf{\Sigma})$  is given as follows.  $\widehat{\mathbf{G}} = \{\mathbf{v}_1, \dots, \mathbf{v}_r\}$ , the generalized eigenvectors of  $\widetilde{\mathbf{X}}' \widetilde{\mathbf{Y}} \left( \widetilde{\mathbf{Y}}' \widetilde{\mathbf{Y}} \right)^{-1} \widetilde{\mathbf{Y}}' \widetilde{\mathbf{X}}$  with respect to  $\widetilde{\mathbf{X}}' \widetilde{\mathbf{X}}$  corresponding to the  $r$  largest eigenvalues.  $\widehat{\mathbf{A}}$ ,  $\widehat{\mathbf{C}}$  and  $\widehat{\mathbf{\Sigma}}$  are obtained by the least-squares regression*

$$\begin{aligned} \mathbf{y}_i &= \widehat{\mathbf{A}} \widehat{\mathbf{G}}' \mathbf{x}_i + \widehat{\mathbf{C}}' \mathbf{z}_i + \widehat{\mathbf{e}}_i \\ \widehat{\mathbf{\Sigma}} &= \frac{1}{n} \sum_{i=1}^n \widehat{\mathbf{e}}_i \widehat{\mathbf{e}}_i'. \end{aligned}$$

$\widehat{\mathbf{A}}_\perp$  equals the generalized eigenvectors of  $\widetilde{\mathbf{e}}' \widetilde{\mathbf{e}}$  with respect to  $\widetilde{\mathbf{Y}}' \widetilde{\mathbf{Y}}$  corresponding to the  $m - r$  smallest eigenvalues.

## Exercises

**Exercise 10.1** Show (10.9) when the errors are conditionally homoskedastic (10.7).

**Exercise 10.2** Show (10.10) when the regressors are common across equations  $\mathbf{x}_{ji} = \mathbf{x}_i$

**Exercise 10.3** Show (10.11) when the regressors are common across equations  $\mathbf{x}_{ji} = \mathbf{x}_i$  and the errors are conditionally homoskedastic (10.7).

**Exercise 10.4** Prove Theorem 10.5.1.

**Exercise 10.5** Show (10.12) when the regressors are common across equations  $\mathbf{x}_{ji} = \mathbf{x}_i$

**Exercise 10.6** Show (10.13) when the regressors are common across equations  $\mathbf{x}_{ji} = \mathbf{x}_i$  and the errors are conditionally homoskedastic (10.7).

**Exercise 10.7** Prove Theorem 10.5.2.

**Exercise 10.8** Prove Theorem 10.6.1.

**Exercise 10.9** Show that (10.15) follows from the steps described.

**Exercise 10.10** Show that (10.16) follows from the steps described.

**Exercise 10.11** Prove Theorem 10.7.1.

**Exercise 10.12** Prove Theorem 10.7.2.

Hint: First, show that it is sufficient to show that

$$\mathbb{E}(\mathbf{X}_i \mathbf{X}_i') (\mathbb{E}(\mathbf{X}_i \boldsymbol{\Sigma}^{-1} \mathbf{X}_i'))^{-1} \mathbb{E}(\mathbf{X}_i \mathbf{X}_i') \leq \mathbb{E}(\mathbf{X}_i \boldsymbol{\Sigma} \mathbf{X}_i').$$

Second, rewrite this equation using the transformations  $\mathbf{U}_i = \mathbf{X}_i \boldsymbol{\Sigma}^{1/2}$  and  $\mathbf{V}_i = \mathbf{X}_i \boldsymbol{\Sigma}^{-1/2}$ , and then apply the matrix Cauchy-Schwarz inequality (B.11).

**Exercise 10.13** Prove Theorem 10.7.3

**Exercise 10.14** Take the model

$$\begin{aligned} y_i &= \boldsymbol{\pi}_i' \boldsymbol{\beta} + e_i \\ \boldsymbol{\pi}_i &= \mathbb{E}(\mathbf{x}_i | \mathbf{z}_i) = \boldsymbol{\Gamma}' \mathbf{z}_i \\ \mathbb{E}(e_i | \mathbf{z}_i) &= 0 \end{aligned}$$

where  $y_i$ ,  $i$  is scalar,  $\mathbf{x}_i$  is a  $k$  vector and  $\mathbf{z}_i$  is an  $\ell$  vector.  $\boldsymbol{\beta}$  and  $\boldsymbol{\pi}_i$  are  $k \times 1$  and  $\boldsymbol{\Gamma}$  is  $\ell \times k$ . The sample is  $(y_i, \mathbf{x}_i, \mathbf{z}_i : i = 1, \dots, n)$  with  $\boldsymbol{\pi}_i$  unobserved.

Consider the estimator  $\hat{\boldsymbol{\beta}}$  for  $\boldsymbol{\beta}$  by OLS of  $y_i$  on  $\hat{\boldsymbol{\pi}}_i = \hat{\boldsymbol{\Gamma}}' \mathbf{z}_i$  where  $\hat{\boldsymbol{\Gamma}}$  is the OLS coefficient from the multivariate regression of  $\mathbf{x}_i$  on  $\mathbf{z}_i$

- Show that  $\hat{\boldsymbol{\beta}}$  is consistent for  $\boldsymbol{\beta}$
- Find the asymptotic distribution  $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$  as  $n \rightarrow \infty$  assuming that  $\boldsymbol{\beta} = \mathbf{0}$ .
- Why is the assumption  $\boldsymbol{\beta} = \mathbf{0}$  an important simplifying condition in part (b)?
- Using the result in (c), construct an appropriate asymptotic test for the hypothesis  $\mathbb{H}_0 : \boldsymbol{\beta} = \mathbf{0}$ .

**Exercise 10.15** The observations are iid,  $(y_{1i}, y_{2i}, \mathbf{x}_i : i = 1, \dots, n)$ . The dependent variables  $y_{1i}$  and  $y_{2i}$  are real-valued. The regressor  $\mathbf{x}_i$  is a  $k$ -vector. The model is the two-equation system

$$\begin{aligned}y_{1i} &= \mathbf{x}_i' \boldsymbol{\beta}_1 + e_{1i} \\ \mathbb{E}(\mathbf{x}_i e_{1i}) &= 0 \\ y_{2i} &= \mathbf{x}_i' \boldsymbol{\beta}_2 + e_{2i} \\ \mathbb{E}(\mathbf{x}_i e_{2i}) &= 0\end{aligned}$$

- (a) What are the appropriate estimators  $\hat{\boldsymbol{\beta}}_1$  and  $\hat{\boldsymbol{\beta}}_2$  for  $\boldsymbol{\beta}_1$  and  $\boldsymbol{\beta}_2$ ?
- (b) Find the joint asymptotic distribution of  $\hat{\boldsymbol{\beta}}_1$  and  $\hat{\boldsymbol{\beta}}_2$ .
- (c) Describe a test for  $\mathbb{H}_0 : \boldsymbol{\beta}_1 = \boldsymbol{\beta}_2$ .



# Chapter 11

## Instrumental Variables

### 11.1 Introduction

We say that there is **endogeneity** in the linear model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i \quad (11.1)$$

if  $\boldsymbol{\beta}$  is the parameter of interest and

$$\mathbb{E}(\mathbf{x}_i e_i) \neq \mathbf{0}. \quad (11.2)$$

This is a core problem in econometrics and largely differentiates econometrics from many branches of statistics. To distinguish (11.1) from the regression and projection models, we will call (11.1) a **structural equation** and  $\boldsymbol{\beta}$  a **structural parameter**. When (11.2) holds, it is typical to say that  $\mathbf{x}_i$  is **endogenous** for  $\boldsymbol{\beta}$ .

Endogeneity cannot happen if the coefficient is defined by linear projection. Indeed, we can define the linear projection coefficient  $\boldsymbol{\beta}^* = \mathbb{E}(\mathbf{x}_i \mathbf{x}_i')^{-1} \mathbb{E}(\mathbf{x}_i y_i)$  and linear projection equation

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta}^* + e_i^* \\ \mathbb{E}(\mathbf{x}_i e_i^*) &= \mathbf{0}. \end{aligned}$$

However, under endogeneity (11.2) the projection coefficient  $\boldsymbol{\beta}^*$  does not equal the structural parameter. Indeed,

$$\begin{aligned} \boldsymbol{\beta}^* &= (\mathbb{E}(\mathbf{x}_i \mathbf{x}_i'))^{-1} \mathbb{E}(\mathbf{x}_i y_i) \\ &= (\mathbb{E}(\mathbf{x}_i \mathbf{x}_i'))^{-1} \mathbb{E}(\mathbf{x}_i (\mathbf{x}_i' \boldsymbol{\beta} + e_i)) \\ &= \boldsymbol{\beta} + (\mathbb{E}(\mathbf{x}_i \mathbf{x}_i'))^{-1} \mathbb{E}(\mathbf{x}_i e_i) \\ &\neq \boldsymbol{\beta} \end{aligned}$$

the final relation since  $\mathbb{E}(\mathbf{x}_i e_i) \neq \mathbf{0}$ .

Thus endogeneity requires that the coefficient be defined differently than projection. We describe such definitions as **structural**. We will present three examples in the following section.

Endogeneity implies that the least-squares estimator is inconsistent for the structural parameter. Indeed, under i.i.d. sampling, least-squares is consistent for the projection coefficient, and thus is inconsistent for  $\boldsymbol{\beta}$ .

$$\hat{\boldsymbol{\beta}} \xrightarrow{p} (\mathbb{E}(\mathbf{x}_i \mathbf{x}_i'))^{-1} \mathbb{E}(\mathbf{x}_i y_i) = \boldsymbol{\beta}^* \neq \boldsymbol{\beta}.$$

The inconsistency of least-squares is typically referred to as **endogeneity bias** or **estimation bias** due to endogeneity. (This is an imperfect label as the actual issue is inconsistency, not bias.)

As the structural parameter  $\boldsymbol{\beta}$  is the parameter of interest, endogeneity requires the development of alternative estimation methods. We discuss those in later sections.

## 11.2 Examples

The concept of endogeneity may be easiest to understand by example. We discuss three distinct examples. In each case it is important to see how the structural parameter  $\beta$  is defined independently from the linear projection model.

**Example: Measurement error in the regressor.** Suppose that  $(y_i, z_i)$  are joint random variables,  $\mathbb{E}(y_i | z_i) = z_i' \beta$  is linear,  $\beta$  is the structural parameter, and  $z_i$  is not observed. Instead we observe  $x_i = z_i + u_i$  where  $u_i$  is a  $k \times 1$  measurement error, independent of  $e_i$  and  $z_i$ . This is an example of a latent variable model, where “latent” refers to a structural variable which is unobserved.

The model  $x_i = z_i + u_i$  with  $z_i$  and  $u_i$  independent and  $\mathbb{E}(u_i) = \mathbf{0}$  is known as **classical measurement error**. This means that  $x_i$  is a noisy but unbiased measure of  $z_i$ .

By substitution we can express  $y_i$  as a function of the observed variable  $x_i$ .

$$\begin{aligned} y_i &= z_i' \beta + e_i \\ &= (x_i - u_i)' \beta + e_i \\ &= x_i' \beta + v_i \end{aligned}$$

where  $v_i = e_i - u_i' \beta$ . This means that  $(y_i, x_i)$  satisfy the linear equation

$$y_i = x_i' \beta + v_i$$

with an error  $v_i$ . But this error is not a projection error. Indeed,

$$\mathbb{E}(x_i v_i) = \mathbb{E}[(z_i + u_i)(e_i - u_i' \beta)] = -\mathbb{E}(u_i u_i') \beta \neq \mathbf{0}$$

if  $\beta \neq \mathbf{0}$  and  $\mathbb{E}(u_i u_i') \neq 0$ . As we learned in the previous section, if  $\mathbb{E}(x_i v_i) \neq 0$  then least-squares estimation will be inconsistent.

We can calculate the form of the projection coefficient (which is consistently estimated by least-squares). For simplicity suppose that  $k = 1$ . We find

$$\beta^* = \beta + \frac{\mathbb{E}(x_i v_i)}{\mathbb{E}(x_i^2)} = \beta \left( 1 - \frac{\mathbb{E}(u_i^2)}{\mathbb{E}(x_i^2)} \right).$$

Since  $\mathbb{E}(u_i^2) / \mathbb{E}(x_i^2) < 1$  the projection coefficient shrinks the structural parameter  $\beta$  towards zero. This is called **measurement error bias** or **attenuation bias**.

**Example: Supply and Demand.** The variables  $q_i$  and  $p_i$  (quantity and price) are determined jointly by the demand equation

$$q_i = -\beta_1 p_i + e_{1i}$$

and the supply equation

$$q_i = \beta_2 p_i + e_{2i}.$$

Assume that  $e_i = \begin{pmatrix} e_{1i} \\ e_{2i} \end{pmatrix}$  is i.i.d.,  $\mathbb{E}(e_i) = \mathbf{0}$  and  $\mathbb{E}(e_i e_i') = \mathbf{I}_2$  (the latter for simplicity). The question is: if we regress  $q_i$  on  $p_i$ , what happens?

It is helpful to solve for  $q_i$  and  $p_i$  in terms of the errors. In matrix notation,

$$\begin{bmatrix} 1 & \beta_1 \\ 1 & -\beta_2 \end{bmatrix} \begin{pmatrix} q_i \\ p_i \end{pmatrix} = \begin{pmatrix} e_{1i} \\ e_{2i} \end{pmatrix}$$

so

$$\begin{aligned} \begin{pmatrix} q_i \\ p_i \end{pmatrix} &= \begin{bmatrix} 1 & \beta_1 \\ 1 & -\beta_2 \end{bmatrix}^{-1} \begin{pmatrix} e_{1i} \\ e_{2i} \end{pmatrix} \\ &= \begin{bmatrix} \beta_2 & \beta_1 \\ 1 & -1 \end{bmatrix} \begin{pmatrix} e_{1i} \\ e_{2i} \end{pmatrix} \begin{pmatrix} 1 \\ \beta_1 + \beta_2 \end{pmatrix} \\ &= \begin{pmatrix} (\beta_2 e_{1i} + \beta_1 e_{2i}) / (\beta_1 + \beta_2) \\ (e_{1i} - e_{2i}) / (\beta_1 + \beta_2) \end{pmatrix}. \end{aligned}$$

The projection of  $q_i$  on  $p_i$  yields

$$\begin{aligned} q_i &= \beta^* p_i + e_i^* \\ \mathbb{E}(p_i e_i^*) &= 0 \end{aligned}$$

where

$$\beta^* = \frac{\mathbb{E}(p_i q_i)}{\mathbb{E}(p_i^2)} = \frac{\beta_2 - \beta_1}{2}.$$

Thus the projection coefficient  $\beta^*$  equals neither the demand slope  $\beta_1$  nor the supply slope  $\beta_2$ , but equals an average of the two. (The fact that it is a simple average is an artifact of the simple covariance structure.)

Hence the OLS estimate satisfies  $\hat{\beta} \xrightarrow{p} \beta^*$ , and the limit does not equal either  $\beta_1$  or  $\beta_2$ . The fact that the limit is neither the supply nor demand slope is called **simultaneous equations bias**. This occurs generally when  $y_i$  and  $x_i$  are jointly determined, as in a market equilibrium.

Generally, when both the dependent variable and a regressor are simultaneously determined, then the variables should be treated as endogenous.

**Example: Choice Variables as Regressors.** Take the classic wage equation

$$\log(\text{wage}) = \beta \text{education} + e$$

with  $\beta$  the average causal effect of education on wages. If wages are affected by unobserved ability, and individuals with high ability self-select into higher education, then  $e$  contains unobserved ability, so *education* and  $e$  will be positively correlated. Hence *education* is endogenous. The positive correlation means that the linear projection coefficient  $\beta^*$  will be upward biased relative to the structural coefficient  $\beta$ . Thus least-squares (which is estimating the projection coefficient) will tend to over-estimate the causal effect of education on wages.

This type of endogeneity occurs generally when  $y$  and  $x$  are both choices made by an economic agent, even if they are made at different points in time.

Generally, when both the dependent variable and a regressor are choice variables made by the same agent, the variables should be treated as endogenous.

### 11.3 Instrumental Variables

We have defined endogeneity as the context where the regressor is correlated with the equation error. In most applications we only treat a subset of the regressors as endogenous; most of the regressors will be treated as **exogenous**, meaning that they are assumed uncorrelated with the equation error. To be specific, we make the partition

$$\mathbf{x}_i = \begin{pmatrix} \mathbf{x}_{1i} \\ \mathbf{x}_{2i} \end{pmatrix} \begin{matrix} k_1 \\ k_2 \end{matrix} \quad (11.3)$$

and similarly

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} \begin{matrix} k_1 \\ k_2 \end{matrix}$$

so that the **structural equation** is

$$\begin{aligned} y_i &= \mathbf{x}'_i \boldsymbol{\beta} + e_i \\ &= \mathbf{x}'_{1i} \boldsymbol{\beta}_1 + \mathbf{x}'_{2i} \boldsymbol{\beta}_2 + e_i. \end{aligned} \quad (11.4)$$

The regressors are assumed to satisfy

$$\begin{aligned} \mathbb{E}(\mathbf{x}_{1i} e_i) &= \mathbf{0} \\ \mathbb{E}(\mathbf{x}_{2i} e_i) &\neq \mathbf{0}. \end{aligned}$$

We call  $\mathbf{x}_{1i}$  **exogenous** and  $\mathbf{x}_{2i}$  **endogenous** for the structural parameter  $\boldsymbol{\beta}$ . As the dependent variable  $y_i$  is also endogenous, we sometimes differentiate  $\mathbf{x}_{2i}$  by calling  $\mathbf{x}_{2i}$  the **endogenous right-hand-side variables**.

In matrix notation we can write (11.4) as

$$\begin{aligned} \mathbf{y} &= \mathbf{X} \boldsymbol{\beta} + \mathbf{e} \\ &= \mathbf{X}_1 \boldsymbol{\beta}_1 + \mathbf{X}_2 \boldsymbol{\beta}_2 + \mathbf{e}. \end{aligned} \quad (11.5)$$

The endogenous regressors  $\mathbf{x}_{2i}$  are the critical variables discussed in the examples of the previous section – simultaneous variables, choice variables, mis-measured regressors – that are potentially correlated with the equation error  $e_i$ . In most applications the number  $k_2$  of variables treated as endogenous is small (1 or 2). The exogenous variables  $\mathbf{x}_{1i}$  are the remaining regressors (including the equation intercept) and can be low or high dimensional.

To consistently estimate  $\boldsymbol{\beta}$  we require additional information. One type of information which is commonly used in economic applications are what we call **instruments**.

**Definition 11.3.1** *The  $\ell \times 1$  random vector  $\mathbf{z}_i$  is an **instrumental variable** for (11.4) if*

$$\mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0} \quad (11.6)$$

$$\mathbb{E}(\mathbf{z}_i \mathbf{z}'_i) > 0 \quad (11.7)$$

$$\text{rank}(\mathbb{E}(\mathbf{z}_i \mathbf{z}'_i)) = k. \quad (11.8)$$

There are three components to the definition as given. The first (11.6) is that the instruments are uncorrelated with the regression error. The second (11.7) is a normalization which excludes linearly redundant instruments. The third (11.8) is often called the **relevance condition** and is essential for the identification of the model, as we discuss later. A necessary condition for (11.8) is that  $\ell \geq k$ .

Condition (11.6) – that the instruments are uncorrelated with the equation error, is often described as that they are **exogenous** in the sense that they are determined outside the model for  $y_i$ .

Notice that the regressors  $\mathbf{x}_{1i}$  satisfy condition (11.6) and thus should be included as instrumental variables. It is thus a subset of the variables  $\mathbf{z}_i$ . Notationally we make the partition

$$\mathbf{z}_i = \begin{pmatrix} \mathbf{z}_{1i} \\ \mathbf{z}_{2i} \end{pmatrix} = \begin{pmatrix} \mathbf{x}_{1i} \\ \mathbf{z}_{2i} \end{pmatrix} \begin{pmatrix} k_1 \\ \ell_2 \end{pmatrix}. \quad (11.9)$$

Here,  $\mathbf{x}_{1i} = \mathbf{z}_{1i}$  are the **included exogenous variables**, and  $\mathbf{z}_{2i}$  are the **excluded exogenous variables**. That is,  $\mathbf{z}_{2i}$  are variables which could be included in the equation for  $y_i$  (in the sense

that they are uncorrelated with  $e_i$ ) yet can be excluded, as they would have true zero coefficients in the equation.

Many authors simply label  $\mathbf{x}_{1i}$  as the “exogenous variables”,  $\mathbf{x}_{2i}$  as the “endogenous variables”, and  $\mathbf{z}_{2i}$  as the “instrumental variables”.

We say that the model is **just-identified** if  $\ell = k$  (and  $\ell_2 = k_2$ ) and **over-identified** if  $\ell > k$  (and  $\ell_2 > k_2$ ).

What variables can be used as instrumental variables? From the definition  $\mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0}$  we see that the instrument must be uncorrelated with the equation error, meaning that it is excluded from the structural equation as mentioned above. From the rank condition (11.8) it is also important that the instrumental variable be correlated with the endogenous variables  $\mathbf{x}_{2i}$  after controlling for the other exogenous variables  $\mathbf{x}_{1i}$ . These two requirements are typically interpreted as requiring that the instruments be determined outside the system for  $(y_i, \mathbf{x}_{2i})$ , causally determine  $\mathbf{x}_{2i}$ , but do not causally determine  $y_i$  except through  $\mathbf{x}_{2i}$ .

Let’s take the three examples given above.

**Measurement error in the regressor.** When  $\mathbf{x}_i$  is a mis-measured version of  $\mathbf{z}_i$ , a common choice for an instrument  $\mathbf{z}_{2i}$  is an alternative measurement of  $\mathbf{z}_i$ . For this  $\mathbf{z}_{2i}$  to satisfy the property of an instrumental variable the measurement error in  $\mathbf{z}_{2i}$  must be independent of that in  $\mathbf{x}_i$ .

**Supply and Demand.** An appropriate instrument for price  $p_i$  in a demand equation is a variable  $z_{2i}$  which influences supply but not demand. Such a variable affects the equilibrium values of  $p_i$  and  $q_i$  but does not directly affect price except through quantity. Variables which affect supply but not demand are typically related to production costs.

An appropriate instrument for price in a supply equation is a variable which influences demand but not supply. Such a variable affects the equilibrium values of price and quantity but only affects price through quantity.

**Choice Variable as Regressor.** An ideal instrument affects the choice of the regressor (education) but does not directly influence the dependent variable (wages) except through the indirect effect on the regressor. We will discuss an example in the next section.

## 11.4 Example: College Proximity

In a influential paper, David Card (1995) suggested if a potential student lives close to a college, this reduces the cost of attendance and thereby raises the likelihood that the student will attend college. However, college proximity does not directly affect a student’s skills or abilities, so should not have a direct effect on his or her market wage. These considerations suggest that college proximity can be used as an instrument for education in a wage regression. We use the simplest model reported in Card’s paper to illustrate the concepts of instrumental variables throughout the chapter.

Card used data from the National Longitudinal Survey of Young Men (NLSYM) for 1976. A baseline least-squares wage regression for his data set is reported in the first column of Table 11.1. The dependent variable is the log of weekly earnings. The regressors are *education* (years of schooling), *experience* (years of work experience, calculated as *age* (years) less *education+6*), *experience*<sup>2</sup>/100, *black*, *south* (an indicator for residence in the southern region of the U.S.), and *urban* (an indicator for residence in a standard metropolitan statistical area). We drop observations for which *wage* is missing. The remaining sample has 3,010 observations. His data is the file **Card1995** on the textbook website.

The point estimate obtained by least-squares suggests an 8% increase in earnings for each year of education.

Table 11.1  
Dependent variable  $\log(wage)$

	OLS	IV(a)	IV(b)	2SLS(a)	2SLS(b)	LIML
education	0.074 (0.004)	0.132 (0.049)	0.133 (0.051)	0.161 (0.040)	0.160 (0.041)	0.164 (0.042)
experience	0.084 (0.007)	0.107 (0.021)	0.056 (0.026)	0.119 (0.018)	0.047 (0.025)	0.120 (0.019)
experience <sup>2</sup> /100	-0.224 (0.032)	-0.228 (0.035)	-0.080 (0.133)	-0.231 (0.037)	-0.032 (0.127)	-0.231 (0.037)
black	-0.190 (0.017)	-0.131 (0.051)	-0.103 (0.075)	-0.102 (0.044)	-0.064 (0.061)	-0.099 (0.045)
south	-0.125 (0.015)	-0.105 (0.023)	-0.098 (0.0287)	-0.095 (0.022)	-0.086 (0.026)	-0.094 (0.022)
urban	0.161 (0.015)	0.131 (0.030)	0.108 (0.049)	0.116 (0.026)	0.083 (0.041)	0.115 (0.027)
Sargan				0.82	0.52	0.82
p-value				0.36	0.47	0.37

Notes:

1. IV(a) uses *college* as an instrument for *education*.
2. IV(b) uses *college*, *age*, and *age*<sup>2</sup> as instruments for *education*, *experience*, and *experience*<sup>2</sup>/100.
3. 2SLS(a) uses *public* and *private* as instruments for *education*.
4. 2SLS(b) uses *public*, *private*, *age*, and *age*<sup>2</sup> as instruments for *education*, *experience*, and *experience*<sup>2</sup>/100.
5. LIML uses *public* and *private* as instruments for *education*.

As discussed in the previous sections, it is reasonable to view years of education as a choice made by an individual, and thus is likely endogenous for the structural return to education. This means that least-squares is an estimate of a linear projection, but is inconsistent for coefficient of a structural equation representing the causal impact of years of education on expected wages. Labor economics predicts that ability, education, and wages will be positively correlated. This suggests that the population projection coefficient estimated by least-squares will be higher than the structural parameter (and hence upwards biased). However, the sign of the bias is uncertain since there are multiple regressors and there are other potential sources of endogeneity.

To instrument for the endogeneity of education, Card suggested that a reasonable instrument is a dummy variable indicating if the individual grew up near a college. We will consider three measures:

- college*    Grew up in same county as a 4-year college
- public*    Grew up in same county as a 4-year public college
- private*   Grew up in same county as a 4-year private college.

### David Card

David Card (1956- ) is a Canadian-American labor economist whose research has changed our understanding of labor markets, the impact of minimum wage legislation, and immigration. His methodological innovations in applied econometrics have transformed empirical microeconomics.

## 11.5 Reduced Form

The reduced form is the relationship between the regressors  $\mathbf{x}_i$  and the instruments  $\mathbf{z}_i$ . A linear reduced form model for  $\mathbf{x}_i$  is

$$\mathbf{x}_i = \mathbf{\Gamma}' \mathbf{z}_i + \mathbf{u}_i. \quad (11.10)$$

This is a multivariate regression as introduced in Chapter 10. The  $\ell \times k$  coefficient matrix  $\mathbf{\Gamma}$  can be defined by linear projection. Thus

$$\mathbf{\Gamma} = \mathbb{E}(\mathbf{z}_i \mathbf{z}_i')^{-1} \mathbb{E}(\mathbf{z}_i \mathbf{x}_i') \quad (11.11)$$

so that

$$\mathbb{E}(\mathbf{z}_i \mathbf{u}_i') = \mathbf{0}.$$

In matrix notation, we can write (11.10) as

$$\mathbf{X} = \mathbf{Z}\mathbf{\Gamma} + \mathbf{U} \quad (11.12)$$

where  $\mathbf{U}$  is  $n \times k$ . Notice that the projection coefficient (11.11) is well defined and unique under (11.7).

Since  $\mathbf{z}_i$  and  $\mathbf{x}_i$  have the common variables  $\mathbf{x}_{1i}$ , we can focus on the reduced form for the endogenous regressors  $\mathbf{x}_{2i}$ . Recalling the partitions (11.3) and (11.9) we can partition  $\mathbf{\Gamma}$  conformably as

$$\begin{aligned} \mathbf{\Gamma} &= \begin{bmatrix} & k_1 & k_2 \\ \mathbf{\Gamma}_{11} & \mathbf{\Gamma}_{12} \\ \mathbf{\Gamma}_{21} & \mathbf{\Gamma}_{22} \end{bmatrix} \begin{matrix} \ell_1 \\ \ell_2 \end{matrix} \\ &= \begin{bmatrix} \mathbf{I} & \mathbf{\Gamma}_{12} \\ \mathbf{0} & \mathbf{\Gamma}_{22} \end{bmatrix} \end{aligned} \quad (11.13)$$

and similarly partition  $\mathbf{u}_i$ . Then (11.10) can be rewritten as two equation systems

$$\mathbf{x}_{1i} = \mathbf{z}_{1i} \quad (11.14)$$

$$\mathbf{x}_{2i} = \mathbf{\Gamma}_{12}' \mathbf{z}_{1i} + \mathbf{\Gamma}_{22}' \mathbf{z}_{2i} + \mathbf{u}_{2i}. \quad (11.15)$$

The first equation (11.14) is a tautology. The second equation (11.15) is the primary reduced form equation of interest. It is a multivariate linear regression for  $\mathbf{x}_{2i}$  as a function of the included and excluded exogenous variables  $\mathbf{z}_{1i}$  and  $\mathbf{z}_{2i}$ .

We can also construct a reduced form equation for  $y_i$ . Substituting (11.10) into (11.4), we find

$$\begin{aligned} y_i &= (\mathbf{\Gamma}' \mathbf{z}_i + \mathbf{u}_i)' \boldsymbol{\beta} + e_i \\ &= \mathbf{z}_i' \boldsymbol{\lambda} + v_i \end{aligned} \quad (11.16)$$

where

$$\boldsymbol{\lambda} = \mathbf{\Gamma} \boldsymbol{\beta} \quad (11.17)$$

and

$$v_i = \mathbf{u}_i' \boldsymbol{\beta} + e_i.$$

Observe that

$$\mathbb{E}(\mathbf{z}_i v_i) = \mathbb{E}(\mathbf{z}_i \mathbf{u}_i') \boldsymbol{\beta} + \mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0}.$$

Thus (11.16) is a projection equation. It is the reduced form for  $y_i$ , as it expresses  $y_i$  as a function of exogenous variables only. Since it is a projection equation we can write the reduced form coefficient as

$$\boldsymbol{\lambda} = \mathbb{E}(\mathbf{z}_i \mathbf{z}_i')^{-1} \mathbb{E}(\mathbf{z}_i y_i) \quad (11.18)$$

which is well defined under (11.7).

Alternatively, we can substitute (11.15) into (11.4) and use  $\mathbf{x}_{1i} = \mathbf{z}_{1i}$  to obtain

$$\begin{aligned} y_i &= \mathbf{x}'_{1i}\boldsymbol{\beta}_1 + (\boldsymbol{\Gamma}'_{12}\mathbf{z}_{1i} + \boldsymbol{\Gamma}'_{22}\mathbf{z}_{2i} + \mathbf{u}_{2i})'\boldsymbol{\beta}_2 + e_i \\ &= \mathbf{z}'_{1i}\boldsymbol{\lambda}_1 + \mathbf{z}'_{2i}\boldsymbol{\lambda}_2 + v_i \end{aligned} \quad (11.19)$$

where

$$\boldsymbol{\lambda}_1 = \boldsymbol{\beta}_1 + \boldsymbol{\Gamma}_{12}\boldsymbol{\beta}_2 \quad (11.20)$$

$$\boldsymbol{\lambda}_2 = \boldsymbol{\Gamma}_{22}\boldsymbol{\beta}_2. \quad (11.21)$$

which is an alternative (and equivalent) expression of (11.17) given (11.13).

(11.10) and (11.16) together (or (11.15) and (11.19) together) are the **reduced form equations** for the system

$$\begin{aligned} y_i &= \mathbf{z}'_i\boldsymbol{\lambda} + v_i \\ \mathbf{x}_i &= \boldsymbol{\Gamma}'\mathbf{z}_i + \mathbf{u}_i. \end{aligned}$$

The relationships (11.17) and (11.20)-(11.21) are critically important for understanding the identification of the structural parameters  $\boldsymbol{\beta}_1$  and  $\boldsymbol{\beta}_2$ , as we discuss below. These equations show the tight relationship between the parameters of the structural equations ( $\boldsymbol{\beta}_1$  and  $\boldsymbol{\beta}_2$ ) and those of the reduced form equations ( $\boldsymbol{\lambda}_1$ ,  $\boldsymbol{\lambda}_2$ ,  $\boldsymbol{\Gamma}_{12}$  and  $\boldsymbol{\Gamma}_{22}$ ).

## 11.6 Reduced Form Estimation

The reduced form equations are projections, so the coefficient matrices may be estimated by least-squares (see Chapter 10). The least-squares estimate of (11.10) is

$$\hat{\boldsymbol{\Gamma}} = \left( \sum_{i=1}^n \mathbf{z}_i \mathbf{z}'_i \right)^{-1} \left( \sum_{i=1}^n \mathbf{z}_i \mathbf{x}'_i \right). \quad (11.22)$$

The estimates of equation (11.10) can be written as

$$\mathbf{x}_i = \hat{\boldsymbol{\Gamma}}'\mathbf{z}_i + \hat{\mathbf{u}}_i. \quad (11.23)$$

In matrix notation, these can be written as

$$\hat{\boldsymbol{\Gamma}} = (\mathbf{Z}'\mathbf{Z})^{-1}(\mathbf{Z}'\mathbf{X})$$

and

$$\mathbf{X} = \mathbf{Z}\hat{\boldsymbol{\Gamma}} + \hat{\mathbf{U}}.$$

Since  $\mathbf{X}$  and  $\mathbf{Z}$  have a common sub-matrix, we have the partition

$$\hat{\boldsymbol{\Gamma}} = \begin{bmatrix} \mathbf{I} & \hat{\boldsymbol{\Gamma}}_{12} \\ \mathbf{0} & \hat{\boldsymbol{\Gamma}}_{22} \end{bmatrix}.$$

The reduced form estimates of equation (11.15) can be written as

$$\mathbf{x}_{2i} = \hat{\boldsymbol{\Gamma}}'_{12}\mathbf{z}_{1i} + \hat{\boldsymbol{\Gamma}}'_{22}\mathbf{z}_{2i} + \hat{\mathbf{u}}_{2i}$$

or in matrix notation as

$$\mathbf{X}_2 = \mathbf{Z}_1\hat{\boldsymbol{\Gamma}}_{12} + \mathbf{Z}_2\hat{\boldsymbol{\Gamma}}_{22} + \hat{\mathbf{U}}_2.$$



We can write the submatrix estimates as

$$\begin{bmatrix} \hat{\Gamma}_{12} \\ \hat{\Gamma}_{22} \end{bmatrix} = \left( \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i' \right)^{-1} \begin{pmatrix} \sum_{i=1}^n \mathbf{z}_i \mathbf{x}_{2i}' \\ \sum_{i=1}^n \mathbf{z}_i y_i \end{pmatrix} = (\mathbf{Z}'\mathbf{Z})^{-1} (\mathbf{Z}'\mathbf{X}_2).$$

The reduced form estimate of equation (11.16) is

$$\begin{aligned} \hat{\lambda} &= \left( \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i' \right)^{-1} \begin{pmatrix} \sum_{i=1}^n \mathbf{z}_i y_i \end{pmatrix} \\ y_i &= \mathbf{z}_i' \hat{\lambda} + \hat{v}_i \\ &= \mathbf{z}_{1i}' \hat{\lambda}_1 + \mathbf{z}_{2i}' \hat{\lambda}_2 + \hat{v}_i \end{aligned}$$

or in matrix notation

$$\begin{aligned} \hat{\lambda} &= (\mathbf{Z}'\mathbf{Z})^{-1} (\mathbf{Z}'\mathbf{y}) \\ \mathbf{y} &= \mathbf{Z}\hat{\lambda} + \hat{\mathbf{v}} \\ &= \mathbf{Z}_1\hat{\lambda}_1 + \mathbf{Z}_2\hat{\lambda}_2 + \hat{\mathbf{v}}. \end{aligned}$$

## 11.7 Identification

A parameter is **identified** if it is a unique function of the probability distribution of the observables. One way to show that a parameter is identified is to write it as an explicit function of population moments. For example, the reduced form coefficient matrices  $\Gamma$  and  $\lambda$  are identified since they can be written as explicit functions of the moments of the observables  $(y_i, \mathbf{x}_i, \mathbf{z}_i)$ . That is,

$$\Gamma = \mathbb{E}(\mathbf{z}_i \mathbf{z}_i')^{-1} \mathbb{E}(\mathbf{z}_i \mathbf{x}_i') \quad (11.24)$$

$$\lambda = \mathbb{E}(\mathbf{z}_i \mathbf{z}_i')^{-1} \mathbb{E}(\mathbf{z}_i y_i). \quad (11.25)$$

These are uniquely determined by the probability distribution of  $(y_i, \mathbf{x}_i, \mathbf{z}_i)$  if Definition 11.3.1 holds, since this includes the requirement that  $\mathbb{E}(\mathbf{z}_i \mathbf{z}_i')$  is invertible.

We are interested in the structural parameter  $\beta$ . It relates to  $(\lambda, \Gamma)$  through (11.17), or

$$\lambda = \Gamma\beta. \quad (11.26)$$

It is identified if it uniquely determined by this relation. This is a set of  $\ell$  equations with  $k$  unknowns with  $\ell \geq k$ . From standard linear algebra we know that there is a unique solution if and only if  $\Gamma$  has full rank  $k$ .

$$\text{rank}(\Gamma) = k. \quad (11.27)$$

Under (11.27),  $\beta$  can be uniquely solved from the linear system  $\lambda = \Gamma\beta$ . On the other hand if  $\text{rank}(\Gamma) < k$  then  $\lambda = \Gamma\beta$  has fewer mutually independent linear equations than coefficients so there is not a unique solution.

From the definitions (11.24)-(11.25) the identification equation (11.26) is the same as

$$\mathbb{E}(\mathbf{z}_i y_i) = \mathbb{E}(\mathbf{z}_i \mathbf{x}_i') \beta$$

which is again a set of  $\ell$  equations with  $k$  unknowns. This has a unique solution if (and only if)

$$\text{rank}(\mathbb{E}(\mathbf{z}_i \mathbf{x}_i')) = k \quad (11.28)$$

which was listed in (11.8) as a conditions of Definition 11.3.1. (Indeed, this is why it was listed as part of the definition.) We can also see that (11.27) and (11.28) are equivalent ways of expressing the

same requirement. If this condition fails then  $\beta$  will not be identified. The condition (11.27)-(11.28) is called the **relevance condition**.

It is useful to have explicit expressions for the solution  $\beta$ . The easiest case is when  $\ell = k$ . Then (11.27) implies  $\Gamma$  is invertible, so the structural parameter equals  $\beta = \Gamma^{-1}\lambda$ . It is a unique solution because  $\Gamma$  and  $\lambda$  are unique and  $\Gamma$  is invertible.

When  $\ell > k$  we can solve for  $\beta$  by applying least-squares to the system of equations  $\lambda = \Gamma\beta$ . This is  $\ell$  equations with  $k$  unknowns and no error. The least-squares solution is  $\beta = (\Gamma'\Gamma)^{-1}\Gamma'\lambda$ . Under (11.27) the matrix  $\Gamma'\Gamma$  is invertible so the solution is unique.

$\beta$  is identified if  $\text{rank}(\Gamma) = k$ , which is true if and only if  $\text{rank}(\Gamma_{22}) = k_2$  (by the upper-diagonal structure of  $\Gamma$ ). Thus the key to identification of the model rests on the  $\ell_2 \times k_2$  matrix  $\Gamma_{22}$  in (11.15). To see this, recall the reduced form relationships (11.20)-(11.21). We can see that  $\beta_2$  is identified from (11.21) alone, and the necessary and sufficient condition is  $\text{rank}(\Gamma_{22}) = k_2$ . If this is satisfied then the solution can be written as  $\beta_2 = (\Gamma'_{22}\Gamma_{22})^{-1}\Gamma'_{22}\lambda_2$ . Then  $\beta_1$  is identified from this and (11.20), with the explicit solution  $\beta_1 = \lambda_1 - \Gamma_{12}(\Gamma'_{22}\Gamma_{22})^{-1}\Gamma'_{22}\lambda_2$ . In the just-identified case ( $\ell_2 = k_2$ ) these equations simplify to take the form  $\beta_2 = \Gamma_{22}^{-1}\lambda_2$  and  $\beta_1 = \lambda_1 - \Gamma_{12}\Gamma_{22}^{-1}\lambda_2$ .

## 11.8 Instrumental Variables Estimator

In this section we consider the special case where the model is just-identified, so that  $\ell = k$ .

The assumption that  $z_i$  is an instrumental variable implies that

$$\mathbb{E}(z_i e_i) = \mathbf{0}.$$

Making the substitution  $e_i = y_i - x'_i\beta$  we find

$$\mathbb{E}(z_i(y_i - x'_i\beta)) = \mathbf{0}.$$

Expanding,

$$\mathbb{E}(z_i y_i) - \mathbb{E}(z_i x'_i)\beta = \mathbf{0}.$$

This is a system of  $\ell = k$  equations and  $k$  unknowns. Solving for  $\beta$  we find

$$\beta = (\mathbb{E}(z_i x'_i))^{-1} \mathbb{E}(z_i y_i).$$

This solution assumes that the matrix  $\mathbb{E}(z_i x'_i)$  is invertible, which holds under (11.8) or equivalently (11.27).

The **instrumental variables** (IV) estimator  $\beta$  replaces the population moments by their sample versions. We find

$$\begin{aligned} \hat{\beta}_{\text{iv}} &= \left( \frac{1}{n} \sum_{i=1}^n z_i x'_i \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n z_i y_i \right) \\ &= \left( \sum_{i=1}^n z_i x'_i \right)^{-1} \left( \sum_{i=1}^n z_i y_i \right) \\ &= (\mathbf{Z}'\mathbf{X})^{-1} (\mathbf{Z}'\mathbf{y}). \end{aligned} \tag{11.29}$$

More generally, it is common to refer to any estimator of the form

$$\hat{\beta}_{\text{iv}} = (\mathbf{W}'\mathbf{X})^{-1} (\mathbf{W}'\mathbf{y})$$

given an  $n \times k$  matrix  $\mathbf{W}$  as an IV estimator for  $\beta$  using the instrument  $\mathbf{W}$ .

Alternatively, recall that when  $\ell = k$  the structural parameter can be written as a function of the reduced form parameters as  $\beta = \Gamma^{-1}\lambda$ . Replacing  $\Gamma$  and  $\lambda$  by their least-squares estimates we can construct what is called the **Indirect Least Squares** (ILS) estimator:

$$\begin{aligned}\hat{\beta}_{\text{ils}} &= \hat{\Gamma}^{-1}\hat{\lambda} \\ &= \left( (\mathbf{Z}'\mathbf{Z})^{-1} (\mathbf{Z}'\mathbf{X}) \right)^{-1} \left( (\mathbf{Z}'\mathbf{Z})^{-1} (\mathbf{Z}'\mathbf{y}) \right) \\ &= (\mathbf{Z}'\mathbf{X})^{-1} (\mathbf{Z}'\mathbf{Z}) (\mathbf{Z}'\mathbf{Z})^{-1} (\mathbf{Z}'\mathbf{y}) \\ &= (\mathbf{Z}'\mathbf{X})^{-1} (\mathbf{Z}'\mathbf{y}).\end{aligned}$$

We see that this equals the IV estimator (11.29). Thus the ILS and IV estimators are equivalent.

Given the IV estimator we define the residual vector

$$\hat{e} = \mathbf{y} - \mathbf{X}\hat{\beta}_{\text{iv}}$$

which satisfies

$$\mathbf{Z}'\hat{e} = \mathbf{Z}'\mathbf{y} - \mathbf{Z}'\mathbf{X} (\mathbf{Z}'\mathbf{X})^{-1} (\mathbf{Z}'\mathbf{y}) = \mathbf{0}. \quad (11.30)$$

Since  $\mathbf{Z}$  includes an intercept, this means that the residuals sum to zero, and are uncorrelated with the included and excluded instruments.

To illustrate, we estimate the reduced form equations corresponding to the college proximity example of Table 11.1, now treating *education* as endogenous and using *college* as an instrumental variable. The reduced form equations for  $\log(\text{wage})$  and *education* are reported in the first and second columns of Table 11.2.

Table 11.2  
Reduced Form Regressions

	$\log(\text{wage})$	<i>education</i>	<i>education</i>	<i>experience</i>	$\text{experience}^2/100$	<i>education</i>
<i>experience</i>	0.053 (0.007)	-0.410 (0.032)				-0.413 (0.032)
$\text{experience}^2/100$	-0.219 (0.033)	0.073 (0.170)				0.093 (0.171)
<i>black</i>	-0.264 (0.018)	-1.006 (0.088)	-1.468 (0.115)	1.468 (0.115)	0.282 (0.026)	-1.006 (0.088)
<i>south</i>	-0.143 (0.017)	-0.291 (0.078)	-0.460 (0.103)	0.460 (0.103)	0.112 (0.022)	-0.267 (0.079)
<i>urban</i>	0.185 (0.017)	0.404 (0.085)	0.835 (0.112)	-0.835 (0.112)	-0.176 (0.025)	0.400 (0.085)
<i>college</i>	0.045 (0.016)	0.337 (0.081)	0.347 (0.109)	-0.347 (0.109)	-0.073 (0.023)	
<i>public</i>						0.430 (0.086)
<i>private</i>						0.123 (0.101)
<i>age</i>			1.061 (0.296)	-0.061 (0.296)	-0.555 (0.065)	
$\text{age}^2/100$			-1.876 (0.516)	1.876 (0.516)	1.313 (0.116)	
<i>F</i>		17.51	8.22	1581	1112	13.87

Of particular interest is the equation for the endogenous regressor (*education*), and the coefficients for the excluded instruments – in this case *college*. The estimated coefficient equals 0.346

with a small standard error. This implies that growing up near a 4-year college increases average educational attainment by 0.3 years. This seems to be a reasonable magnitude.

Since the structural equation is just-identified with one right-hand-side endogenous variable, we can calculate the ILS/IV estimate for the education coefficient as the ratio of the coefficient estimates for the instrument *college* in the two equations, e.g.  $0.346/0.047 = 0.135$ , implying a 13% return to each year of education. This is substantially greater than the 8% least-squares estimate from the first column of Table 11.1.

The IV estimates of the full equation are reported in the second column of Table 11.1.

Card (1995) also points out that if *education* is endogenous, then so is our measure of *experience*, since it is calculated by subtracting *education* from *age*. He suggests that we can use the variables *age* and *age*<sup>2</sup> as instruments for *experience* and *experience*<sup>2</sup>, as they are clearly exogenous and yet highly correlated with *experience* and *experience*<sup>2</sup>. Notice that this approach treats *experience*<sup>2</sup> as a variable separate from *experience*. Indeed, this is the correct approach.

Following this recommendation we now have three endogenous regressors and three instruments. We present the three reduced form equations for the three endogenous regressors in the third through fifth columns of Table 11.2. It is interesting to compare the equations for *education* and *experience*. The two sets of coefficients are simply the sign change of the other, with the exception of the coefficient on *age*. Indeed this must be the case, because the three variables are linearly related. Does this cause a problem for 2SLS? Fortunately, no. The fact that the coefficient on *age* is not simply a sign change means that the equations are not linearly singular. Hence Assumption (11.27) is not violated.

The IV estimates using the three instruments *college*, *age* and *age*<sup>2</sup> for the endogenous regressors *education*, *experience* and *experience*<sup>2</sup> is presented in the third column of Table 11.1. The estimate of the returns to schooling is not affected by this change in the instrument set, but the estimated return to experience profile flattens (the quadratic effect diminishes).

The IV estimator may be calculated in Stata using the `ivregress 2sls` command.

## 11.9 Demeaned Representation

Does the well-known demeaned representation for linear regression (3.20) carry over to the IV estimator? To see this, write the linear projection equation in the format

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + \alpha + e_i$$

where  $\alpha$  is the intercept and  $\mathbf{x}_i$  does not contain a constant. Similarly, partition the instrument as  $(1, \mathbf{z}_i)$  where  $\mathbf{z}_i$  does not contain an intercept. We can write the IV estimates as

$$y_i = \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{\text{iv}} + \hat{\alpha}_{\text{iv}} + \hat{e}_i$$

The orthogonality (11.30) implies the two-equation system

$$\begin{aligned} \sum_{i=1}^n \left( y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{\text{iv}} - \hat{\alpha}_{\text{iv}} \right) &= 0 \\ \sum_{i=1}^n \mathbf{z}_i \left( y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{\text{iv}} - \hat{\alpha}_{\text{iv}} \right) &= \mathbf{0}. \end{aligned}$$

The first equation implies

$$\hat{\alpha}_{\text{iv}} = \bar{y} - \bar{\mathbf{x}}' \hat{\boldsymbol{\beta}}_{\text{iv}}.$$

Substituting into the second equation

$$\sum_{i=1}^n \mathbf{z}_i \left( (y_i - \bar{y}) - (\mathbf{x}_i - \bar{\mathbf{x}})' \hat{\boldsymbol{\beta}}_{\text{iv}} \right)$$

and solving for  $\hat{\beta}_{iv}$  we find

$$\begin{aligned}\hat{\beta}_{iv} &= \left( \sum_{i=1}^n z_i (x_i - \bar{x})' \right)^{-1} \left( \sum_{i=1}^n z_i (y_i - \bar{y}) \right) \\ &= \left( \sum_{i=1}^n (z_i - \bar{z}) (x_i - \bar{x})' \right)^{-1} \left( \sum_{i=1}^n (z_i - \bar{z}) (y_i - \bar{y}) \right).\end{aligned}\quad (11.31)$$

Thus the demeaning equations for least-squares carry over to the IV estimator. The coefficient estimate  $\hat{\beta}_{iv}$  is a function only of the demeaned data.

## 11.10 Wald Estimator

In many cases, including the Card proximity example, the excluded instrument is a binary (dummy) variable. Let's focus on that case, and suppose that the model has just one endogenous regressor and no other regressors beyond the intercept. Thus the model can be written as

$$\begin{aligned}y_i &= x_i \beta + \alpha + e_i \\ \mathbb{E}(e_i \mid z_i) &= 0\end{aligned}$$

with  $z_i$  binary.

Notice that if we take expectations of the structural equation given  $z_i = 1$  and  $z_i = 0$ , respectively, we obtain

$$\begin{aligned}\mathbb{E}(y_i \mid z_i = 1) &= \mathbb{E}(x_i \mid z_i = 1) \beta + \alpha \\ \mathbb{E}(y_i \mid z_i = 0) &= \mathbb{E}(x_i \mid z_i = 0) \beta + \alpha.\end{aligned}$$

Subtracting and dividing, we obtain an expression for the slope coefficient  $\beta$

$$\beta = \frac{\mathbb{E}(y_i \mid z_i = 1) - \mathbb{E}(y_i \mid z_i = 0)}{\mathbb{E}(x_i \mid z_i = 1) - \mathbb{E}(x_i \mid z_i = 0)}.\quad (11.32)$$

The natural moment estimator for  $\beta$  replaces the expectations by the averages within the “grouped data” where  $z_i = 1$  and  $z_i = 0$ , respectively. That is, define the group means

$$\begin{aligned}\bar{y}_1 &= \frac{\sum_{i=1}^n z_i y_i}{\sum_{i=1}^n z_i}, & \bar{y}_0 &= \frac{\sum_{i=1}^n (1 - z_i) y_i}{\sum_{i=1}^n (1 - z_i)} \\ \bar{x}_1 &= \frac{\sum_{i=1}^n z_i x_i}{\sum_{i=1}^n z_i}, & \bar{x}_0 &= \frac{\sum_{i=1}^n (1 - z_i) x_i}{\sum_{i=1}^n (1 - z_i)}\end{aligned}$$

and the moment estimator

$$\hat{\beta} = \frac{\bar{y}_1 - \bar{y}_0}{\bar{x}_1 - \bar{x}_0}.\quad (11.33)$$

This is known as the “Wald estimator” as it was proposed by Wald (1940).

These expressions are rather insightful. (11.32) shows that the structural slope coefficient is the expected change in  $y_i$  due to changing the instrument divided by the expected change in  $x_i$  due to changing the instrument. Informally, it is the change in  $y$  (due to  $z$ ) over the change in  $x$  (due to  $z$ ). Equation (11.33) shows that slope coefficient can be estimated by a simple ratio in means.

The expression (11.33) may appear like a distinct estimator from the IV estimator  $\hat{\beta}_{iv}$ , but it turns out that they are the same. That is,  $\hat{\beta} = \hat{\beta}_{iv}$ . To see this, use (11.31) to find

$$\begin{aligned}\hat{\beta}_{iv} &= \frac{\sum_{i=1}^n z_i (y_i - \bar{y})}{\sum_{i=1}^n z_i (x_i - \bar{x})} \\ &= \frac{\bar{y}_1 - \bar{y}}{\bar{x}_1 - \bar{x}}.\end{aligned}$$

Then notice

$$\bar{y}_1 - \bar{y} = \bar{y}_1 - \left( \frac{1}{n} \sum_{i=1}^n z_i \bar{y}_1 + \frac{1}{n} \sum_{i=1}^n (1 - z_i) \bar{y}_0 \right) = \frac{1}{n} \sum_{i=1}^n (1 - z_i) (\bar{y}_1 - \bar{y}_0)$$

and similarly

$$\bar{x}_1 - \bar{x} = \frac{1}{n} \sum_{i=1}^n (1 - z_i) (\bar{x}_1 - \bar{x}_0)$$

and hence

$$\hat{\beta}_{iv} = \frac{\frac{1}{n} \sum_{i=1}^n (1 - z_i) (\bar{y}_1 - \bar{y}_0)}{\frac{1}{n} \sum_{i=1}^n (1 - z_i) (\bar{x}_1 - \bar{x}_0)} = \hat{\beta}$$

as defined in (11.33). Thus the Wald estimator equals the IV estimator.

We can illustrate using the Card proximity example. If we estimate a simple IV model with no covariates we obtain the estimate  $\hat{\beta}_{iv} = 0.19$ . If we estimate the group-mean log wages and education levels based on the instrument *college*, we find

	near college	not near college
log(wage)	6.311	6.156
education	13.527	12.698

Based on these estimates the Wald estimator of the slope coefficient is  $(6.311 - 6.156) / (13.527 - 12.698) = 0.19$ , the same as the IV estimator.

## 11.11 Two-Stage Least Squares

The IV estimator described in the previous section presumed  $\ell = k$ . Now we allow the general case of  $\ell \geq k$ . Examining the reduced-form equation (11.16) we see

$$y_i = \mathbf{z}_i' \mathbf{\Gamma} \boldsymbol{\beta} + v_i$$

$$\mathbb{E}(\mathbf{z}_i v_i) = \mathbf{0}.$$

Defining  $\mathbf{w}_i = \mathbf{\Gamma}' \mathbf{z}_i$  we can write this as

$$y_i = \mathbf{w}_i' \boldsymbol{\beta} + v_i$$

$$\mathbb{E}(\mathbf{w}_i v_i) = \mathbf{0}.$$

Suppose that  $\mathbf{\Gamma}$  were known. Then we would estimate  $\boldsymbol{\beta}$  by least-squares of  $y_i$  on  $\mathbf{w}_i = \mathbf{\Gamma}' \mathbf{z}_i$

$$\hat{\boldsymbol{\beta}} = (\mathbf{W}' \mathbf{W})^{-1} (\mathbf{W}' \mathbf{y})$$

$$= (\mathbf{\Gamma}' \mathbf{Z}' \mathbf{Z} \mathbf{\Gamma})^{-1} (\mathbf{\Gamma}' \mathbf{Z}' \mathbf{y}).$$

While this is infeasible, we can estimate  $\mathbf{\Gamma}$  from the reduced form regression. Replacing  $\mathbf{\Gamma}$  with its estimate  $\hat{\mathbf{\Gamma}} = (\mathbf{Z}' \mathbf{Z})^{-1} (\mathbf{Z}' \mathbf{X})$  we obtain

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{2sls} &= (\hat{\mathbf{\Gamma}}' \mathbf{Z}' \mathbf{Z} \hat{\mathbf{\Gamma}})^{-1} (\hat{\mathbf{\Gamma}}' \mathbf{Z}' \mathbf{y}) \\ &= \left( \mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{y} \\ &= \left( \mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{y}. \end{aligned} \tag{11.34}$$

This is called the **two-stage-least squares** (2SLS) estimator. It was originally proposed by Theil (1953) and Basman (1957), and is a standard estimator for linear equations with instruments.

If the model is just-identified, so that  $k = \ell$ , then 2SLS simplifies to the IV estimator of the previous section. Since the matrices  $\mathbf{X}'\mathbf{Z}$  and  $\mathbf{Z}'\mathbf{X}$  are square, we can factor

$$\begin{aligned} \left( \mathbf{X}'\mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{X} \right)^{-1} &= (\mathbf{Z}'\mathbf{X})^{-1} \left( (\mathbf{Z}'\mathbf{Z})^{-1} \right)^{-1} (\mathbf{X}'\mathbf{Z})^{-1} \\ &= (\mathbf{Z}'\mathbf{X})^{-1} (\mathbf{Z}'\mathbf{Z}) (\mathbf{X}'\mathbf{Z})^{-1}. \end{aligned}$$

(Once again, this only works when  $k = \ell$ .) Then

$$\begin{aligned} \hat{\beta}_{2\text{sls}} &= \left( \mathbf{X}'\mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{X} \right)^{-1} \mathbf{X}'\mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{y} \\ &= (\mathbf{Z}'\mathbf{X})^{-1} (\mathbf{Z}'\mathbf{Z}) (\mathbf{X}'\mathbf{Z})^{-1} \mathbf{X}'\mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{y} \\ &= (\mathbf{Z}'\mathbf{X})^{-1} (\mathbf{Z}'\mathbf{Z}) (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{y} \\ &= (\mathbf{Z}'\mathbf{X})^{-1} \mathbf{Z}'\mathbf{y} \\ &= \hat{\beta}_{\text{iv}} \end{aligned}$$

as claimed. This shows that the 2SLS estimator as defined in (11.34) is a generalization of the IV estimator defined in (11.29).

There are several alternative representations of the 2SLS estimator which we now describe. First, defining the projection matrix

$$\mathbf{P}_Z = \mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}' \quad (11.35)$$

we can write the 2SLS estimator more compactly as

$$\hat{\beta}_{2\text{sls}} = (\mathbf{X}'\mathbf{P}_Z\mathbf{X})^{-1} \mathbf{X}'\mathbf{P}_Z\mathbf{y}. \quad (11.36)$$

This is useful for representation and derivations, but is not useful for computation as the  $n \times n$  matrix  $\mathbf{P}_Z$  is too large to compute when  $n$  is large.

Second, define the fitted values for  $\mathbf{X}$  from the reduced form

$$\widehat{\mathbf{X}} = \mathbf{P}_Z\mathbf{X} = \mathbf{Z}\widehat{\Gamma}.$$

Then the 2SLS estimator can be written as

$$\hat{\beta}_{2\text{sls}} = \left( \widehat{\mathbf{X}}'\widehat{\mathbf{X}} \right)^{-1} \widehat{\mathbf{X}}'\mathbf{y}.$$

This is an IV estimator as defined in the previous section using  $\widehat{\mathbf{X}}$  as the instrument.

Third, since  $\mathbf{P}_Z$  is idempotent, we can also write the 2SLS estimator as

$$\begin{aligned} \hat{\beta}_{2\text{sls}} &= (\mathbf{X}'\mathbf{P}_Z\mathbf{P}_Z\mathbf{X})^{-1} \mathbf{X}'\mathbf{P}_Z\mathbf{y} \\ &= \left( \widehat{\mathbf{X}}'\widehat{\mathbf{X}} \right)^{-1} \widehat{\mathbf{X}}'\mathbf{y} \end{aligned}$$

which is the least-squares estimator obtained by regressing  $\mathbf{y}$  on the fitted values  $\widehat{\mathbf{X}}$ .

This is the source of the “two-stage” name is since it can be computed as follows.

- First regress  $\mathbf{X}$  on  $\mathbf{Z}$ , vis.,  $\widehat{\Gamma} = (\mathbf{Z}'\mathbf{Z})^{-1} (\mathbf{Z}'\mathbf{X})$  and  $\widehat{\mathbf{X}} = \mathbf{Z}\widehat{\Gamma} = \mathbf{P}_Z\mathbf{X}$ .
- Second, regress  $\mathbf{y}$  on  $\widehat{\mathbf{X}}$ , vis.,  $\hat{\beta}_{2\text{sls}} = \left( \widehat{\mathbf{X}}'\widehat{\mathbf{X}} \right)^{-1} \widehat{\mathbf{X}}'\mathbf{y}$ .

It is useful to scrutinize the projection  $\widehat{\mathbf{X}}$ . Recall,  $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2]$  and  $\mathbf{Z} = [\mathbf{X}_1, \mathbf{Z}_2]$ . Notice  $\widehat{\mathbf{X}}_1 = \mathbf{P}_Z \mathbf{X}_1 = \mathbf{X}_1$  since  $\mathbf{X}_1$  lies in the span of  $\mathbf{Z}$ . Then

$$\widehat{\mathbf{X}} = [\widehat{\mathbf{X}}_1, \widehat{\mathbf{X}}_2] = [\mathbf{X}_1, \widehat{\mathbf{X}}_2].$$

Thus in the second stage, we regress  $\mathbf{y}$  on  $\mathbf{X}_1$  and  $\widehat{\mathbf{X}}_2$ . So only the endogenous variables  $\mathbf{X}_2$  are replaced by their fitted values:

$$\widehat{\mathbf{X}}_2 = \mathbf{X}_1 \widehat{\Gamma}_{12} + \mathbf{Z}_2 \widehat{\Gamma}_{22}.$$

This least squares estimator can be written as

$$\mathbf{y} = \mathbf{X}_1 \widehat{\beta}_1 + \widehat{\mathbf{X}}_2 \widehat{\beta}_2 + \widehat{\varepsilon}.$$

A fourth representation of 2SLS can be obtained from the previous representation for  $\widehat{\beta}_2$ . Set  $\mathbf{P}_1 = \mathbf{X}_1 (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1'$ . Applying the FWL theorem we obtain

$$\begin{aligned} \widehat{\beta}_2 &= \left( \widehat{\mathbf{X}}_2' (\mathbf{I}_n - \mathbf{P}_1) \widehat{\mathbf{X}}_2 \right)^{-1} \left( \widehat{\mathbf{X}}_2' (\mathbf{I}_n - \mathbf{P}_1) \mathbf{y} \right) \\ &= \left( \mathbf{X}_2' \mathbf{P}_Z (\mathbf{I}_n - \mathbf{P}_1) \mathbf{P}_Z \mathbf{X}_2 \right)^{-1} \left( \mathbf{X}_2' \mathbf{P}_Z (\mathbf{I}_n - \mathbf{P}_1) \mathbf{y} \right) \\ &= \left( \mathbf{X}_2' (\mathbf{P}_Z - \mathbf{P}_1) \mathbf{X}_2 \right)^{-1} \left( \mathbf{X}_2' (\mathbf{P}_Z - \mathbf{P}_1) \mathbf{y} \right) \end{aligned}$$

since  $\mathbf{P}_Z \mathbf{P}_1 = \mathbf{P}_1$ .

A fifth representation can be obtained by a further projection. The projection matrix  $\mathbf{P}_Z$  can be replaced by the projection onto the pair  $[\mathbf{X}_1, \widetilde{\mathbf{Z}}_2]$  where  $\widetilde{\mathbf{Z}}_2 = (\mathbf{I}_n - \mathbf{P}_1) \mathbf{Z}_2$  is  $\mathbf{Z}_2$  projected orthogonal to  $\mathbf{X}_1$ . Since  $\mathbf{X}_1$  and  $\widetilde{\mathbf{Z}}_2$  are orthogonal,  $\mathbf{P}_Z = \mathbf{P}_1 + \mathbf{P}_2$  where  $\mathbf{P}_2 = \widetilde{\mathbf{Z}}_2 \left( \widetilde{\mathbf{Z}}_2' \widetilde{\mathbf{Z}}_2 \right)^{-1} \widetilde{\mathbf{Z}}_2'$ . Thus  $\mathbf{P}_Z - \mathbf{P}_1 = \mathbf{P}_2$  and

$$\begin{aligned} \widehat{\beta}_2 &= \left( \mathbf{X}_2' \mathbf{P}_2 \mathbf{X}_2 \right)^{-1} \left( \mathbf{X}_2' \mathbf{P}_2 \mathbf{y} \right) \\ &= \left( \mathbf{X}_2' \widetilde{\mathbf{Z}}_2 \left( \widetilde{\mathbf{Z}}_2' \widetilde{\mathbf{Z}}_2 \right)^{-1} \widetilde{\mathbf{Z}}_2' \mathbf{X}_2 \right)^{-1} \left( \mathbf{X}_2' \widetilde{\mathbf{Z}}_2 \left( \widetilde{\mathbf{Z}}_2' \widetilde{\mathbf{Z}}_2 \right)^{-1} \widetilde{\mathbf{Z}}_2' \mathbf{y} \right) \end{aligned} \quad (11.37)$$

Given the 2SLS estimator we define the residual vector

$$\widehat{\mathbf{e}} = \mathbf{y} - \mathbf{X} \widehat{\beta}_{2\text{sls}}.$$

When the model is overidentified, the instruments and residuals are not orthogonal. That is

$$\mathbf{Z}' \widehat{\mathbf{e}} \neq \mathbf{0}.$$

It does, however, satisfy

$$\begin{aligned} \widehat{\mathbf{X}}' \widehat{\mathbf{e}} &= \widehat{\Gamma}' \mathbf{Z}' \widehat{\mathbf{e}} \\ &= \mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \widehat{\mathbf{e}} \\ &= \mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{y} - \mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X} \widehat{\beta}_{2\text{sls}} \\ &= \mathbf{0}. \end{aligned}$$

Returning to Card's college proximity example, suppose that we treat experience as exogenous, but that instead of using the single instrument *college* (grew up near a 4-year college) we use the two instruments (*public*, *private*) (grew up near a public/private 4-year college, respectively). In this case we have one endogenous variable (*education*) and two instruments (*public*, *private*). The estimated reduced form equation for *education* is presented in the sixth column of Table 11.2. In this specification, the coefficient on *public* – growing up near a public 4-year college – is larger



than that found for the variable *college* in the previous specification (column 2). Furthermore, the coefficient on *private* – growing up near a private 4-year college – is much smaller. This indicates that the key impact of proximity on education is via public colleges rather than private colleges.

The 2SLS estimates obtained using these two instruments are presented in the fourth column of Table 11.1. The coefficient on *education* increases to 0.162, indicating a 16% return to a year of education. This is roughly twice as large as the estimate obtained by least-squares in the first column.

Additionally, if we follow Card and treat *experience* as endogenous and use *age* as an instrument, we now have three endogenous variables (*education*, *experience*, *experience*<sup>2</sup>/100) and four instruments (*public*, *private*, *age*, *age*<sup>2</sup>). We present the 2SLS estimates using this specification in the fifth column of Table 11.1. The estimate of the return to education remains about 16%, but again the return to experience flattens.

You might wonder if we could use all three instruments – *college*, *public*, and *private*. The answer is no. This is because *college* = *public* + *private* so the three variables are colinear. Since the instruments are linearly related, the three together would violate the full-rank condition (11.7).

The 2SLS estimator may be calculated in Stata using the `ivregress 2sls` command.

## 11.12 Limited Information Maximum Likelihood

An alternative method to estimate the parameters of the structural equation is by maximum likelihood. Anderson and Rubin (1949) derived the maximum likelihood estimator for the joint distribution of  $(y_i, \mathbf{x}_{2i})$ . The estimator is known as **limited information maximum likelihood**, or LIML.

This estimator is called “limited information” because it is based on the structural equation for  $y_i$  combined with the reduced form equation for  $\mathbf{x}_{2i}$ . If maximum likelihood is derived based on a structural equation for  $\mathbf{x}_{2i}$  as well, then this leads to what is known as **full information maximum likelihood** (FIML). The advantage of the LIML approach relative to FIML is that the former does not require a structural model for  $\mathbf{x}_{2i}$ , and thus allows the researcher to focus on the structural equation of interest – that for  $y_i$ . We do not describe the FIML estimator here as it is not commonly used in applied econometric practice.

While the LIML estimator is less widely used among economists than 2SLS, it has received a resurgence of attention from econometric theorists.

To derive the LIML estimator, start by writing the joint reduced form equations (11.19) and (11.15) as

$$\begin{aligned} \mathbf{w}_i &= \begin{pmatrix} y_i \\ \mathbf{x}_{2i} \end{pmatrix} \\ &= \begin{bmatrix} \boldsymbol{\lambda}'_1 & \boldsymbol{\lambda}'_2 \\ \boldsymbol{\Gamma}'_{12} & \boldsymbol{\Gamma}'_{22} \end{bmatrix} \begin{pmatrix} \mathbf{z}_{1i} \\ \mathbf{z}_{2i} \end{pmatrix} + \begin{pmatrix} v_i \\ \mathbf{u}_{2i} \end{pmatrix} \\ &= \boldsymbol{\Pi}'_1 \mathbf{z}_{1i} + \boldsymbol{\Pi}'_2 \mathbf{z}_{2i} + \boldsymbol{\xi}_i \end{aligned} \tag{11.38}$$

where  $\boldsymbol{\Pi}_1 = [\boldsymbol{\lambda}_1 \quad \boldsymbol{\Gamma}_{12}]$ ,  $\boldsymbol{\Pi}_2 = [\boldsymbol{\lambda}_2 \quad \boldsymbol{\Gamma}_{22}]$  and  $\boldsymbol{\xi}'_i = [v_i \quad \mathbf{u}'_{2i}]$ . The LIML estimator is derived under the assumption that  $\boldsymbol{\xi}_i$  is multivariate normal.

Define  $\boldsymbol{\gamma}' = [1 \quad -\boldsymbol{\beta}'_2]$ . From (11.21) we find

$$\boldsymbol{\Pi}_2 \boldsymbol{\gamma} = \boldsymbol{\lambda}_2 - \boldsymbol{\Gamma}_{22} \boldsymbol{\beta}_2 = \mathbf{0}.$$

Thus the  $\ell_2 \times (k_2 + 1)$  coefficient matrix  $\boldsymbol{\Pi}_2$  in (11.38) has deficient rank. Indeed, its rank must be  $k_2$ , since  $\boldsymbol{\Gamma}_{22}$  has full rank.

This means that the model (11.38) is precisely the reduced rank regression model of Section 10.9. Theorem 10.9.1 presents the maximum likelihood estimators for the reduced rank parameters.

In particular, the MLE for  $\gamma$  is

$$\hat{\gamma} = \underset{\gamma}{\operatorname{argmin}} \frac{\gamma' \mathbf{W}' \mathbf{M}_1 \mathbf{W} \gamma}{\gamma' \mathbf{W}' \mathbf{M}_Z \mathbf{W} \gamma} \quad (11.39)$$

where  $\mathbf{W}$  is the  $n \times (1 + k_2)$  matrix of the stacked  $\mathbf{w}'_i = (y_i \quad \mathbf{x}'_{2i})$ ,  $\mathbf{M}_1 = \mathbf{I}_n - \mathbf{Z}_1 (\mathbf{Z}'_1 \mathbf{Z}_1)^{-1} \mathbf{Z}'_1$  and  $\mathbf{M}_Z = \mathbf{I}_n - \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}'$ . The minimization (11.39) is sometimes called the “least variance ratio” problem.

The minimization problem (11.39) is invariant to the scale of  $\gamma$  (that is,  $\hat{\gamma}c$  is equivalently the argmin for any  $c$ ) so a normalization is required. For estimation of the structural parameters a convenient normalization is  $\gamma' = [1 \quad -\beta'_2]$ . Another is to set  $\gamma' \mathbf{W}' \mathbf{M}_Z \mathbf{W} \gamma = 1$ . In this case, from the theory of the minimum of quadratic forms (Section A.11),  $\hat{\gamma}$  is the generalized eigenvector of  $\mathbf{W}' \mathbf{M}_1 \mathbf{W}$  with respect to  $\mathbf{W}' \mathbf{M}_Z \mathbf{W}$  associated with the smallest generalized eigenvalue. (See Section A.10 for the definition of generalized eigenvalues and eigenvectors.) Computationally this is straightforward. For example, in MATLAB, the generalized eigenvalues and eigenvectors of the matrix  $\mathbf{A}$  with respect to  $\mathbf{B}$  is found by the command `eig(A,B)`. Once  $\hat{\gamma}$  is found, to obtain the MLE for  $\beta_2$ , make the partition  $\hat{\gamma}' = [\hat{\gamma}_1 \quad \hat{\gamma}'_2]$  and set  $\hat{\beta}_2 = -\hat{\gamma}_2/\hat{\gamma}_1$ .

To obtain the MLE for  $\beta_1$ , recall the structural equation  $y_i = \mathbf{x}'_{1i} \beta_1 + \mathbf{x}'_{2i} \beta_2 + e_i$ . Replacing  $\beta_2$  with the MLE  $\hat{\beta}_2$  and then applying regression we obtain the MLE for  $\beta_1$ . Thus

$$\hat{\beta}_1 = (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 (\mathbf{Y} - \mathbf{X}_2 \hat{\beta}_2). \quad (11.40)$$

These solutions are the MLE (known as the LIML estimator) for the structural parameters  $\beta_1$  and  $\beta_2$ .

Many previous econometrics textbooks do not present a derivation of the LIML estimator as the original derivation by Anderson and Rubin (1949) is lengthy and not particularly insightful. In contrast, the derivation given here based on reduced rank regression is relatively simple.

There is an alternative (and traditional) expression for the LIML estimator. Define the minimum obtained in (11.39)

$$\hat{\kappa} = \min_{\gamma} \frac{\gamma' \mathbf{W}' \mathbf{M}_1 \mathbf{W} \gamma}{\gamma' \mathbf{W}' \mathbf{M}_Z \mathbf{W} \gamma} \quad (11.41)$$

which is the smallest generalized eigenvalue of  $\mathbf{W}' \mathbf{M}_1 \mathbf{W}$  with respect to  $\mathbf{W}' \mathbf{M}_Z \mathbf{W}$ . The LIML estimator then can be written as

$$\hat{\beta}_{\text{liml}} = (\mathbf{X}' (\mathbf{I}_n - \hat{\kappa} \mathbf{M}_Z) \mathbf{X})^{-1} (\mathbf{X}' (\mathbf{I}_n - \hat{\kappa} \mathbf{M}_Z) \mathbf{y}). \quad (11.42)$$

We defer the derivation of (11.42) until the end of this section. Expression (11.42) does not simplify the computation (since  $\hat{\kappa}$  requires solving the same eigenvector problem that yields  $\hat{\beta}_2$ ). However (11.42) is important for the distribution theory of the LIML estimator, and to reveal the algebraic connection between LIML, least-squares, and 2SLS.

The estimator class (11.42) with arbitrary  $\kappa$  is known as a  $k$  class estimator of  $\beta$ . While the LIML estimator obtains by setting  $\kappa = \hat{\kappa}$ , the least-squares estimator is obtained by setting  $\kappa = 0$  and 2SLS is obtained by setting  $\kappa = 1$ . It is worth observing that the LIML solution to (11.41) satisfies  $\hat{\kappa} \geq 1$ .

When the model is just-identified, the LIML estimator is identical to the IV and 2SLS estimators. They are only different in the over-identified setting. (One corollary is that under just-identification the IV estimator is MLE under normality.)

For inference, it is useful to observe that (11.42) shows that  $\hat{\beta}_{\text{liml}}$  can be written as an IV estimator

$$\hat{\beta}_{\text{liml}} = (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} (\tilde{\mathbf{X}}' \mathbf{y})$$

using the instrument

$$\tilde{\mathbf{X}} = (\mathbf{I}_n - \hat{\kappa} \mathbf{M}_Z) \mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 - \hat{\kappa} \hat{\mathbf{U}}_2 \end{pmatrix}$$

where  $\widehat{\mathbf{U}}_2 = \mathbf{M}_Z \mathbf{X}_2$  are the (reduced-form) residuals from the multivariate regression of the endogenous regressors  $\mathbf{x}_{2i}$  on the instruments  $\mathbf{z}_i$ . Expressing LIML using this IV formula is useful for variance estimation.

Asymptotically the LIML estimator has the same distribution as 2SLS. However, they can have quite different behaviors in finite samples. There is considerable evidence that the LIML estimator has superior finite sample performance to 2SLS when there are many instruments or the reduced form is weak. (We review these cases in the following sections.) However, on the other hand there is worry that since the LIML estimator is derived under normality it may not be robust in non-normal settings.

We now derive the expression (11.42). Use the normaliaation  $\gamma' = [1 \quad -\beta_2']$  to write (11.39) as

$$\widehat{\beta}_2 = \underset{\beta_2}{\operatorname{argmin}} \frac{(\mathbf{Y} - \mathbf{X}_2 \beta_2)' \mathbf{M}_1 (\mathbf{Y} - \mathbf{X}_2 \beta_2)}{(\mathbf{Y} - \mathbf{X}_2 \beta_2)' \mathbf{M}_Z (\mathbf{Y} - \mathbf{X}_2 \beta_2)}$$

The first-order-condition for minimization

$$2 \frac{\mathbf{X}_2' \mathbf{M}_1 (\mathbf{Y} - \mathbf{X}_2 \widehat{\beta}_2)}{(\mathbf{Y} - \mathbf{X}_2 \widehat{\beta}_2)' \mathbf{M}_Z (\mathbf{Y} - \mathbf{X}_2 \widehat{\beta}_2)} - 2 \frac{(\mathbf{Y} - \mathbf{X}_2 \widehat{\beta}_2)' \mathbf{M}_1 (\mathbf{Y} - \mathbf{X}_2 \widehat{\beta}_2)}{(\mathbf{Y} - \mathbf{X}_2 \widehat{\beta}_2)' \mathbf{M}_Z (\mathbf{Y} - \mathbf{X}_2 \widehat{\beta}_2)^2} \mathbf{X}_2' \mathbf{M}_Z (\mathbf{Y} - \mathbf{X}_2 \widehat{\beta}_2) = 0.$$

Multiplying by  $(\mathbf{Y} - \mathbf{X}_2 \widehat{\beta}_2)' \mathbf{M}_Z (\mathbf{Y} - \mathbf{X}_2 \widehat{\beta}_2) / 2$  and using definition (11.41) we find

$$\mathbf{X}_2' \mathbf{M}_1 (\mathbf{Y} - \mathbf{X}_2 \widehat{\beta}_2) - \widehat{\kappa} \mathbf{X}_2' \mathbf{M}_Z (\mathbf{Y} - \mathbf{X}_2 \widehat{\beta}_2) = 0.$$

Rewriting,

$$\mathbf{X}_2' (\mathbf{M}_1 - \widehat{\kappa} \mathbf{M}_Z) \mathbf{X}_2 \widehat{\beta}_2 = \mathbf{X}_2' (\mathbf{M}_1 - \widehat{\kappa} \mathbf{M}_Z) \mathbf{y}. \quad (11.43)$$

Equation (11.42) is the same as the two equation system

$$\begin{aligned} \mathbf{X}_1' \mathbf{X}_1 \widehat{\beta}_1 + \mathbf{X}_1' \mathbf{X}_2 \widehat{\beta}_2 &= \mathbf{X}_1' \mathbf{y} \\ \mathbf{X}_2' \mathbf{X}_1 \widehat{\beta}_1 + (\mathbf{X}_2' (\mathbf{I}_n - \widehat{\kappa} \mathbf{M}_Z) \mathbf{X}_2) \widehat{\beta}_2 &= \mathbf{X}_2' (\mathbf{I}_n - \widehat{\kappa} \mathbf{M}_Z) \mathbf{y}. \end{aligned}$$

The first equation is (11.40). Using (11.40), the second is

$$\mathbf{X}_2' \mathbf{X}_1 (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' (\mathbf{Y} - \mathbf{X}_2 \widehat{\beta}_2) + (\mathbf{X}_2' (\mathbf{I}_n - \widehat{\kappa} \mathbf{M}_Z) \mathbf{X}_2) \widehat{\beta}_2 = \mathbf{X}_2' (\mathbf{I}_n - \widehat{\kappa} \mathbf{M}_Z) \mathbf{y}$$

which is (11.43) when rearranged. We have thus shown that (11.42) is equivalent to (11.40) and (11.43) and is thus a valid expression for the LIML estimator.

Returning to the Card college proximity example, we now present the LIML estimates of the equation with the two instruments (*public*, *private*). They are reported in the final column of Table 11.1. They are quite similar to the 2SLS estimates in this application.

The LIML estimator may be calculated in Stata using the `ivregress liml` command.

### Theodore Anderson

Theodore (Ted) Anderson (1918-2016) was a American statistician and econometrician, who made fundamental contributions to multivariate statistical theory. Important contributions include the Anderson-Darling distribution test, the Anderson-Rubin statistic, the method of reduced rank regression, and his most famous econometrics contribution – the LIML estimator. He continued working throughout his long life, even publishing theoretical work at the age of 97!

### 11.13 Consistency of 2SLS

We now present a demonstration of the consistency of the 2SLS estimate for the structural parameter. The following is a set of regularity conditions.

**Assumption 11.13.1**

1. The observations  $(y_i, \mathbf{x}_i, \mathbf{z}_i)$ ,  $i = 1, \dots, n$ , are independent and identically distributed.
2.  $\mathbb{E}(y^2) < \infty$ .
3.  $\mathbb{E}\|\mathbf{x}\|^2 < \infty$ .
4.  $\mathbb{E}\|\mathbf{z}\|^2 < \infty$ .
5.  $\mathbb{E}(\mathbf{z}\mathbf{z}')$  is positive definite.
6.  $\mathbb{E}(\mathbf{z}\mathbf{x}')$  has full rank  $k$ .
7.  $\mathbb{E}(\mathbf{z}\mathbf{e}) = 0$ .

Assumptions 11.13.1.2-4 state that all variables have finite variances. Assumption 11.13.1.5 states that the instrument vector has an invertible design matrix, which is identical to the core assumption about regressors in the linear regression model. This excludes linearly redundant instruments. Assumptions 11.13.1.6 and 11.13.1.7 are the key identification conditions for instrumental variables. Assumption 11.13.1.6 states that the instruments and regressors have a full-rank cross-moment matrix. This is often called the relevance condition. Assumption 11.13.1.7 states that the instrumental variables and structural error are uncorrelated. Assumptions 11.13.1.5-7 are identical to Definition 11.3.1.

**Theorem 11.13.1** Under Assumption 11.13.1,  $\hat{\beta}_{2\text{sls}} \xrightarrow{p} \beta$  as  $n \rightarrow \infty$ .

The proof of the theorem is provided below

This theorem shows that the 2SLS estimator is consistent for the structural coefficient  $\beta$  under similar moment conditions as the least-squares estimator. The key differences are the instrumental variables assumption  $\mathbb{E}(\mathbf{z}\mathbf{e}) = 0$  and the identification assumption  $\text{rank}(\mathbb{E}(\mathbf{z}\mathbf{x}')) = k$ .

The result includes the IV estimator (when  $\ell = k$ ) as a special case.

The proof of this consistency result is similar to that for the least-squares estimator. Take the structural equation  $\mathbf{y} = \mathbf{X}\beta + \mathbf{e}$  in matrix format and substitute it into the expression for the estimator. We obtain

$$\begin{aligned} \hat{\beta}_{2\text{sls}} &= \left( \mathbf{X}'\mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{X} \right)^{-1} \mathbf{X}'\mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'(\mathbf{X}\beta + \mathbf{e}) \\ &= \beta + \left( \mathbf{X}'\mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{X} \right)^{-1} \mathbf{X}'\mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{e}. \end{aligned} \quad (11.44)$$

This separates out the stochastic component. Re-writing and applying the WLLN and CMT

$$\begin{aligned}\widehat{\beta}_{2\text{sls}} - \beta &= \left( \left( \frac{1}{n} \mathbf{X}' \mathbf{Z} \right) \left( \frac{1}{n} \mathbf{Z}' \mathbf{Z} \right)^{-1} \left( \frac{1}{n} \mathbf{Z}' \mathbf{X} \right) \right)^{-1} \\ &\quad \cdot \left( \frac{1}{n} \mathbf{X}' \mathbf{Z} \right) \left( \frac{1}{n} \mathbf{Z}' \mathbf{Z} \right)^{-1} \left( \frac{1}{n} \mathbf{Z}' \mathbf{e} \right) \\ &\xrightarrow{p} (\mathbf{Q}_{xz} \mathbf{Q}_{zz}^{-1} \mathbf{Q}_{zx})^{-1} \mathbf{Q}_{xz} \mathbf{Q}_{zz}^{-1} \mathbb{E}(z_i e_i) = 0\end{aligned}$$

where

$$\begin{aligned}\mathbf{Q}_{xz} &= \mathbb{E}(x_i z_i') \\ \mathbf{Q}_{zz} &= \mathbb{E}(z_i z_i') \\ \mathbf{Q}_{zx} &= \mathbb{E}(z_i x_i').\end{aligned}$$

The WLLN holds under the i.i.d. Assumption 11.13.1.1 and the finite second moment Assumptions 11.13.1.2-4. The continuous mapping theorem applies if the matrices  $\mathbf{Q}_{zz}$  and  $\mathbf{Q}_{xz} \mathbf{Q}_{zz}^{-1} \mathbf{Q}_{zx}$  are invertible, which hold under the identification Assumptions 11.13.1.5 and 11.13.1.6. The final equality uses Assumption 11.13.1.7.

## 11.14 Asymptotic Distribution of 2SLS

We now show that the 2SLS estimator satisfies a central limit theorem. We first state a set of sufficient regularity conditions.

**Assumption 11.14.1** *In addition to Assumption 11.13.1,*

1.  $\mathbb{E}(y^4) < \infty$ .
2.  $\mathbb{E}\|z\|^4 < \infty$ .

Assumption 11.14.1 strengthens Assumption 11.13.1 by requiring that the dependent variable and instruments have finite fourth moments. This is used to establish the central limit theorem.

**Theorem 11.14.1** *Under Assumption 11.14.1, as  $n \rightarrow \infty$ .*

$$\sqrt{n}(\widehat{\beta}_{2\text{sls}} - \beta) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_\beta)$$

where

$$\mathbf{V}_\beta = (\mathbf{Q}_{xz} \mathbf{Q}_{zz}^{-1} \mathbf{Q}_{zx})^{-1} (\mathbf{Q}_{xz} \mathbf{Q}_{zz}^{-1} \mathbf{\Omega} \mathbf{Q}_{zz}^{-1} \mathbf{Q}_{zx}) (\mathbf{Q}_{xz} \mathbf{Q}_{zz}^{-1} \mathbf{Q}_{zx})^{-1}$$

and

$$\mathbf{\Omega} = \mathbb{E}(z_i z_i' e_i^2).$$

This shows that the 2SLS estimator converges at a  $\sqrt{n}$  rate to a normal random vector. It shows as well the form of the covariance matrix. The latter takes a substantially more complicated form than the least-squares estimator.

As in the case of least-squares estimation, the asymptotic variance simplifies under a conditional homoskedasticity condition. For 2SLS the simplification occurs when  $\mathbb{E}(e_i^2 | \mathbf{z}_i) = \sigma^2$ . This holds when  $\mathbf{z}_i$  and  $e_i$  are independent. It may be reasonable in some contexts to conceive that the error  $e_i$  is independent of the excluded instruments  $\mathbf{z}_{2i}$ , since by assumption the impact of  $\mathbf{z}_{2i}$  on  $y_i$  is only through  $\mathbf{x}_i$ , but there is no reason to expect  $e_i$  to be independent of the included exogenous variables  $\mathbf{x}_{1i}$ . Hence heteroskedasticity should be equally expected in 2SLS and least-squares regression. Nevertheless, under the homoskedasticity condition then we have the simplifications  $\mathbf{\Omega} = \mathbf{Q}_{zz}\sigma^2$  and  $\mathbf{V}_\beta = \mathbf{V}_\beta^0 \stackrel{\text{def}}{=} (\mathbf{Q}_{xz}\mathbf{Q}_{zz}^{-1}\mathbf{Q}_{zx})^{-1}\sigma^2$ .

The derivation of the asymptotic distribution builds on the proof of consistency. Using equation (11.44) we have

$$\begin{aligned} \sqrt{n}(\hat{\beta}_{2\text{sls}} - \beta) &= \left( \left( \frac{1}{n} \mathbf{X}' \mathbf{Z} \right) \left( \frac{1}{n} \mathbf{Z}' \mathbf{Z} \right)^{-1} \left( \frac{1}{n} \mathbf{Z}' \mathbf{X} \right) \right)^{-1} \\ &\quad \cdot \left( \frac{1}{n} \mathbf{X}' \mathbf{Z} \right) \left( \frac{1}{n} \mathbf{Z}' \mathbf{Z} \right)^{-1} \left( \frac{1}{\sqrt{n}} \mathbf{Z}' \mathbf{e} \right). \end{aligned}$$

We apply the WLLN and CMT for the moment matrices involving  $\mathbf{X}$  and  $\mathbf{Z}$  the same as in the proof of consistency. In addition, by the CLT for i.i.d. observations

$$\frac{1}{\sqrt{n}} \mathbf{Z}' \mathbf{e} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{z}_i e_i \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{\Omega})$$

because the vector  $\mathbf{z}_i e_i$  is i.i.d. and mean zero under Assumptions 11.13.1.1 and 11.13.1.7, and has a finite second moment as we verify below.

We obtain

$$\begin{aligned} \sqrt{n}(\hat{\beta}_{2\text{sls}} - \beta) &= \left( \left( \frac{1}{n} \mathbf{X}' \mathbf{Z} \right) \left( \frac{1}{n} \mathbf{Z}' \mathbf{Z} \right)^{-1} \left( \frac{1}{n} \mathbf{Z}' \mathbf{X} \right) \right)^{-1} \\ &\quad \cdot \left( \frac{1}{n} \mathbf{X}' \mathbf{Z} \right) \left( \frac{1}{n} \mathbf{Z}' \mathbf{Z} \right)^{-1} \left( \frac{1}{\sqrt{n}} \mathbf{Z}' \mathbf{e} \right) \\ &\xrightarrow{d} (\mathbf{Q}_{xz}\mathbf{Q}_{zz}^{-1}\mathbf{Q}_{zx})^{-1} \mathbf{Q}_{xz}\mathbf{Q}_{zz}^{-1} \mathbf{N}(\mathbf{0}, \mathbf{\Omega}) = \mathbf{N}(\mathbf{0}, \mathbf{V}_\beta) \end{aligned}$$

as stated.

For completeness, we demonstrate that  $\mathbf{z}_i e_i$  has a finite second moment under Assumption 11.14.1. To see this, note that by Minkowski's inequality

$$\begin{aligned} (\mathbb{E}(e^4))^{1/4} &= \left( \mathbb{E}((y - \mathbf{x}'\beta)^4) \right)^{1/4} \\ &\leq (\mathbb{E}(y^4))^{1/4} + \|\beta\| \left( \mathbb{E}\|\mathbf{x}\|^4 \right)^{1/4} < \infty \end{aligned}$$

under Assumptions 11.14.1.1 and 11.14.1.2. Then by the Cauchy-Schwarz inequality

$$\mathbb{E}\|\mathbf{z}e\|^2 \leq \left( \mathbb{E}\|\mathbf{z}\|^4 \right)^{1/2} (\mathbb{E}(e^4))^{1/2} < \infty$$

using Assumptions 11.14.1.3.

### 11.15 Determinants of 2SLS Variance

It is instructive to examine the asymptotic variance of the 2SLS estimator to understand the factors which determine the precision (or lack thereof) of the estimator. As in the least-squares case, it is more transparent to examine the variance under the assumption of homoskedasticity. In this case the asymptotic variance takes the form

$$\begin{aligned}\mathbf{V}_\beta^0 &= (\mathbf{Q}_{xz} \mathbf{Q}_{zz}^{-1} \mathbf{Q}_{zx})^{-1} \sigma^2 \\ &= \left( \mathbb{E}(\mathbf{x}_i \mathbf{z}_i') (\mathbb{E}(\mathbf{z}_i \mathbf{z}_i'))^{-1} \mathbb{E}(\mathbf{z}_i \mathbf{x}_i') \right)^{-1} \mathbb{E}(e_i^2).\end{aligned}$$

As in the least-squares case, we can see that the variance is increasing in the variance of the error  $e_i$ , and decreasing in the variance of  $\mathbf{x}_i$ . What is different is that the variance is decreasing in the (matrix-valued) correlation between  $\mathbf{x}_i$  and  $\mathbf{z}_i$ .

It is also useful to observe that the variance expression is not affected by the variance structure of  $\mathbf{z}_i$ . Indeed,  $\mathbf{V}_\beta^0$  is invariant to rotations of  $\mathbf{z}_i$  (if you replace  $\mathbf{z}_i$  with  $\mathbf{C}\mathbf{z}_i$  for invertible  $\mathbf{C}$  the expression does not change). This means that the variance expression is not affected by the scaling of  $\mathbf{z}_i$ , and is not directly affected by correlation among the  $\mathbf{z}_i$ .

We can also use this expression to examine the impact of increasing the instrument set. Suppose we partition  $\mathbf{z}_i = (\mathbf{z}_{ai}, \mathbf{z}_{bi})$  where  $\dim(\mathbf{z}_{ai}) \geq k$  so we can construct the 2SLS estimator using  $\mathbf{z}_{ai}$ . Let  $\hat{\beta}_a$  and  $\hat{\beta}$  denote the 2SLS estimators constructed using the instrument sets  $\mathbf{z}_{ai}$  and  $(\mathbf{z}_{ai}, \mathbf{z}_{bi})$ , respectively. Without loss of generality we can assume that  $\mathbf{z}_{ai}$  and  $\mathbf{z}_{bi}$  are uncorrelated (if not, replace  $\mathbf{z}_{bi}$  with the projection error after projecting onto  $\mathbf{z}_{ai}$ ). In this case both  $\mathbb{E}(\mathbf{z}_i \mathbf{z}_i')$  and  $(\mathbb{E}(\mathbf{z}_i \mathbf{z}_i'))^{-1}$  are block diagonal, so

$$\begin{aligned}\text{avar}(\hat{\beta}) &= \left( \mathbb{E}(\mathbf{x}_i \mathbf{z}_i') (\mathbb{E}(\mathbf{z}_i \mathbf{z}_i'))^{-1} \mathbb{E}(\mathbf{z}_i \mathbf{x}_i') \right)^{-1} \sigma^2 \\ &= \left( \mathbb{E}(\mathbf{x}_i \mathbf{z}_{ai}') (\mathbb{E}(\mathbf{z}_{ai} \mathbf{z}_{ai}'))^{-1} \mathbb{E}(\mathbf{z}_{ai} \mathbf{x}_i') + \mathbb{E}(\mathbf{x}_i \mathbf{z}_{bi}') (\mathbb{E}(\mathbf{z}_{bi} \mathbf{z}_{bi}'))^{-1} \mathbb{E}(\mathbf{z}_{bi} \mathbf{x}_i') \right)^{-1} \sigma^2 \\ &\leq \left( \mathbb{E}(\mathbf{x}_i \mathbf{z}_{ai}') (\mathbb{E}(\mathbf{z}_{ai} \mathbf{z}_{ai}'))^{-1} \mathbb{E}(\mathbf{z}_{ai} \mathbf{x}_i') \right)^{-1} \sigma^2 \\ &= \text{avar}(\hat{\beta}_a)\end{aligned}$$

with strict inequality if  $\mathbb{E}(\mathbf{x}_i \mathbf{z}_{bi}') \neq \mathbf{0}$ . Thus the 2SLS estimator with the full instrument set has a smaller asymptotic variance than the estimator with the smaller instrument set.

What we have shown is that the asymptotic variance of the 2SLS estimator is decreasing as the number of instruments increases. From the viewpoint of asymptotic efficiency, this means that it is better to use more instruments (when they are available and are all known to be valid instruments) rather than less.

Unfortunately, there is always a catch. In this case it turns out that the finite sample bias of the 2SLS estimator (which cannot be calculated exactly, but can be approximated using asymptotic expansions) is generically increasing linearly as the number of instruments increases. We will see some calculations illustrating this phenomenon in Section 11.33. Thus the choice of instruments in practice induces a trade-off between bias and variance.

### 11.16 Covariance Matrix Estimation

Estimation of the asymptotic variance matrix  $\mathbf{V}_\beta$  is done using similar techniques as for least-squares estimation. The estimator is constructed by replacing the population moment matrices by sample counterparts. Thus

$$\hat{\mathbf{V}}_\beta = \left( \hat{\mathbf{Q}}_{xz} \hat{\mathbf{Q}}_{zz}^{-1} \hat{\mathbf{Q}}_{zx} \right)^{-1} \left( \hat{\mathbf{Q}}_{xz} \hat{\mathbf{Q}}_{zz}^{-1} \hat{\Omega} \hat{\mathbf{Q}}_{zz}^{-1} \hat{\mathbf{Q}}_{zx} \right) \left( \hat{\mathbf{Q}}_{xz} \hat{\mathbf{Q}}_{zz}^{-1} \hat{\mathbf{Q}}_{zx} \right)^{-1} \quad (11.45)$$

where

$$\begin{aligned}\widehat{\mathbf{Q}}_{zz} &= \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i' = \frac{1}{n} \mathbf{Z}' \mathbf{Z} \\ \widehat{\mathbf{Q}}_{xz} &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{z}_i' = \frac{1}{n} \mathbf{X}' \mathbf{Z} \\ \widehat{\Omega} &= \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i' \widehat{e}_i^2 \\ \widehat{e}_i &= y_i - \mathbf{x}_i' \widehat{\boldsymbol{\beta}}_{2\text{sls}}.\end{aligned}$$

The homoskedastic variance matrix can be estimated by

$$\begin{aligned}\widehat{\mathbf{V}}_{\boldsymbol{\beta}}^0 &= \left( \widehat{\mathbf{Q}}_{xz} \widehat{\mathbf{Q}}_{zz}^{-1} \widehat{\mathbf{Q}}_{zx} \right)^{-1} \widehat{\sigma}^2 \\ \widehat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n \widehat{e}_i^2.\end{aligned}$$

Standard errors for the coefficients are obtained as the square roots of the diagonal elements of  $n^{-1} \widehat{\mathbf{V}}_{\boldsymbol{\beta}}$ . Confidence intervals, t-tests, and Wald tests may all be constructed from the coefficient estimates and covariance matrix estimate exactly as for least-squares regression.

In Stata, the `ivregress` command by default calculates the covariance matrix estimator using the homoskedastic variance matrix. To obtain covariance matrix estimation and standard errors with the robust estimator  $\widehat{\mathbf{V}}_{\boldsymbol{\beta}}$ , use the “,r” option.

**Theorem 11.16.1** *Under Assumption 11.14.1, as  $n \rightarrow \infty$ ,*

$$\begin{aligned}\widehat{\mathbf{V}}_{\boldsymbol{\beta}}^0 &\xrightarrow{p} \mathbf{V}_{\boldsymbol{\beta}}^0 \\ \widehat{\mathbf{V}}_{\boldsymbol{\beta}} &\xrightarrow{p} \mathbf{V}_{\boldsymbol{\beta}}.\end{aligned}$$

To prove Theorem 11.16.1 the key is to show  $\widehat{\Omega} \xrightarrow{p} \Omega$  as the other convergence results were established in the proof of consistency. We defer this to Exercise 11.6.

It is important that the covariance matrix be constructed using the correct residual formula  $\widehat{e}_i = y_i - \mathbf{x}_i' \widehat{\boldsymbol{\beta}}_{2\text{sls}}$ . This is different than what would be obtained if the “two-stage” computation method is used. To see this, let’s walk through the two-stage method. First, we estimate the reduced form

$$\mathbf{x}_i = \widehat{\boldsymbol{\Gamma}}' \mathbf{z}_i + \widehat{\mathbf{u}}_i$$

to obtain the predicted values  $\widehat{\mathbf{x}}_i = \widehat{\boldsymbol{\Gamma}}' \mathbf{z}_i$ . Second, we regress  $y_i$  on  $\widehat{\mathbf{x}}_i$  to obtain the 2SLS estimator  $\widehat{\boldsymbol{\beta}}_{2\text{sls}}$ . This latter regression takes the form

$$y_i = \widehat{\mathbf{x}}_i' \widehat{\boldsymbol{\beta}}_{2\text{sls}} + \widehat{v}_i \tag{11.46}$$

where  $\widehat{v}_i$  are least-squares residuals. The covariance matrix (and standard errors) reported by this regression are constructed using the residual  $\widehat{v}_i$ . For example, the homoskedastic formula is

$$\begin{aligned}\widehat{\mathbf{V}}_{\boldsymbol{\beta}} &= \left( \frac{1}{n} \widehat{\mathbf{X}}' \widehat{\mathbf{X}} \right)^{-1} \widehat{\sigma}_v^2 = \left( \widehat{\mathbf{Q}}_{xz} \widehat{\mathbf{Q}}_{zz}^{-1} \widehat{\mathbf{Q}}_{zx} \right)^{-1} \widehat{\sigma}_v^2 \\ \widehat{\sigma}_v^2 &= \frac{1}{n} \sum_{i=1}^n \widehat{v}_i^2\end{aligned}$$



which is proportional to the variance estimate  $\hat{\sigma}_v^2$  rather than  $\hat{\sigma}^2$ . This is important because the residual  $\hat{v}_i$  differs from  $\hat{e}_i$ . We can see this because the regression (11.46) uses the regressor  $\hat{\mathbf{x}}_i$  rather than  $\mathbf{x}_i$ . Indeed, we can calculate that

$$\begin{aligned}\hat{v}_i &= y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{2\text{sls}} + (\mathbf{x}_i - \hat{\mathbf{x}}_i)' \hat{\boldsymbol{\beta}}_{2\text{sls}} \\ &= \hat{e}_i + \hat{\mathbf{u}}_i' \hat{\boldsymbol{\beta}}_{2\text{sls}} \\ &\neq \hat{e}_i\end{aligned}$$

This means that standard errors reported by the regression (11.46) will be incorrect.

This problem is avoided if the 2SLS estimator is constructed directly and the standard errors calculated with the correct formula rather than taking the “two-step” shortcut.

## 11.17 Asymptotic Distribution and Covariance Estimation for LIML

Recall, the LIML estimator has several representations, including

$$\begin{aligned}\hat{\boldsymbol{\beta}}_{\text{lml}} &= (\mathbf{X}'(\mathbf{I}_n - \hat{\kappa}\mathbf{M}_Z)\mathbf{X})^{-1}(\mathbf{X}'(\mathbf{I}_n - \hat{\kappa}\mathbf{M}_Z)\mathbf{y}) \\ &= (\mathbf{X}'\mathbf{P}_Z\mathbf{X} - \hat{\mu}\mathbf{X}'\mathbf{M}_Z\mathbf{X})^{-1}(\mathbf{X}'\mathbf{P}_Z\mathbf{y} - \hat{\mu}\mathbf{X}'\mathbf{M}_Z\mathbf{y})\end{aligned}$$

where  $\hat{\mu} = \hat{\kappa} - 1$  and

$$\hat{\kappa} = \min_{\gamma} \frac{\gamma' \mathbf{W}' \mathbf{M}_1 \mathbf{W} \gamma}{\gamma' \mathbf{W}' \mathbf{M}_Z \mathbf{W} \gamma}.$$

Using multivariate regression analysis, we can show that  $\hat{\kappa} \xrightarrow{p} 1$  and thus  $\hat{\mu} \xrightarrow{p} 0$ . It follows that

$$\begin{aligned}\sqrt{n}(\hat{\boldsymbol{\beta}}_{\text{lml}} - \boldsymbol{\beta}) &= \left( \frac{1}{n} \mathbf{X}' \mathbf{P}_Z \mathbf{X} - \hat{\mu} \frac{1}{n} \mathbf{X}' \mathbf{M}_Z \mathbf{X} \right)^{-1} \left( \frac{1}{\sqrt{n}} \mathbf{X}' \mathbf{P}_Z \mathbf{e} - \hat{\mu} \frac{1}{\sqrt{n}} \mathbf{X}' \mathbf{M}_Z \mathbf{e} \right) \\ &= \left( \frac{1}{n} \mathbf{X}' \mathbf{P}_Z \mathbf{X} - o_p(1) \right)^{-1} \left( \frac{1}{\sqrt{n}} \mathbf{X}' \mathbf{P}_Z \mathbf{e} - o_p(1) \right) \\ &= \sqrt{n}(\hat{\boldsymbol{\beta}}_{2\text{sls}} - \boldsymbol{\beta}) + o_p(1)\end{aligned}$$

which means that LIML and 2SLS have the same asymptotic distribution. This holds under the same assumptions as for 2SLS, and in particular does not require normality of the errors.

Consequently, one method to obtain an asymptotically valid covariance estimate for LIML is to use the same formula as for 2SLS. However, this is not the best choice. Rather, consider the IV representation for LIML

$$\hat{\boldsymbol{\beta}}_{\text{lml}} = (\widetilde{\mathbf{X}}' \mathbf{X})^{-1} (\widetilde{\mathbf{X}}' \mathbf{y})$$

where

$$\widetilde{\mathbf{X}} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 - \hat{\kappa} \hat{\mathbf{U}}_2 \end{pmatrix}$$

and  $\hat{\mathbf{U}}_2 = \mathbf{M}_Z \mathbf{X}_2$ . The asymptotic covariance matrix formula for an IV estimator is

$$\hat{\mathbf{V}}_{\boldsymbol{\beta}} = \left( \frac{1}{n} \widetilde{\mathbf{X}}' \mathbf{X} \right)^{-1} \hat{\boldsymbol{\Omega}} \left( \frac{1}{n} \mathbf{X}' \widetilde{\mathbf{X}} \right)^{-1} \quad (11.47)$$

where

$$\begin{aligned}\hat{\boldsymbol{\Omega}} &= \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i' \hat{e}_i^2 \\ \hat{e}_i &= y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{\text{lml}}.\end{aligned}$$

This simplifies to the 2SLS formula when  $\hat{\kappa} = 1$  but otherwise differs. The estimator (11.47) is a better choice than the 2SLS formula for covariance matrix estimation as it takes advantage of the LIML estimator structure.

## 11.18 Functions of Parameters

Given the distribution theory in Theorems 11.14.1 and 11.16.1 it is straightforward to derive the asymptotic distribution of smooth nonlinear functions of the coefficients.

Specifically, given a function  $\mathbf{r}(\boldsymbol{\beta}) : \mathbb{R}^k \rightarrow \Theta \subset \mathbb{R}^q$  we define the parameter

$$\boldsymbol{\theta} = \mathbf{r}(\boldsymbol{\beta})$$

Given  $\hat{\boldsymbol{\beta}}_{2\text{sls}}$  a natural estimator of  $\boldsymbol{\theta}$  is  $\hat{\boldsymbol{\theta}}_{2\text{sls}} = \mathbf{r}(\hat{\boldsymbol{\beta}}_{2\text{sls}})$ .

Consistency follows from Theorem 11.13.1 and the continuous mapping theorem.

**Theorem 11.18.1** *Under Assumption 11.13.1, if  $\mathbf{r}(\boldsymbol{\beta})$  is continuous at  $\boldsymbol{\beta}$ , then  $\hat{\boldsymbol{\theta}}_{2\text{sls}} \xrightarrow{p} \boldsymbol{\theta}$  as  $n \rightarrow \infty$ .*

If  $\mathbf{r}(\boldsymbol{\beta})$  is differentiable then an estimator of the asymptotic covariance matrix for  $\hat{\boldsymbol{\theta}}$  is

$$\begin{aligned}\hat{\mathbf{V}}_{\boldsymbol{\theta}} &= \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\boldsymbol{\beta}} \hat{\mathbf{R}} \\ \hat{\mathbf{R}} &= \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{r}(\hat{\boldsymbol{\beta}}_{2\text{sls}})'\end{aligned}$$

We similarly define the homoskedastic variance estimator as

$$\hat{\mathbf{V}}_{\boldsymbol{\theta}}^0 = \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\boldsymbol{\beta}}^0 \hat{\mathbf{R}}.$$

The asymptotic distribution theory follows from Theorems 11.14.1 and 11.16.1, and the delta method.

**Theorem 11.18.2** *Under Assumption 11.14.1, if  $\mathbf{r}(\boldsymbol{\beta})$  is continuously differentiable at  $\boldsymbol{\beta}$ , then as  $n \rightarrow \infty$*

$$\sqrt{n} \left( \hat{\boldsymbol{\theta}}_{2\text{sls}} - \boldsymbol{\theta} \right) \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{V}_{\boldsymbol{\theta}})$$

where

$$\begin{aligned}\mathbf{V}_{\boldsymbol{\theta}} &= \mathbf{R}' \mathbf{V}_{\boldsymbol{\beta}} \mathbf{R} \\ \mathbf{R} &= \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{r}(\boldsymbol{\beta})'\end{aligned}$$

and

$$\hat{\mathbf{V}}_{\boldsymbol{\theta}} \xrightarrow{p} \mathbf{V}_{\boldsymbol{\theta}}.$$

When  $q = 1$ , a standard error for  $\hat{\boldsymbol{\theta}}_{2\text{sls}}$  is  $s(\hat{\boldsymbol{\theta}}_{2\text{sls}}) = \sqrt{n^{-1} \hat{\mathbf{V}}_{\boldsymbol{\theta}}}$ .

For example, let's take the parameter estimates from the fifth column of Table 11.1, which are the 2SLS estimates with three endogenous regressors and four excluded instruments. Suppose we are interested in the return to experience, which depends on the level of experience. The estimated return at *experience* = 10 is  $0.0473 - 0.032 * 2 * 10/100 = 0.041$  and its standard error is 0.003. This implies a 4% increase in wages per year of experience and is precisely estimated. Or suppose we are interested in the level of experience at which the function maximizes. The estimate is  $50 * 0.047/0.032 = 73$ . This has a standard error of 249. The large standard error implies that the estimate (73 years of experience) is without precision and is thus uninformative.

## 11.19 Hypothesis Tests

As in the previous section, for a given function  $\mathbf{r}(\boldsymbol{\beta}) : \mathbb{R}^k \rightarrow \Theta \subset \mathbb{R}^q$  we define the parameter  $\boldsymbol{\theta} = \mathbf{r}(\boldsymbol{\beta})$  and consider tests of hypotheses of the form

$$\mathbb{H}_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$$

against

$$\mathbb{H}_1 : \boldsymbol{\theta} \neq \boldsymbol{\theta}_0.$$

The Wald statistic for  $\mathbb{H}_0$  is

$$W = n \left( \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \right)' \hat{\mathbf{V}}_{\hat{\boldsymbol{\theta}}}^{-1} \left( \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \right).$$

From Theorem 11.18.2 we deduce that  $W$  is asymptotically chi-square distributed. Let  $G_q(u)$  denote the  $\chi_q^2$  distribution function.

**Theorem 11.19.1** *Under Assumption 11.14.1, if  $\mathbf{r}(\boldsymbol{\beta})$  is continuously differentiable at  $\boldsymbol{\beta}$ , and  $\mathbb{H}_0$  holds, then as  $n \rightarrow \infty$ ,*

$$W \xrightarrow{d} \chi_q^2.$$

*For  $c$  satisfying  $\alpha = 1 - G_q(c)$ ,*

$$\Pr(W > c \mid \mathbb{H}_0) \longrightarrow \alpha$$

*so the test “Reject  $\mathbb{H}_0$  if  $W > c$ ” has asymptotic size  $\alpha$ .*

In linear regression we often report the  $F$  version of the Wald statistic (by dividing by degrees of freedom) and use the  $F$  distribution for inference, as this is justified in the normal sampling model. For 2SLS estimation, however, this is not done as there is no finite sample  $F$  justification for the  $F$  version of the Wald statistic.

To illustrate, once again let's take the parameter estimates from the fifth column of Table 11.1 and again consider the return to experience which is determined by the coefficients on *experience* and *experience*<sup>2</sup>/100. Neither coefficient is statistically significant at the 5% level, so it is unclear from a casual look if the overall effect is statistically significant. We can assess this by testing the joint hypothesis that both coefficients are zero. The Wald statistic for this hypothesis is  $W = 254$ , which is highly significant with an asymptotic p-value of 0.0000. Thus by examining the joint test, in contrast to the individual tests, is quite clear that experience has a non-zero effect.

## 11.20 Finite Sample Theory

In Chapter 5 we reviewed the rich exact distribution available for the linear regression model under the assumption of normal innovations. There was a similarly rich literature in econometrics which developed a distribution theory for IV, 2SLS and LIML estimators. This theory is reviewed by Peter Phillips (1983), and much of the theory was developed by Peter Phillips in a series of papers in the 1970s and early 1980s.

This theory was developed under the assumption that the structural error vector  $\mathbf{e}$  and reduced form error  $\mathbf{u}_2$  are multivariate normally distributed. The challenge is that the IV estimators are non-linear functions of  $\mathbf{u}_2$  and are thus non-normally distributed. Formulae for the exact distributions have been derived, but are unfortunately functions of model parameters and hence are not directly useful for finite sample inference.

One important implication of this literature is that it is quite clear that even in this optimal context of exact normal innovations, the finite sample distributions of the IV estimators are non-normal and the finite sample distributions of test statistics are not chi-squared. The normal and chi-squared approximations hold asymptotically, but there is no reason to expect these approximations to be accurate in finite samples.

## 11.21 Clustered Dependence

In Section 4.20 we introduced clustered dependence. We can also use the methods of clustered dependence for 2SLS estimation. Recall, the  $g^{th}$  cluster has the observations  $\mathbf{y}_g = (y_{1g}, \dots, y_{n_{gg}})'$ ,  $\mathbf{X}_g = (\mathbf{x}_{1g}, \dots, \mathbf{x}_{n_{gg}})'$ , and  $\mathbf{Z}_g = (\mathbf{z}_{1g}, \dots, \mathbf{z}_{n_{gg}})'$ . The structural equation for the  $g^{th}$  cluster can be written as the matrix system

$$\mathbf{y}_g = \mathbf{X}_g \boldsymbol{\beta} + \mathbf{e}_g.$$

Using this notation the center 2SLS estimator can be written as

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{2sls} - \boldsymbol{\beta} &= \left( \mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{e} \\ &= \left( \mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \left( \sum_{g=1}^G \mathbf{Z}'_g \mathbf{e}_g \right). \end{aligned}$$

The cluster-robust covariance matrix estimator for  $\hat{\boldsymbol{\beta}}_{2sls}$  thus takes the form

$$\hat{\mathbf{V}}_{\boldsymbol{\beta}} = \left( \mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \hat{\mathbf{S}} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X} \left( \mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X} \right)^{-1}$$

with

$$\hat{\mathbf{S}} = \sum_{g=1}^G \mathbf{Z}'_g \hat{\mathbf{e}}_g \hat{\mathbf{e}}'_g \mathbf{Z}_g$$

and the clustered residuals

$$\hat{\mathbf{e}}_g = \mathbf{y}_g - \mathbf{X}_g \hat{\boldsymbol{\beta}}_{2sls}.$$

The difference between the heteroskedasticity-robust estimator and the cluster-robust estimator is the covariance estimator  $\hat{\mathbf{S}}$ .

## 11.22 Generated Regressors

The “two-stage” form of the 2SLS estimator is an example of what is called “estimation with generated regressors”. We say a regressor is a **generated** if it is an estimate of an idealized regressor, or if it is a function of estimated parameters. Typically, a generated regressor  $\hat{\mathbf{w}}_i$  is an estimate of an unobserved ideal regressor  $\mathbf{w}_i$ . As an estimate,  $\hat{\mathbf{w}}_i$  is a function of the sample, not just observation  $i$ . Hence it is not “i.i.d.” as it is dependent across observations, which invalidates the conventional regression assumptions. Consequently, the sampling distribution of regression estimates is affected. Unless this is incorporated into our inference methods, covariance matrix estimates and standard errors will be incorrect.

The econometric theory of generated regressors was developed by Pagan (1984) for linear models, and extended to non-linear models and more general two-step estimators by Pagan (1986). Here we focus on the linear model:

$$\begin{aligned} y_i &= \mathbf{w}'_i \boldsymbol{\beta} + v_i \\ \mathbf{w}_i &= \mathbf{A}' \mathbf{z}_i \\ \mathbb{E}(\mathbf{z}_i v_i) &= \mathbf{0}. \end{aligned} \tag{11.48}$$

The observables are  $(y_i, \mathbf{z}_i)$ . We also have an estimate  $\hat{\mathbf{A}}$  of  $\mathbf{A}$ .

Given  $\hat{\mathbf{A}}$  we construct the estimate  $\hat{\mathbf{w}}_i = \hat{\mathbf{A}}' \mathbf{z}_i$  of  $\mathbf{w}_i$ , replace  $\mathbf{w}_i$  in (11.48) with  $\hat{\mathbf{w}}_i$ , and then estimate  $\beta$  by least-squares, resulting in the estimator

$$\hat{\beta} = \left( \sum_{i=1}^n \hat{\mathbf{w}}_i \hat{\mathbf{w}}_i' \right)^{-1} \left( \sum_{i=1}^n \hat{\mathbf{w}}_i y_i \right). \quad (11.49)$$

The regressors  $\hat{\mathbf{w}}_i$  are called **generated regressors**. The properties of  $\hat{\beta}$  are different than least-squares with i.i.d. observations, since the generated regressors are themselves estimates.

This framework includes the 2SLS estimator as well as other common estimators. The 2SLS model can be written as (11.48) by looking at the reduced form equation (11.16), with  $\mathbf{w}_i = \mathbf{\Gamma}' \mathbf{z}_i$ ,  $\mathbf{A} = \mathbf{\Gamma}$ , and  $\hat{\mathbf{A}} = \hat{\mathbf{\Gamma}}$  is (11.22).

The examples which motivated Pagan (1984) emerged from the macroeconomics literature, in particular the work of Barro (1977) which examined the impact of inflation expectations and expectation errors on economic output. For example, let  $\pi_i$  denote realized inflation and  $\mathbf{z}_i$  be the information available to economic agents. A model of inflation expectations sets  $w_i = \mathbb{E}(\pi_i | \mathbf{z}_i) = \gamma' \mathbf{z}_i$  and a model of expectation error sets  $w_i = \pi_i - \mathbb{E}(\pi_i | \mathbf{z}_i) = \pi_i - \gamma' \mathbf{z}_i$ . Since expectations and errors are not observed they are replaced in applications with the fitted values  $\hat{w}_i = \hat{\gamma}' \mathbf{z}_i$  or residuals  $\hat{w}_i = \pi_i - \hat{\gamma}' \mathbf{z}_i$  where  $\hat{\gamma}$  is a coefficient estimate from a regression of  $\pi_i$  on  $\mathbf{z}_i$ .

The generated regressor framework includes all of these examples.

The goal is to obtain a distributional approximation for  $\hat{\beta}$  in order to construct standard errors, confidence intervals and conduct tests. Start by substituting equation (11.48) into (11.49). We obtain

$$\hat{\beta} = \left( \sum_{i=1}^n \hat{\mathbf{w}}_i \hat{\mathbf{w}}_i' \right)^{-1} \left( \sum_{i=1}^n \hat{\mathbf{w}}_i (\mathbf{w}_i' \beta + v_i) \right).$$

Next, substitute  $\mathbf{w}_i' \beta = \hat{\mathbf{w}}_i' \beta + (\mathbf{w}_i - \hat{\mathbf{w}}_i)' \beta$ . We obtain

$$\hat{\beta} - \beta = \left( \sum_{i=1}^n \hat{\mathbf{w}}_i \hat{\mathbf{w}}_i' \right)^{-1} \left( \sum_{i=1}^n \hat{\mathbf{w}}_i ((\mathbf{w}_i - \hat{\mathbf{w}}_i)' \beta + v_i) \right). \quad (11.50)$$

Effectively, this shows that the distribution of  $\hat{\beta} - \beta$  has two random components, one due to the conventional regression component  $\hat{\mathbf{w}}_i v_i$ , and the second due to the generated regressor  $(\mathbf{w}_i - \hat{\mathbf{w}}_i)' \beta$ . Conventional variance estimators do not address this second component and thus will be biased.

Interestingly, the distribution in (11.50) dramatically simplifies in the special case that the “generated regressor term”  $(\mathbf{w}_i - \hat{\mathbf{w}}_i)' \beta$  disappears. This occurs when the slope coefficients on the generated regressors are zero. To be specific, partition  $\mathbf{w}_i = (\mathbf{w}_{1i}, \mathbf{w}_{2i})$ ,  $\hat{\mathbf{w}}_i = (\hat{\mathbf{w}}_{1i}, \hat{\mathbf{w}}_{2i})$ , and  $\beta = (\beta_1, \beta_2)$  so that  $\mathbf{w}_{1i}$  are the conventional observed regressors and  $\hat{\mathbf{w}}_{2i}$  are the generated regressors. Then  $(\mathbf{w}_i - \hat{\mathbf{w}}_i)' \beta = (\mathbf{w}_{2i} - \hat{\mathbf{w}}_{2i})' \beta_2$ . Thus if  $\beta_2 = \mathbf{0}$  this term disappears. In this case (11.50) equals

$$\hat{\beta} - \beta = \left( \sum_{i=1}^n \hat{\mathbf{w}}_i \hat{\mathbf{w}}_i' \right)^{-1} \left( \sum_{i=1}^n \hat{\mathbf{w}}_i v_i \right).$$

This is a dramatic simplification.

Furthermore, since  $\hat{\mathbf{w}}_i = \hat{\mathbf{A}}' \mathbf{z}_i$  we can write the estimator as a function of sample moments:

$$\sqrt{n} (\hat{\beta} - \beta) = \left( \hat{\mathbf{A}}' \left( \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i' \right) \hat{\mathbf{A}} \right)^{-1} \hat{\mathbf{A}}' \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{z}_i v_i \right).$$

If  $\hat{\mathbf{A}} \xrightarrow{p} \mathbf{A}$  we find from standard manipulations that

$$\sqrt{n} (\hat{\beta} - \beta) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_\beta)$$

where

$$\mathbf{V}_\beta = (\mathbf{A}' \mathbb{E}(\mathbf{z}_i \mathbf{z}_i') \mathbf{A})^{-1} (\mathbf{A}' \mathbb{E}(\mathbf{z}_i \mathbf{z}_i' v_i^2) \mathbf{A}) (\mathbf{A}' \mathbb{E}(\mathbf{z}_i \mathbf{z}_i') \mathbf{A})^{-1}. \quad (11.51)$$

The conventional asymptotic covariance matrix estimator for  $\hat{\beta}$  takes the form

$$\hat{\mathbf{V}}_\beta = \left( \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{w}}_i \hat{\mathbf{w}}_i' \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{w}}_i \hat{\mathbf{w}}_i' \hat{v}_i^2 \right) \left( \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{w}}_i \hat{\mathbf{w}}_i' \right)^{-1} \quad (11.52)$$

where  $\hat{v}_i = y_i - \hat{\mathbf{w}}_i' \hat{\beta}$ . Under the given assumptions,  $\hat{\mathbf{V}}_\beta \xrightarrow{p} \mathbf{V}_\beta$ . Thus inference using  $\hat{\mathbf{V}}_\beta$  is asymptotically valid. This is useful when we are interested in tests of  $\beta_2 = \mathbf{0}$ . Often this is of major interest in applications.

To test  $\mathbb{H}_0 : \beta_2 = \mathbf{0}$  we partition  $\hat{\beta} = (\hat{\beta}_1, \hat{\beta}_2)$  and construct a conventional Wald statistic

$$W = n \hat{\beta}_2' \left( \left[ \hat{\mathbf{V}}_\beta \right]_{22} \right)^{-1} \hat{\beta}_2.$$

**Theorem 11.22.1** Take model (11.48) with  $\mathbb{E}(y_i^4) < \infty$ ,  $\mathbb{E}\|\mathbf{z}_i\|^4 < \infty$ ,  $\mathbf{A}' \mathbb{E}(\mathbf{z}_i \mathbf{z}_i') \mathbf{A} > 0$ ,  $\hat{\mathbf{A}} \xrightarrow{p} \mathbf{A}$  and  $\hat{\mathbf{w}}_i = (\mathbf{w}_{1i}, \mathbf{w}_{2i})$ . Under  $\mathbb{H}_0 : \beta_2 = \mathbf{0}$ , then as  $n \rightarrow \infty$ ,

$$\sqrt{n} (\hat{\beta} - \beta) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_\beta)$$

where  $\mathbf{V}_\beta$  is given in (11.51). For  $\hat{\mathbf{V}}_\beta$  given in (11.52),

$$\hat{\mathbf{V}}_\beta \xrightarrow{p} \mathbf{V}_\beta.$$

Furthermore,

$$W \xrightarrow{d} \chi_q^2$$

where  $q = \dim(\beta_2)$ . For  $c$  satisfying  $\alpha = 1 - G_q(c)$

$$\Pr(W > c \mid \mathbb{H}_0) \longrightarrow \alpha$$

so the test “Reject  $\mathbb{H}_0$  if  $W > c$ ” has asymptotic size  $\alpha$ .

In the special case that  $\hat{\mathbf{A}} = \mathbf{A}(\mathbf{X}, \mathbf{Z})$  and  $v_i | \mathbf{x}_i, \mathbf{z}_i \sim N(0, \sigma^2)$  then there is a finite sample version of the previous result. Let  $W^0$  be the Wald statistic constructed with a homoskedastic variance matrix estimator, and let

$$F = W/q \quad (11.53)$$

be the the  $F$  statistic, where  $q = \dim(\beta_2)$ .

**Theorem 11.22.2** Take model (11.48) with  $\hat{\mathbf{A}} = \mathbf{A}(\mathbf{X}, \mathbf{Z})$ ,  $v_i | \mathbf{x}_i, \mathbf{z}_i \sim N(0, \sigma^2)$  and  $\hat{\mathbf{w}}_i = (\mathbf{w}_{1i}, \mathbf{w}_{2i})$ . Under  $\mathbb{H}_0 : \beta_2 = \mathbf{0}$ ,  $t$ -statistics have exact  $N(0, 1)$  distributions, and the  $F$  statistic (11.53) has an exact  $F_{q, n-k}$  distribution, where  $q = \dim(\beta_2)$  and  $k = \dim(\beta)$ .

The theory introduced above allows tests of  $\mathbb{H}_0 : \beta_2 = \mathbf{0}$  but does not lead to methods to construct standard errors or confidence intervals. For this, we need to work out the distribution without imposing the simplification  $\beta_2 = \mathbf{0}$ . This often needs to be worked out case-by-case, or by using methods based on the generalized method of moments to be introduced in Chapter 12. However, in some important set of examples it is straightforward to work out the asymptotic distribution.

For the remainder of this section we examine the setting where the estimators  $\hat{\mathbf{A}}$  take a least-squares form, so for some  $\mathbf{X}$  can be written as  $\hat{\mathbf{A}} = (\mathbf{Z}'\mathbf{Z})^{-1}(\mathbf{Z}'\mathbf{X})$ . Such estimators correspond to the multivariate projection model

$$\begin{aligned} \mathbf{x}_i &= \mathbf{A}'\mathbf{z}_i + \mathbf{u}_i \\ \mathbb{E}(\mathbf{z}_i\mathbf{u}_i') &= \mathbf{0}. \end{aligned} \tag{11.54}$$

This class of estimators directly includes 2SLS and the expectation model described above. We can write the matrix of generated regressors as  $\widehat{\mathbf{W}} = \mathbf{Z}\hat{\mathbf{A}}$  and then (11.50) as

$$\begin{aligned} \hat{\beta} - \beta &= (\widehat{\mathbf{W}}'\widehat{\mathbf{W}})^{-1} (\widehat{\mathbf{W}}'((\mathbf{W} - \widehat{\mathbf{W}})\beta + \mathbf{v})) \\ &= (\hat{\mathbf{A}}'\mathbf{Z}'\mathbf{Z}\hat{\mathbf{A}})^{-1} (\hat{\mathbf{A}}'\mathbf{Z}'(-\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}(\mathbf{Z}'\mathbf{U})\beta + \mathbf{v})) \\ &= (\hat{\mathbf{A}}'\mathbf{Z}'\mathbf{Z}\hat{\mathbf{A}})^{-1} (\hat{\mathbf{A}}'\mathbf{Z}'(-\mathbf{U}\beta + \mathbf{v})) \\ &= (\hat{\mathbf{A}}'\mathbf{Z}'\mathbf{Z}\hat{\mathbf{A}})^{-1} (\hat{\mathbf{A}}'\mathbf{Z}'\mathbf{e}) \end{aligned}$$

where

$$e_i = v_i - \mathbf{u}_i'\beta = y_i - \mathbf{x}_i'\beta. \tag{11.55}$$

This estimator has the asymptotic distribution

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_\beta)$$

where

$$\mathbf{V}_\beta = (\mathbf{A}'\mathbb{E}(\mathbf{z}_i\mathbf{z}_i')\mathbf{A})^{-1}(\mathbf{A}'\mathbb{E}(\mathbf{z}_i\mathbf{z}_i'e_i^2)\mathbf{A})(\mathbf{A}'\mathbb{E}(\mathbf{z}_i\mathbf{z}_i')\mathbf{A})^{-1}. \tag{11.56}$$

Under conditional homoskedasticity the covariance matrix simplifies to

$$\mathbf{V}_\beta = (\mathbf{A}'\mathbb{E}(\mathbf{z}_i\mathbf{z}_i')\mathbf{A})^{-1}\mathbb{E}(e_i^2).$$

An appropriate estimator of  $\mathbf{V}_\beta$  is

$$\begin{aligned} \hat{\mathbf{V}}_\beta &= \left(\frac{1}{n}\widehat{\mathbf{W}}'\widehat{\mathbf{W}}\right)^{-1} \left(\frac{1}{n}\sum_{i=1}^n \widehat{\mathbf{w}}_i\widehat{\mathbf{w}}_i'\widehat{e}_i^2\right) \left(\frac{1}{n}\widehat{\mathbf{W}}'\widehat{\mathbf{W}}\right)^{-1} \\ \widehat{e}_i &= y_i - \mathbf{x}_i'\hat{\beta}. \end{aligned} \tag{11.57}$$

Under the assumption of conditional homoskedasticity this can be simplified as usual.

This appears to be the usual covariance matrix estimator, but it is not, because the least-squares residuals  $\widehat{v}_i = y_i - \widehat{\mathbf{w}}_i'\hat{\beta}$  have been replaced with  $\widehat{e}_i = y_i - \mathbf{x}_i'\hat{\beta}$ . This is exactly the substitution made by the 2SLS covariance matrix formula. Indeed, the covariance matrix estimator  $\hat{\mathbf{V}}_\beta$  precisely equals the estimator (11.45).

**Theorem 11.22.3** Take model (11.48) and (11.54) with  $\mathbb{E}(y_i^4) < \infty$ ,  $\mathbb{E}\|z_i\|^4 < \infty$ ,  $\mathbf{A}'\mathbb{E}(z_i z_i')\mathbf{A} > 0$ , and  $\hat{\mathbf{A}} = (\mathbf{Z}'\mathbf{Z})^{-1}(\mathbf{Z}'\mathbf{X})$ . As  $n \rightarrow \infty$ ,

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_{\boldsymbol{\beta}})$$

where  $\mathbf{V}_{\boldsymbol{\beta}}$  is given in (11.56) with  $e_i$  defined in (11.55). For  $\hat{\mathbf{V}}_{\boldsymbol{\beta}}$  given in (11.57),

$$\hat{\mathbf{V}}_{\boldsymbol{\beta}} \xrightarrow{p} \mathbf{V}_{\boldsymbol{\beta}}.$$

Since the parameter estimates are asymptotically normal and the covariance matrix is consistently estimated, standard errors and test statistics constructed from  $\hat{\mathbf{V}}_{\boldsymbol{\beta}}$  are asymptotically valid with conventional interpretations.

We now summarize the results of this section. In general, care needs to be exercised when estimating models with generated regressors. As a general rule, generated regressors and two-step estimation affects sampling distributions and variance matrices. An important simplification occurs for tests that the generated regressors have zero slopes. In this case conventional tests have conventional distributions, both asymptotically and in finite samples. Another important special case occurs when the generated regressors are least-squares fitted values. In this case the asymptotic distribution takes a conventional form, but the conventional residual needs to be replaced by one constructed with the forecasted variable. With this one modification asymptotic inference using the generated regressors is conventional.

## 11.23 Regression with Expectation Errors

In this section we examine a generated regressor model which includes expectation errors in the regression. This is an important class of generated regressor models, and is relatively straightforward to characterize.

The model is

$$\begin{aligned} y_i &= \mathbf{w}_i' \boldsymbol{\beta} + \mathbf{u}_i' \boldsymbol{\alpha} + v_i \\ \mathbf{w}_i &= \mathbf{A}' \mathbf{z}_i \\ \mathbf{x}_i &= \mathbf{w}_i + \mathbf{u}_i \\ \mathbb{E}(\mathbf{z}_i \varepsilon_i) &= \mathbf{0} \\ \mathbb{E}(\mathbf{u}_i \varepsilon_i) &= \mathbf{0} \\ \mathbb{E}(\mathbf{z}_i \mathbf{u}_i') &= \mathbf{0}. \end{aligned}$$

The observables are  $(y_i, \mathbf{x}_i, \mathbf{z}_i)$ . This model states that  $\mathbf{w}_i$  is the expectation of  $\mathbf{x}_i$  (or more generally, the projection of  $\mathbf{x}_i$  on  $\mathbf{z}_i$ ) and  $\mathbf{u}_i$  is its expectation error. The model allows for exogenous regressors as in the standard IV model if they are listed in  $\mathbf{w}_i$ ,  $\mathbf{x}_i$  and  $\mathbf{z}_i$ . This model is used, for example, to decompose the effect of expectations from expectation errors. In some cases it is desired to include only the expectation error  $\mathbf{u}_i$ , not the expectation  $\mathbf{w}_i$ . This does not change the results described here.

The model is estimated as follows. First,  $\mathbf{A}$  is estimated by multivariate least-squares of  $\mathbf{x}_i$  on  $\mathbf{z}_i$ ,  $\hat{\mathbf{A}} = (\mathbf{Z}'\mathbf{Z})^{-1}(\mathbf{Z}'\mathbf{X})$ , which yields as by-products the fitted values  $\hat{\mathbf{W}} = \mathbf{Z}\hat{\mathbf{A}}$  and residuals  $\hat{\mathbf{U}} = \mathbf{X} - \hat{\mathbf{W}}$ . Second, the coefficients are estimated by least-squares of  $y_i$  on the fitted values  $\hat{\mathbf{w}}_i$  and residuals  $\hat{\mathbf{u}}_i$

$$y_i = \hat{\mathbf{w}}_i' \hat{\boldsymbol{\beta}} + \hat{\mathbf{u}}_i' \hat{\boldsymbol{\alpha}} + \hat{v}_i.$$

We now examine the asymptotic distributions of these estimates.



By the first-step regression  $\mathbf{Z}'\hat{\mathbf{U}} = \mathbf{0}$ ,  $\widehat{\mathbf{W}}'\hat{\mathbf{U}} = \mathbf{0}$  and  $\mathbf{W}'\hat{\mathbf{U}} = \mathbf{0}$ . This means that  $\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\alpha}}$  can be computed separately. Notice that

$$\hat{\boldsymbol{\beta}} = \left( \widehat{\mathbf{W}}' \widehat{\mathbf{W}} \right)^{-1} \widehat{\mathbf{W}}' \mathbf{y}$$

and

$$\mathbf{y} = \widehat{\mathbf{W}}\boldsymbol{\beta} + \mathbf{U}\boldsymbol{\alpha} + \left( \mathbf{W} - \widehat{\mathbf{W}} \right) \boldsymbol{\beta} + \mathbf{v}.$$

Substituting, using  $\widehat{\mathbf{W}}'\hat{\mathbf{U}} = \mathbf{0}$  and  $\mathbf{W} - \widehat{\mathbf{W}} = -\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{U}$  we find

$$\begin{aligned} \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} &= \left( \widehat{\mathbf{W}}' \widehat{\mathbf{W}} \right)^{-1} \widehat{\mathbf{W}}' \left( \mathbf{U}\boldsymbol{\alpha} + \left( \mathbf{W} - \widehat{\mathbf{W}} \right) \boldsymbol{\beta} + \mathbf{v} \right) \\ &= \left( \hat{\mathbf{A}}' \mathbf{Z}' \mathbf{Z} \hat{\mathbf{A}} \right)^{-1} \hat{\mathbf{A}}' \mathbf{Z}' (\mathbf{U}\boldsymbol{\alpha} - \mathbf{U}\boldsymbol{\beta} + \mathbf{v}) \\ &= \left( \hat{\mathbf{A}}' \mathbf{Z}' \mathbf{Z} \hat{\mathbf{A}} \right)^{-1} \hat{\mathbf{A}}' \mathbf{Z}' \mathbf{e} \end{aligned}$$

where

$$e_i = v_i + \mathbf{u}_i'(\boldsymbol{\alpha} - \boldsymbol{\beta}) = y_i - \mathbf{x}_i'\boldsymbol{\beta}.$$

We also find

$$\hat{\boldsymbol{\alpha}} = \left( \hat{\mathbf{U}}' \hat{\mathbf{U}} \right)^{-1} \hat{\mathbf{U}}' \mathbf{y}.$$

Since  $\hat{\mathbf{U}}'\mathbf{W} = \mathbf{0}$ ,  $\mathbf{U} - \hat{\mathbf{U}} = \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{U}$  and  $\hat{\mathbf{U}}'\mathbf{Z} = \mathbf{0}$  then

$$\begin{aligned} \hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha} &= \left( \hat{\mathbf{U}}' \hat{\mathbf{U}} \right)^{-1} \hat{\mathbf{U}}' \left( \mathbf{W}\boldsymbol{\beta} + \left( \mathbf{U} - \hat{\mathbf{U}} \right) \boldsymbol{\alpha} + \mathbf{v} \right) \\ &= \left( \hat{\mathbf{U}}' \hat{\mathbf{U}} \right)^{-1} \hat{\mathbf{U}}' \mathbf{v}. \end{aligned}$$

Together, we establish the following distributional result.

**Theorem 11.23.1** *For the model and estimates described in this section, with  $\mathbb{E}(y_i^4) < \infty$ ,  $\mathbb{E}\|\mathbf{z}_i\|^4 < \infty$ ,  $\mathbb{E}\|\mathbf{x}_i\|^4 < \infty$ ,  $\mathbf{A}'\mathbb{E}(\mathbf{z}_i\mathbf{z}_i')\mathbf{A} > 0$ , and  $\mathbb{E}(\mathbf{u}_i\mathbf{u}_i') > 0$ , as  $n \rightarrow \infty$*

$$\sqrt{n} \begin{pmatrix} \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \\ \hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha} \end{pmatrix} \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{V}) \quad (11.58)$$

where

$$\mathbf{V} = \begin{pmatrix} \mathbf{V}_{\beta\beta} & \mathbf{V}_{\beta\alpha} \\ \mathbf{V}_{\alpha\beta} & \mathbf{V}_{\alpha\alpha} \end{pmatrix}$$

and

$$\begin{aligned} \mathbf{V}_{\beta\beta} &= (\mathbf{A}'\mathbb{E}(\mathbf{z}_i\mathbf{z}_i')\mathbf{A})^{-1} (\mathbf{A}'\mathbb{E}(\mathbf{z}_i\mathbf{z}_i'e_i^2)\mathbf{A}) (\mathbf{A}'\mathbb{E}(\mathbf{z}_i\mathbf{z}_i')\mathbf{A})^{-1} \\ \mathbf{V}_{\alpha\beta} &= (\mathbb{E}(\mathbf{u}_i\mathbf{u}_i'))^{-1} (\mathbb{E}(\mathbf{u}_i\mathbf{z}_i'e_i v_i)\mathbf{A}) (\mathbf{A}'\mathbb{E}(\mathbf{z}_i\mathbf{z}_i')\mathbf{A})^{-1} \\ \mathbf{V}_{\alpha\alpha} &= (\mathbb{E}(\mathbf{u}_i\mathbf{u}_i'))^{-1} \mathbb{E}(\mathbf{u}_i\mathbf{u}_i'v_i^2) (\mathbb{E}(\mathbf{u}_i\mathbf{u}_i'))^{-1}. \end{aligned}$$

The asymptotic covariance matrix is estimated by

$$\begin{aligned}\widehat{\mathbf{V}}_{\beta\beta} &= \left( \frac{1}{n} \widehat{\mathbf{W}}' \widehat{\mathbf{W}} \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \widehat{\mathbf{w}}_i \widehat{\mathbf{w}}_i' \widehat{e}_i^2 \right) \left( \frac{1}{n} \widehat{\mathbf{W}}' \widehat{\mathbf{W}} \right)^{-1} \\ \widehat{\mathbf{V}}_{\alpha\beta} &= \left( \frac{1}{n} \widehat{\mathbf{U}}' \widehat{\mathbf{U}} \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \widehat{\mathbf{u}}_i \widehat{\mathbf{w}}_i' \widehat{e}_i \widehat{v}_i \right) \left( \frac{1}{n} \widehat{\mathbf{W}}' \widehat{\mathbf{W}} \right)^{-1} \\ \widehat{\mathbf{V}}_{\alpha\alpha} &= \left( \frac{1}{n} \widehat{\mathbf{U}}' \widehat{\mathbf{U}} \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \widehat{\mathbf{u}}_i \widehat{\mathbf{u}}_i' \widehat{v}_i^2 \right) \left( \frac{1}{n} \widehat{\mathbf{U}}' \widehat{\mathbf{U}} \right)^{-1}\end{aligned}$$

where

$$\begin{aligned}\widehat{\mathbf{w}}_i &= \widehat{\mathbf{A}}' \mathbf{z}_i \\ \widehat{\mathbf{u}}_i &= \widehat{\mathbf{x}}_i - \widehat{\mathbf{w}}_i \\ \widehat{e}_i &= y_i - \mathbf{x}_i' \widehat{\boldsymbol{\beta}} \\ \widehat{v}_i &= y_i - \widehat{\mathbf{w}}_i' \widehat{\boldsymbol{\beta}} - \widehat{\mathbf{u}}_i' \widehat{\boldsymbol{\alpha}}.\end{aligned}$$

Under conditional homoskedasticity, specifically

$$\mathbb{E} \left( \begin{pmatrix} e_i^2 & e_i v_i \\ e_i v_i & v_i^2 \end{pmatrix} \middle| \mathbf{z}_i \right) = \mathbf{C}$$

then  $\mathbf{V}_{\alpha\beta} = \mathbf{0}$  and the coefficient estimates  $\widehat{\boldsymbol{\beta}}$  and  $\widehat{\boldsymbol{\alpha}}$  are asymptotically independent. The variance components also simplify to

$$\begin{aligned}\mathbf{V}_{\beta\beta} &= (\mathbf{A}' \mathbb{E}(\mathbf{z}_i \mathbf{z}_i') \mathbf{A})^{-1} \mathbb{E}(e_i^2) \\ \mathbf{V}_{\alpha\alpha} &= (\mathbb{E}(\mathbf{u}_i \mathbf{u}_i'))^{-1} \mathbb{E}(v_i^2).\end{aligned}$$

In this case we have the covariance matrix estimators

$$\begin{aligned}\widehat{\mathbf{V}}_{\beta\beta}^0 &= \left( \frac{1}{n} \widehat{\mathbf{W}}' \widehat{\mathbf{W}} \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \widehat{e}_i^2 \right) \\ \widehat{\mathbf{V}}_{\alpha\alpha}^0 &= \left( \frac{1}{n} \widehat{\mathbf{U}}' \widehat{\mathbf{U}} \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \widehat{v}_i^2 \right)\end{aligned}$$

and  $\widehat{\mathbf{V}}_{\alpha\beta}^0 = \mathbf{0}$ .

## 11.24 Control Function Regression

In this section we present an alternative way of computing the 2SLS estimator by least squares. It is useful in more complicated nonlinear contexts, and also in the linear model to construct tests for endogeneity.

The structural and reduced form equations for the standard IV model are

$$\begin{aligned}y_i &= \mathbf{x}_{1i}' \boldsymbol{\beta}_1 + \mathbf{x}_{2i}' \boldsymbol{\beta}_2 + e_i \\ \mathbf{x}_{2i} &= \boldsymbol{\Gamma}_{12}' \mathbf{z}_{1i} + \boldsymbol{\Gamma}_{22}' \mathbf{z}_{2i} + \mathbf{u}_{2i}.\end{aligned}$$

Since the instrumental variable assumption specifies that  $\mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0}$ ,  $\mathbf{x}_{2i}$  is endogenous (correlated with  $e_i$ ) if and only if  $\mathbf{u}_{2i}$  and  $e_i$  are correlated. We can therefore consider the linear projection of

$e_i$  on  $\mathbf{u}_{2i}$

$$\begin{aligned} e_i &= \mathbf{u}_{2i}' \boldsymbol{\alpha} + \varepsilon_i \\ \boldsymbol{\alpha} &= \left( \mathbb{E}(\mathbf{u}_{2i} \mathbf{u}_{2i}') \right)^{-1} \mathbb{E}(\mathbf{u}_{2i} e_i) \\ \mathbb{E}(\mathbf{u}_{2i} \varepsilon_i) &= \mathbf{0}. \end{aligned}$$

Substituting this into the structural form equation we find

$$\begin{aligned} y_i &= \mathbf{x}_{1i}' \boldsymbol{\beta}_1 + \mathbf{x}_{2i}' \boldsymbol{\beta}_2 + \mathbf{u}_{2i}' \boldsymbol{\alpha} + \varepsilon_i \\ \mathbb{E}(\mathbf{x}_{1i} \varepsilon_i) &= \mathbf{0} \\ \mathbb{E}(\mathbf{x}_{2i} \varepsilon_i) &= \mathbf{0} \\ \mathbb{E}(\mathbf{u}_{2i} \varepsilon_i) &= \mathbf{0}. \end{aligned} \tag{11.59}$$

Notice that  $\mathbf{x}_{2i}$  is uncorrelated with  $\varepsilon_i$ . This is because  $\mathbf{x}_{2i}$  is correlated with  $e_i$  only through  $\mathbf{u}_{2i}$ , and  $\varepsilon_i$  is the error after  $e_i$  has been projected orthogonal to  $\mathbf{u}_{2i}$ .

If  $\mathbf{u}_{2i}$  were observed we could then estimate (11.59) by least-squares. While it is not observed, we can estimate  $\mathbf{u}_{2i}$  by the reduced-form residual

$$\hat{\mathbf{u}}_{2i} = \mathbf{x}_{2i} - \hat{\boldsymbol{\Gamma}}_{12}' \mathbf{z}_{1i} - \hat{\boldsymbol{\Gamma}}_{22}' \mathbf{z}_{2i}$$

as defined in (11.23). Then the coefficients  $(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \boldsymbol{\alpha})$  can be estimated by least-squares of  $y_i$  on  $(\mathbf{x}_{1i}, \mathbf{x}_{2i}, \hat{\mathbf{u}}_{2i})$ . We can write this as

$$y_i = \mathbf{x}_{1i}' \hat{\boldsymbol{\beta}} + \hat{\mathbf{u}}_{2i}' \hat{\boldsymbol{\alpha}} + \hat{\varepsilon}_i \tag{11.60}$$

or in matrix notation as

$$\mathbf{y} = \mathbf{X} \hat{\boldsymbol{\beta}} + \hat{\mathbf{U}}_2 \hat{\boldsymbol{\alpha}} + \hat{\boldsymbol{\varepsilon}}.$$

This turns out to be an alternative algebraic expression for the 2SLS estimator.

Indeed, we now show that  $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_{2\text{sls}}$ . First, note that the reduced form residual can be written as

$$\hat{\mathbf{U}}_2 = (\mathbf{I}_n - \mathbf{P}_Z) \mathbf{X}_2$$

where  $\mathbf{P}_Z$  is defined in (11.35). By the FWL representation

$$\hat{\boldsymbol{\beta}} = \left( \widetilde{\mathbf{X}}' \widetilde{\mathbf{X}} \right)^{-1} \left( \widetilde{\mathbf{X}}' \mathbf{y} \right) \tag{11.61}$$

where  $\widetilde{\mathbf{X}} = [\widetilde{\mathbf{X}}_1, \widetilde{\mathbf{X}}_2]$ , with

$$\widetilde{\mathbf{X}}_1 = \mathbf{X}_1 - \hat{\mathbf{U}}_2 \left( \hat{\mathbf{U}}_2' \hat{\mathbf{U}}_2 \right)^{-1} \hat{\mathbf{U}}_2' \mathbf{X}_1 = \mathbf{X}_1$$

(since  $\hat{\mathbf{U}}_2' \mathbf{X}_1 = 0$ ) and

$$\begin{aligned} \widetilde{\mathbf{X}}_2 &= \mathbf{X}_2 - \hat{\mathbf{U}}_2 \left( \hat{\mathbf{U}}_2' \hat{\mathbf{U}}_2 \right)^{-1} \hat{\mathbf{U}}_2' \mathbf{X}_2 \\ &= \mathbf{X}_2 - \hat{\mathbf{U}}_2 \left( \mathbf{X}_2' (\mathbf{I}_n - \mathbf{P}_Z) \mathbf{X}_2 \right)^{-1} \mathbf{X}_2' (\mathbf{I}_n - \mathbf{P}_Z) \mathbf{X}_2 \\ &= \mathbf{X}_2 - \hat{\mathbf{U}}_2 \\ &= \mathbf{P}_Z \mathbf{X}_2. \end{aligned}$$

Thus  $\widetilde{\mathbf{X}} = [\mathbf{X}_1, \mathbf{P}_Z \mathbf{X}_2] = \mathbf{P}_Z \mathbf{X}$ . Substituted into (11.61) we find

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}' \mathbf{P}_Z \mathbf{X})^{-1} (\mathbf{X}' \mathbf{P}_Z \mathbf{y}) = \hat{\boldsymbol{\beta}}_{2\text{sls}}$$

which is (11.36) as claimed.

Again, what we have found is that OLS estimation of equation (11.60) yields algebraically the 2SLS estimator  $\hat{\beta}_{2\text{sls}}$ .

We now consider the distribution of the control function estimates. It is a generated regression model, and in fact is covered by the model examined in Section 11.23 after a slight reparametrization. Let  $\mathbf{w}_i = \mathbf{\Gamma}' \mathbf{z}_i$  and  $\mathbf{u}_i = \mathbf{x}_i - \mathbf{\Gamma}' \mathbf{z}_i = (\mathbf{0}', \mathbf{u}_{2i}')'$ . Then the main equation (11.59) can be written as

$$y_i = \mathbf{w}_i' \boldsymbol{\beta} + \mathbf{u}_{2i}' \boldsymbol{\gamma} + \varepsilon_i$$

where  $\boldsymbol{\gamma} = \boldsymbol{\alpha} + \boldsymbol{\beta}_2$ . This is the model in Section 11.23.

Set  $\hat{\boldsymbol{\gamma}} = \hat{\boldsymbol{\alpha}} + \hat{\boldsymbol{\beta}}_2$ . It follows from (11.58) that as  $n \rightarrow \infty$  we have the joint distribution

$$\sqrt{n} \begin{pmatrix} \hat{\boldsymbol{\beta}}_2 - \boldsymbol{\beta}_2 \\ \hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma} \end{pmatrix} \xrightarrow{d} \text{N}(\mathbf{0}, \mathbf{V})$$

where

$$\mathbf{V} = \begin{pmatrix} \mathbf{V}_{22} & \mathbf{V}_{2\gamma} \\ \mathbf{V}_{\gamma 2} & \mathbf{V}_{\gamma\gamma} \end{pmatrix}$$

$$\begin{aligned} \mathbf{V}_{22} &= \left[ (\mathbf{\Gamma}' \mathbb{E}(\mathbf{z}_i \mathbf{z}_i') \mathbf{\Gamma})^{-1} (\mathbf{\Gamma}' \mathbb{E}(\mathbf{z}_i \mathbf{z}_i' e_i^2 \mathbf{\Gamma})) (\mathbf{\Gamma}' \mathbb{E}(\mathbf{z}_i \mathbf{z}_i') \mathbf{\Gamma})^{-1} \right]_{22} \\ \mathbf{V}_{\gamma 2} &= \left[ (\mathbb{E}(\mathbf{u}_{2i} \mathbf{u}_{2i}'))^{-1} (\mathbb{E}(\mathbf{u}_{2i} \mathbf{z}_i' e_i \varepsilon_i) \mathbf{\Gamma}) (\mathbf{\Gamma}' \mathbb{E}(\mathbf{z}_i \mathbf{z}_i') \mathbf{\Gamma})^{-1} \right]_{.2} \\ \mathbf{V}_{\gamma\gamma} &= (\mathbb{E}(\mathbf{u}_{2i} \mathbf{u}_{2i}'))^{-1} \mathbb{E}(\mathbf{u}_{2i} \mathbf{u}_{2i}' \varepsilon_i^2) (\mathbb{E}(\mathbf{u}_{2i} \mathbf{u}_{2i}'))^{-1} \\ e_i &= y_i - \mathbf{x}_i' \boldsymbol{\beta}. \end{aligned}$$

The asymptotic distribution of  $\hat{\boldsymbol{\gamma}} = \hat{\boldsymbol{\alpha}} - \hat{\boldsymbol{\beta}}_2$  can then be deduced.

**Theorem 11.24.1** *If  $\mathbb{E}(y_i^4) < \infty$ ,  $\mathbb{E}\|\mathbf{z}_i\|^4 < \infty$ ,  $\mathbb{E}\|\mathbf{x}_i\|^4 < \infty$ ,  $\mathbf{A}' \mathbb{E}(\mathbf{z}_i \mathbf{z}_i') \mathbf{A} > 0$ , and  $\mathbb{E}(\mathbf{u}_i \mathbf{u}_i') > 0$ , as  $n \rightarrow \infty$*

$$\sqrt{n}(\hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}) \xrightarrow{d} \text{N}(\mathbf{0}, \mathbf{V}_{\boldsymbol{\alpha}})$$

where

$$\mathbf{V}_{\boldsymbol{\alpha}} = \mathbf{V}_{22} + \mathbf{V}_{\gamma\gamma} - \mathbf{V}_{2\gamma} - \mathbf{V}_{\gamma 2}.$$

$$\sqrt{n}(\hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}) \xrightarrow{d} \text{N}(\mathbf{0}, \mathbf{V}_{\boldsymbol{\alpha}})$$

where

$$\mathbf{V}_{\boldsymbol{\alpha}} = \mathbf{V}_{22} + \mathbf{V}_{\gamma\gamma} - \mathbf{V}_{2\gamma} - \mathbf{V}_{\gamma 2}.$$

Under conditional homoskedasticity we have the important simplifications

$$\begin{aligned} \mathbf{V}_{22} &= \left[ (\mathbf{\Gamma}' \mathbb{E}(\mathbf{z}_i \mathbf{z}_i') \mathbf{\Gamma})^{-1} \right]_{22} \mathbb{E}(e_i^2) \\ \mathbf{V}_{\gamma\gamma} &= (\mathbb{E}(\mathbf{u}_{2i} \mathbf{u}_{2i}'))^{-1} \mathbb{E}(\varepsilon_i^2) \\ \mathbf{V}_{\gamma 2} &= \mathbf{0} \\ \mathbf{V}_{\boldsymbol{\alpha}} &= \mathbf{V}_{22} + \mathbf{V}_{\gamma\gamma}. \end{aligned}$$

An estimator for  $\mathbf{V}_{\boldsymbol{\alpha}}$  in the general case is

$$\hat{\mathbf{V}}_{\boldsymbol{\alpha}} = \hat{\mathbf{V}}_{22} + \hat{\mathbf{V}}_{\gamma\gamma} - \hat{\mathbf{V}}_{2\gamma} - \hat{\mathbf{V}}_{\gamma 2} \quad (11.62)$$

where

$$\begin{aligned}\widehat{\mathbf{V}}_{22} &= \left[ \frac{1}{n} (\mathbf{X}'\mathbf{P}_Z\mathbf{X})^{-1} \mathbf{X}'\mathbf{Z} (\mathbf{Z}'\mathbf{Z})^{-1} \left( \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i' \widehat{e}_i^2 \right) (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{X} (\mathbf{X}'\mathbf{P}_Z\mathbf{X})^{-1} \right]_{22} \\ \widehat{\mathbf{V}}_{\gamma 2} &= \left[ \frac{1}{n} (\widehat{\mathbf{U}}'\widehat{\mathbf{U}})^{-1} \left( \sum_{i=1}^n \widehat{\mathbf{u}}_i \widehat{\mathbf{w}}_i' \widehat{e}_i \widehat{\varepsilon}_i \right) (\mathbf{X}'\mathbf{P}_Z\mathbf{X})^{-1} \right]_{.2} \\ \widehat{e}_i &= y_i - \mathbf{x}_i' \widehat{\boldsymbol{\beta}} \\ \widehat{\varepsilon}_i &= y_i - \mathbf{x}_i' \widehat{\boldsymbol{\beta}} - \widehat{\mathbf{u}}_{2i}' \widehat{\boldsymbol{\alpha}}.\end{aligned}$$

Under the assumption of conditional homoskedasticity we have the estimator

$$\begin{aligned}\widehat{\mathbf{V}}_{\alpha}^0 &= \widehat{\mathbf{V}}_{\beta\beta}^0 + \widehat{\mathbf{V}}_{\gamma\gamma}^0 \\ \widehat{\mathbf{V}}_{\beta\beta} &= \left[ (\mathbf{X}'\mathbf{P}_Z\mathbf{X})^{-1} \right]_{22} \left( \sum_{i=1}^n \widehat{e}_i^2 \right) \\ \widehat{\mathbf{V}}_{\gamma\gamma} &= (\widehat{\mathbf{U}}'\widehat{\mathbf{U}})^{-1} \left( \sum_{i=1}^n \widehat{\varepsilon}_i^2 \right).\end{aligned}$$

## 11.25 Endogeneity Tests

The 2SLS estimator allows the regressor  $\mathbf{x}_{2i}$  to be endogenous, meaning that  $\mathbf{x}_{2i}$  is correlated with the structural error  $e_i$ . If this correlation is zero, then  $\mathbf{x}_{2i}$  is exogenous and the structural equation can be estimated by least-squares. This is a testable restriction. Effectively, the null hypothesis is

$$\mathbb{H}_0 : \mathbb{E}(\mathbf{x}_{2i}e_i) = \mathbf{0}$$

with the alternative

$$\mathbb{H}_1 : \mathbb{E}(\mathbf{x}_{2i}e_i) \neq \mathbf{0}.$$

The maintained hypothesis is  $\mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0}$ . Since  $\mathbf{x}_{1i}$  is a component of  $\mathbf{z}_i$ , this implies  $\mathbb{E}(\mathbf{x}_{1i}e_i) = \mathbf{0}$ . Consequently we could alternatively write the null as  $\mathbb{H}_0 : \mathbb{E}(\mathbf{x}_i e_i) = \mathbf{0}$  (and some authors do so).

Recall the control function regression (11.59)

$$\begin{aligned}y_i &= \mathbf{x}_{1i}'\boldsymbol{\beta}_1 + \mathbf{x}_{2i}'\boldsymbol{\beta}_2 + \mathbf{u}_{2i}'\boldsymbol{\alpha} + \varepsilon_i \\ \boldsymbol{\alpha} &= (\mathbb{E}(\mathbf{u}_{2i}\mathbf{u}_{2i}'))^{-1} \mathbb{E}(\mathbf{u}_{2i}e_i).\end{aligned}$$

Notice that  $\mathbb{E}(\mathbf{x}_{2i}e_i) = \mathbf{0}$  if and only if  $\mathbb{E}(\mathbf{u}_{2i}e_i) = \mathbf{0}$ , so the hypothesis can be restated as  $\mathbb{H}_0 : \boldsymbol{\alpha} = \mathbf{0}$  against  $\mathbb{H}_1 : \boldsymbol{\alpha} \neq \mathbf{0}$ . Thus a natural test is based on the Wald statistic  $W$  for  $\boldsymbol{\alpha} = \mathbf{0}$  in the control function regression (11.24). Under Theorem 11.22.1 and Theorem 11.22.2, under  $\mathbb{H}_0$ ,  $W$  is asymptotically chi-square with  $k_2$  degrees of freedom. In addition, under the normal regression assumptions the  $F$  statistic has an exact  $F(k_2, n - k_1 - 2k_2)$  distribution. We accept the null hypothesis that  $\mathbf{x}_{2i}$  is exogenous if  $W$  (or  $F$ ) is smaller than the critical value, and reject in favor of the hypothesis that  $\mathbf{x}_{2i}$  is endogenous if the statistic is larger than the critical value.

Specifically, estimate the reduced form by least squares

$$\mathbf{x}_{2i} = \widehat{\boldsymbol{\Gamma}}'_{12} \mathbf{z}_{1i} + \widehat{\boldsymbol{\Gamma}}'_{22} \mathbf{z}_{2i} + \widehat{\mathbf{u}}_{2i}$$

to obtain the residuals. Then estimate the control function by least squares

$$y_i = \mathbf{x}_i' \widehat{\boldsymbol{\beta}} + \widehat{\mathbf{u}}_{2i}' \widehat{\boldsymbol{\alpha}} + \widehat{\varepsilon}_i. \quad (11.63)$$

Let  $W$ ,  $W^0$  and  $F = W^0/k_2$  denote the Wald statistic, homoskedastic Wald statistic, and  $F$  statistic for  $\boldsymbol{\alpha} = \mathbf{0}$ .

**Theorem 11.25.1** Under  $\mathbb{H}_0$ ,  $W \xrightarrow{d} \chi_{k_2}^2$ . Let  $c_{1-\alpha}$  solve  $\Pr(\chi_{k_2}^2 \leq c_{1-\alpha}) = 1 - \alpha$ . The test “Reject  $\mathbb{H}_0$  if  $W > c_{1-\alpha}$ ” has asymptotic size  $\alpha$ .

**Theorem 11.25.2** Suppose  $e_i | \mathbf{x}_i, \mathbf{z}_i \sim N(0, \sigma^2)$ . Under  $\mathbb{H}_0$ ,  $F \sim F(k_2, n - k_1 - 2k_2)$ . Let  $c_{1-\alpha}$  solve  $\Pr(F(k_2, n - k_1 - 2k_2) \leq c_{1-\alpha}) = 1 - \alpha$ . The test “Reject  $\mathbb{H}_0$  if  $F > c_{1-\alpha}$ ” has exact size  $\alpha$ .

Since in general we do not want to impose homoskedasticity, these results suggest that the most appropriate test is the Wald statistic constructed with the robust heteroskedastic covariance matrix. This can be computed in Stata using the command `estat endogenous` after `ivregress` when the latter uses a robust covariance option. Stata reports the Wald statistic in  $F$  form (and thus uses the  $F$  distribution to calculate the p-value) as “Robust regression F”. Using the  $F$  rather than the  $\chi^2$  distribution is not formally justified but is a reasonable finite sample adjustment. If the command `estat endogenous` is applied after `ivregress` without a robust covariance option, Stata reports the  $F$  statistic as “Wu-Hausman F”.

There is an alternative (and traditional) way to derive a test for endogeneity. Under  $\mathbb{H}_0$ , both OLS and 2SLS are consistent estimators. But under  $\mathbb{H}_1$ , they converge to different values. Thus the difference between the OLS and 2SLS estimators is a valid test statistic for endogeneity. It also measures what we often care most about – the impact of endogeneity on the parameter estimates. This literature was developed under the assumption of conditional homoskedasticity (and it is important for these results) so we assume this condition for the development of the statistics.

Let  $\hat{\beta} = (\hat{\beta}_1, \hat{\beta}_2)$  be the OLS estimator and let  $\tilde{\beta} = (\tilde{\beta}_1, \tilde{\beta}_2)$  be the 2SLS estimator. Under  $\mathbb{H}_0$  (and homoskedasticity) the OLS estimator is Gauss-Markov efficient, so by the Hausman equality

$$\begin{aligned} \text{var}(\hat{\beta}_2 - \tilde{\beta}_2) &= \text{var}(\tilde{\beta}_2) - \text{var}(\hat{\beta}_2) \\ &= \left( (\mathbf{X}_2'(\mathbf{P}_Z - \mathbf{P}_1)\mathbf{X}_2)^{-1} - (\mathbf{X}_2'\mathbf{M}_1\mathbf{X}_2)^{-1} \right) \sigma^2 \end{aligned}$$

where  $\mathbf{P}_Z = \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'$ ,  $\mathbf{P}_1 = \mathbf{X}_1(\mathbf{X}_1'\mathbf{X}_1)^{-1}\mathbf{X}_1'$ , and  $\mathbf{M}_1 = \mathbf{I}_n - \mathbf{P}_1$ . Thus a valid test statistic for  $\mathbb{H}_0$  is

$$T = \frac{(\hat{\beta}_2 - \tilde{\beta}_2)' \left( (\mathbf{X}_2'(\mathbf{P}_Z - \mathbf{P}_1)\mathbf{X}_2)^{-1} - (\mathbf{X}_2'\mathbf{M}_1\mathbf{X}_2)^{-1} \right)^{-1} (\hat{\beta}_2 - \tilde{\beta}_2)}{\hat{\sigma}^2} \quad (11.64)$$

for some estimate  $\hat{\sigma}^2$  of  $\sigma^2$ . Durbin (1954) first proposed  $T$  as a test for endogeneity in the context of IV estimation, setting  $\hat{\sigma}^2$  to be the least-squares estimate of  $\sigma^2$ . Wu (1973) proposed  $T$  as a test for endogeneity in the context of 2SLS estimation, considering a set of possible estimates  $\hat{\sigma}^2$ , including the regression estimate from (11.63). Hausman (1978) proposed a version of  $T$  based on the full contrast  $\hat{\beta} - \tilde{\beta}$ , and observed that it equals the regression Wald statistic  $W^0$  described earlier. In fact, when  $\hat{\sigma}^2$  is the regression estimate from (11.63), the statistic (11.64) algebraically equals both  $W^0$  and the version of (11.64) based on the full contrast  $\hat{\beta} - \tilde{\beta}$ . We show these equalities below. Thus these three approaches yield exactly the same statistic except for possible differences regarding the choice of  $\hat{\sigma}^2$ . Since the regression  $F$  test described earlier has an exact  $F$  distribution in the normal sampling model, and thus can exactly control test size, this is the

preferred version of the test. The general class of tests are called **Durbin-Wu-Hausman** tests, **Wu-Hausman** tests, or **Hausman** tests, depending on the author.

When  $k_2 = 1$  (there is one right-hand-side endogenous variable) which is quite common in applications, the endogeneity test can be equivalently expressed at the t-statistic for  $\hat{\alpha}$  in the estimated control function. Thus it is sufficient to estimate the control function regression and check the t-statistic for  $\hat{\alpha}$ . If  $|\hat{\alpha}| > 2$  then we can reject the hypothesis that  $\mathbf{x}_{2i}$  is exogenous for  $\beta$ .

We illustrate using the Card proximity example using the two instruments *public* and *private*. We first estimate the reduced form for *education*, obtain the residual, and then estimate the control function regression. The residual has a coefficient  $-0.088$  with a standard error of  $0.037$  and a t-statistic of  $2.4$ . Since the latter exceeds the 5% critical value (its p-value is  $0.017$ ) we reject exogeneity. This means that the 2SLS estimates are statistically different from the least-squares estimates of the structural equation and supports our decision to treat education as an endogenous variable. (Alternatively, the  $F$  statistic is  $2.4^2 = 5.7$  with the same p-value).

We now show the equality of the various statistics.

We first show that the statistic (11.64) is not altered if based on the full contrast  $\hat{\beta} - \tilde{\beta}$ . Indeed,  $\hat{\beta}_1 - \tilde{\beta}_1$  is a linear function of  $\hat{\beta}_2 - \tilde{\beta}_2$ , so there is no extra information in the full contrast. To see this, observe that given  $\hat{\beta}_2$ , we can solve by least-squares to find

$$\hat{\beta}_1 = (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \left( \mathbf{X}'_1 (\mathbf{y} - \mathbf{X}_2 \hat{\beta}_2) \right)$$

and similarly

$$\begin{aligned} \tilde{\beta}_1 &= (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \left( \mathbf{X}'_1 (\mathbf{y} - \mathbf{P}_Z \mathbf{X}_2 \tilde{\beta}) \right) \\ &= (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \left( \mathbf{X}'_1 (\mathbf{y} - \mathbf{X}_2 \tilde{\beta}) \right) \end{aligned}$$

the second equality since  $\mathbf{P}_Z \mathbf{X}_1 = \mathbf{X}_1$ . Thus

$$\begin{aligned} \hat{\beta}_1 - \tilde{\beta}_1 &= (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 (\mathbf{y} - \mathbf{X}_2 \hat{\beta}_2) - (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 (\mathbf{y} - \mathbf{P}_Z \mathbf{X}_2 \tilde{\beta}) \\ &= (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{X}_2 (\tilde{\beta}_2 - \hat{\beta}_2) \end{aligned}$$

as claimed.

We next show that  $T$  in (11.64) equals the homoskedastic Wald statistic  $W^0$  for  $\hat{\alpha}$  from the regression (11.63). Consider the latter regression. Since  $\mathbf{X}_2$  is contained in  $\mathbf{X}$ , the coefficient estimate  $\hat{\alpha}$  is invariant to replacing  $\hat{\mathbf{U}}_2 = \mathbf{X}_2 - \hat{\mathbf{X}}_2$  with  $-\hat{\mathbf{X}}_2 = -\mathbf{P}_Z \mathbf{X}_2$ . By the FWL representation, setting  $\mathbf{M}_X = \mathbf{I}_n - \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'$

$$\begin{aligned} \hat{\alpha} &= - \left( \hat{\mathbf{X}}_2' \mathbf{M}_X \hat{\mathbf{X}}_2 \right)^{-1} \hat{\mathbf{X}}_2' \mathbf{M}_X \mathbf{y} \\ &= - \left( \mathbf{X}_2' \mathbf{P}_Z \mathbf{M}_X \mathbf{P}_Z \mathbf{X}_2 \right)^{-1} \mathbf{X}_2' \mathbf{P}_Z \mathbf{M}_X \mathbf{y}. \end{aligned} \tag{11.65}$$

It follows that

$$W^0 = \frac{\mathbf{y}' \mathbf{M}_X \mathbf{P}_Z \mathbf{X}_2 (\mathbf{X}_2' \mathbf{P}_Z \mathbf{M}_X \mathbf{P}_Z \mathbf{X}_2)^{-1} \mathbf{X}_2' \mathbf{P}_Z \mathbf{M}_X \mathbf{y}}{\hat{\sigma}^2}.$$

Our goal is to show that  $T = W^0$ . Define  $\tilde{\mathbf{X}}_2 = (\mathbf{I}_n - \mathbf{P}_1) \mathbf{X}_2$  so  $\tilde{\beta}_2 = \left( \tilde{\mathbf{X}}_2' \tilde{\mathbf{X}}_2 \right)^{-1} \tilde{\mathbf{X}}_2' \mathbf{y}$ . Then

defining using  $(\mathbf{P}_Z - \mathbf{P}_1)(\mathbf{I}_n - \mathbf{P}_1) = (\mathbf{P}_Z - \mathbf{P}_1)$  and defining  $\mathbf{Q} = \widetilde{\mathbf{X}}_2 \left( \widetilde{\mathbf{X}}_2' \widetilde{\mathbf{X}}_2 \right)^{-1} \widetilde{\mathbf{X}}_2'$

$$\begin{aligned}
\Delta &\stackrel{def}{=} (\mathbf{X}_2' (\mathbf{P}_Z - \mathbf{P}_1) \mathbf{X}_2) (\widetilde{\beta}_2 - \widehat{\beta}_2) \\
&= \mathbf{X}_2' (\mathbf{P}_Z - \mathbf{P}_1) \mathbf{y} - (\mathbf{X}_2' (\mathbf{P}_Z - \mathbf{P}_1) \mathbf{X}_2) \left( \widetilde{\mathbf{X}}_2' \widetilde{\mathbf{X}}_2 \right)^{-1} \widetilde{\mathbf{X}}_2' \mathbf{y} \\
&= \mathbf{X}_2' (\mathbf{P}_Z - \mathbf{P}_1) (\mathbf{I}_n - \mathbf{Q}) \mathbf{y} \\
&= \mathbf{X}_2' (\mathbf{P}_Z - \mathbf{P}_1 - \mathbf{P}_Z \mathbf{Q}) \mathbf{y} \\
&= \mathbf{X}_2' \mathbf{P}_Z (\mathbf{I}_n - \mathbf{P}_1 - \mathbf{Q}) \mathbf{y} \\
&= \mathbf{X}_2' \mathbf{P}_Z \mathbf{M}_X \mathbf{y}.
\end{aligned}$$

The third-to-last equality is  $\mathbf{P}_1 \mathbf{Q} = \mathbf{0}$  and the final uses  $\mathbf{M}_X = \mathbf{I}_n - \mathbf{P}_1 - \mathbf{Q}$ . We also calculate that

$$\begin{aligned}
\mathbf{Q}^* &\stackrel{def}{=} (\mathbf{X}_2' (\mathbf{P}_Z - \mathbf{P}_1) \mathbf{X}_2) \left( (\mathbf{X}_2' (\mathbf{P}_Z - \mathbf{P}_1) \mathbf{X}_2)^{-1} - (\mathbf{X}_2' \mathbf{M}_1 \mathbf{X}_2)^{-1} \right) \\
&\cdot (\mathbf{X}_2' (\mathbf{P}_Z - \mathbf{P}_1) \mathbf{X}_2) \\
&= \mathbf{X}_2' (\mathbf{P}_Z - \mathbf{P}_1 - (\mathbf{P}_Z - \mathbf{P}_1) \mathbf{Q} (\mathbf{P}_Z - \mathbf{P}_1)) \mathbf{X}_2 \\
&= \mathbf{X}_2' (\mathbf{P}_Z - \mathbf{P}_1 - \mathbf{P}_Z \mathbf{Q} \mathbf{P}_Z) \mathbf{X}_2 \\
&= \mathbf{X}_2' \mathbf{P}_Z \mathbf{M}_X \mathbf{P}_Z \mathbf{X}_2.
\end{aligned}$$

Thus

$$\begin{aligned}
T &= \frac{\Delta' \mathbf{Q}^{*-1} \Delta}{\widehat{\sigma}^2} \\
&= \frac{\mathbf{y}' \mathbf{M}_X \mathbf{P}_Z \mathbf{X}_2 (\mathbf{X}_2' \mathbf{P}_Z \mathbf{M}_X \mathbf{P}_Z \mathbf{X}_2)^{-1} \mathbf{X}_2' \mathbf{P}_Z \mathbf{M}_X \mathbf{y}}{\widehat{\sigma}^2} \\
&= W^0
\end{aligned}$$

as claimed.

## 11.26 Subset Endogeneity Tests

In some cases we may only wish to test the endogeneity of a subset of the variables. In the Card proximity example, we may wish test the exogeneity of *education* separately from *experience* and its square. To execute a subset endogeneity test it is useful to partition the regressors into three groups, so that the structural model is

$$\begin{aligned}
y_i &= \mathbf{x}_{1i}' \beta_1 + \mathbf{x}_{2i}' \beta_2 + \mathbf{x}_{3i}' \beta_3 + e_i \\
\mathbb{E}(\mathbf{z}_i e_i) &= \mathbf{0}.
\end{aligned}$$

As before, the instrument vector  $\mathbf{z}_i$  includes  $\mathbf{x}_{1i}$ . The variables  $\mathbf{x}_{3i}$  is treated as endogenous, and  $\mathbf{x}_{2i}$  is treated as potentially endogenous. The hypothesis to test is that  $\mathbf{x}_{2i}$  is exogenous, or

$$\mathbb{H}_0 : \mathbb{E}(\mathbf{x}_{2i} e_i) = \mathbf{0}$$

against

$$\mathbb{H}_1 : \mathbb{E}(\mathbf{x}_{2i} e_i) \neq \mathbf{0}.$$

Under homoskedasticity, a straightforward test can be constructed by the Durbin-Wu-Hausman principle. Under  $\mathbb{H}_0$ , the appropriate estimator is 2SLS using the instruments  $(\mathbf{z}_i, \mathbf{x}_{2i})$ . Let this estimator of  $\beta_2$  be denoted  $\widehat{\beta}_2$ . Under  $\mathbb{H}_1$ , the appropriate estimator is 2SLS using the smaller



instrument set  $\mathbf{z}_i$ . Let this estimator of  $\beta_2$  be denoted  $\tilde{\beta}_2$ . A Durbin-Wu-Hausman-type test of  $\mathbb{H}_0$  against  $\mathbb{H}_1$  is

$$T = \left( \hat{\beta}_2 - \tilde{\beta}_2 \right)' \left( \widehat{\text{var}} \left( \tilde{\beta}_2 \right) - \widehat{\text{var}} \left( \hat{\beta}_2 \right) \right)^{-1} \left( \hat{\beta}_2 - \tilde{\beta}_2 \right).$$

The asymptotic distribution under  $\mathbb{H}_0$  is  $\chi_{k_2}^2$  where  $k_2 = \dim(\mathbf{x}_{2i})$ , so we reject the hypothesis that the variables  $\mathbf{x}_{2i}$  are exogenous if  $T$  exceeds an upper critical value from the  $\chi_{k_2}^2$  distribution.

Instead of using the Wald statistic, one could use the  $F$  version of the test by dividing by  $k_2$  and using the  $F$  distribution for critical values. There is no finite sample justification for this modification, however, since  $\mathbf{x}_{3i}$  is endogenous under the null hypothesis.

In Stata, the command `estat endogenous` (adding the variable name to specify which variable to test for exogeneity) after `ivregress` without a robust covariance option reports the  $F$  version of this statistic as “Wu-Hausman F”. For example, in the Card proximity example using the four instruments *public*, *private*, *age* and *age*<sup>2</sup>, if we estimate the equation by 2SLS with a non-robust covariance matrix, and then compute the endogeneity test for education, we find  $F = 272$  with a p-value of 0.0000, but if we compute the test for experience and its square we find  $F = 2.98$  with a p-value of 0.051. In this equation, education is clearly endogenous but the experience variables are unclear.

A heteroskedasticity or cluster-robust test cannot be constructed easily by the Durbin-Wu-Hausman approach, since the covariance matrix does not take a simple form. Instead, we can use the regression approach if we account for the generated regressor problem. The ideal control function regression takes the form

$$y_i = \mathbf{x}'_i \beta + \mathbf{u}'_{2i} \alpha_2 + \mathbf{u}'_{3i} \alpha_3 + \varepsilon_i$$

where  $\mathbf{u}_{2i}$  and  $\mathbf{u}_{3i}$  are the reduced-form errors from the projections of  $\mathbf{x}_{2i}$  and  $\mathbf{x}_{3i}$  on the instruments  $\mathbf{z}_i$ . The coefficients  $\alpha_2$  and  $\alpha_3$  solve the equations

$$\begin{pmatrix} \mathbb{E}(\mathbf{u}_{2i} \mathbf{u}'_{2i}) & \mathbb{E}(\mathbf{u}_{2i} \mathbf{u}'_{3i}) \\ \mathbb{E}(\mathbf{u}_{3i} \mathbf{u}'_{2i}) & \mathbb{E}(\mathbf{u}_{3i} \mathbf{u}'_{3i}) \end{pmatrix} \begin{pmatrix} \alpha_2 \\ \alpha_3 \end{pmatrix} = \begin{pmatrix} \mathbb{E}(\mathbf{u}_{2i} \varepsilon_i) \\ \mathbb{E}(\mathbf{u}_{3i} \varepsilon_i) \end{pmatrix}.$$

The null hypothesis  $\mathbb{E}(\mathbf{x}_{2i} \varepsilon_i) = \mathbf{0}$  is equivalent to  $\mathbb{E}(\mathbf{u}_{2i} \varepsilon_i) = \mathbf{0}$ . This implies

$$\Psi' \begin{pmatrix} \alpha_2 \\ \alpha_3 \end{pmatrix} = \mathbf{0} \quad (11.66)$$

where

$$\Psi = \begin{pmatrix} \mathbb{E}(\mathbf{u}_{2i} \mathbf{u}'_{2i}) \\ \mathbb{E}(\mathbf{u}_{3i} \mathbf{u}'_{2i}) \end{pmatrix}.$$

This suggests that an appropriate regression-based test of  $\mathbb{H}_0$  versus  $\mathbb{H}_1$  is to construct a Wald statistic for the restriction (11.66) in the control function regression

$$y_i = \mathbf{x}'_i \hat{\beta} + \hat{\mathbf{u}}'_{2i} \hat{\alpha}_2 + \hat{\mathbf{u}}'_{3i} \hat{\alpha}_3 + \hat{\varepsilon}_i \quad (11.67)$$

where  $\hat{\mathbf{u}}_{2i}$  and  $\hat{\mathbf{u}}_{3i}$  are the least-squares residuals from the regressions of  $\mathbf{x}_{2i}$  and  $\mathbf{x}_{3i}$  on the instruments  $\mathbf{z}_i$ , respectively, and  $\Psi$  is estimated by

$$\hat{\Psi} = \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{u}}_{2i} \hat{\mathbf{u}}'_{2i} \\ \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{u}}_{3i} \hat{\mathbf{u}}'_{2i} \end{pmatrix}.$$

A complication is that the regression (11.67) has generated regressors which have non-zero coefficients under  $\mathbb{H}_0$ . The solution is to use the control-function-robust covariance matrix estimator (11.62) for  $(\hat{\alpha}_2, \hat{\alpha}_3)$ . This yields a valid Wald statistic for  $\mathbb{H}_0$  versus  $\mathbb{H}_1$ . The asymptotic distribution of the statistic under  $\mathbb{H}_0$  is  $\chi_{k_2}^2$  where  $k_2 = \dim(\mathbf{x}_{2i})$ , so the null hypothesis that  $\mathbf{x}_{2i}$  is exogenous is rejected if the Wald statistic exceeds the upper critical value from the  $\chi_{k_2}^2$  distribution.

Heteroskedasticity-robust and cluster-robust subset endogeneity tests are not currently implemented in Stata.

## 11.27 OverIdentification Tests

When  $\ell > k$  the model is **overidentified** meaning that there are more moments than free parameters. This is a restriction and is testable. Such tests are called **overidentification tests**.

The instrumental variables model specifies that

$$\mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0}.$$

Equivalently, since  $e_i = y_i - \mathbf{x}_i' \boldsymbol{\beta}$ , this is the same as

$$\mathbb{E}(\mathbf{z}_i y_i) - \mathbb{E}(\mathbf{z}_i \mathbf{x}_i') \boldsymbol{\beta} = \mathbf{0}.$$

This is an  $\ell \times 1$  vector of restrictions on the moment matrices  $\mathbb{E}(\mathbf{z}_i y_i)$  and  $\mathbb{E}(\mathbf{z}_i \mathbf{x}_i')$ . Yet since  $\boldsymbol{\beta}$  is of dimension  $k$  which is less than  $\ell$ , it is not certain if indeed such a  $\boldsymbol{\beta}$  exists.

To make things a bit more concrete, suppose there is a single endogenous regressor  $x_{2i}$ , no  $x_{1i}$ , and two instruments  $z_{1i}$  and  $z_{2i}$ . Then the model specifies that

$$\mathbb{E}(z_{1i} y_i) = \mathbb{E}(z_{1i} x_{2i}) \boldsymbol{\beta}$$

and

$$\mathbb{E}(z_{2i} y_i) = \mathbb{E}(z_{2i} x_{2i}) \boldsymbol{\beta}.$$

Thus  $\boldsymbol{\beta}$  solves both equations. This is rather special.

Another way of thinking about this is that in this context we could solve for  $\boldsymbol{\beta}$  using either one equation or the other. In terms of estimation, this is equivalent to estimating by IV using just the instrument  $z_1$  or instead just using the instrument  $z_2$ . These two estimators (in finite samples) will be different. But if the overidentification hypothesis is correct, both are estimating the same parameter, and both are consistent for  $\boldsymbol{\beta}$  (if the instruments are relevant). In contrast, if the overidentification hypothesis is false, then the two estimators will converge to different probability limits and it is unclear if either probability limit is interesting.

For example, take the 2SLS estimates in the fourth column of Table 11.1, which use *public* and *private* as instruments for *education*. Suppose we instead estimate by IV, using just *public* as an instrument, and then repeat using *private*. The IV coefficient for *education* in the first case is 0.17, and in the second case 0.27. These appear to be quite different. However, the second estimate has quite a large standard error (0.17) so perhaps the difference is sampling variation. An overidentification test addresses this question formally.

For a general overidentification test, the null and alternative hypotheses are

$$\begin{aligned} \mathbb{H}_0 &: \mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0} \\ \mathbb{H}_1 &: \mathbb{E}(\mathbf{z}_i e_i) \neq \mathbf{0}. \end{aligned}$$

We will also add the conditional homoskedasticity assumption

$$\mathbb{E}(e_i^2 | \mathbf{z}_i) = \sigma^2. \quad (11.68)$$

To avoid imposing (11.68), it is best to take a GMM approach, which we defer until Chapter 12.

To implement a test of  $\mathbb{H}_0$ , consider a linear regression of the error  $e_i$  on the instruments  $\mathbf{z}_i$

$$e_i = \mathbf{z}_i' \boldsymbol{\alpha} + \varepsilon_i \quad (11.69)$$

with

$$\boldsymbol{\alpha} = (\mathbb{E}(\mathbf{z}_i \mathbf{z}_i'))^{-1} \mathbb{E}(\mathbf{z}_i e_i).$$

We can rewrite  $\mathbb{H}_0$  as  $\boldsymbol{\alpha} = \mathbf{0}$ . While  $e_i$  is not observed we can replace it with the 2SLS residual  $\hat{e}_i$ , and estimate  $\boldsymbol{\alpha}$  by least-squares regression

$$\hat{\boldsymbol{\alpha}} = (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \hat{\mathbf{e}}.$$

Sargan (1958) proposed testing  $\mathbb{H}_0$  via a score test, which takes the form

$$S = \hat{\alpha}' (\widehat{\text{var}}(\hat{\alpha}))^{-1} \hat{\alpha} = \frac{\hat{e}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \hat{e}}{\hat{\sigma}^2}. \quad (11.70)$$

where  $\hat{\sigma}^2 = \frac{1}{n} \hat{e}' \hat{e}$ . Basman (1960) independently proposed a Wald statistic for  $\mathbb{H}_0$ , which is  $S$  with  $\hat{\sigma}^2$  replaced with  $\tilde{\sigma}^2 = n^{-1} \hat{e}' \hat{\varepsilon}$  where  $\hat{\varepsilon} = \hat{e} - \mathbf{Z} \hat{\alpha}$ . By the equivalence of homoskedastic score and Wald tests (see Section 9.16), Basman's statistic is a monotonic function of Sargan's statistic and hence they yield equivalent tests. Sargan's version is more typically reported.

The Sargan test rejects  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  if  $S > c$  for some critical value  $c$ . An asymptotic test sets  $c$  as the  $1 - \alpha$  quantile of the  $\chi_{\ell-k}^2$  distribution. This is justified by the asymptotic null distribution of  $S$  which we now derive.

**Theorem 11.27.1** *Under Assumption 11.14.1 and  $\mathbb{E}(e_i^2 | \mathbf{z}_i) = \sigma^2$ , then as  $n \rightarrow \infty$*

$$S \xrightarrow{d} \chi_{\ell-k}^2.$$

*For  $c$  satisfying  $\alpha = 1 - G_{\ell-k}(c)$ ,*

$$\Pr(S > c \mid \mathbb{H}_0) \rightarrow \alpha$$

*so the test “Reject  $\mathbb{H}_0$  if  $S > c$ ” has asymptotic size  $\alpha$ .*

We prove Theorem 11.27.1 below.

The Sargan statistic  $S$  is an asymptotic test of the overidentifying restrictions under the assumption of conditional homoskedasticity. It has some limitations. First, it is an asymptotic test, and does not have a finite sample (e.g.  $F$ ) counterpart. Simulation evidence suggests that the test can be oversized (reject too frequently) in small and moderate sample sizes. Consequently, p-values should be interpreted cautiously. Second, the assumption of conditional homoskedasticity is unrealistic in applications. The best way to generalize the Sargan statistic to allow heteroskedasticity is to use the GMM overidentification statistic – which we will examine in Chapter 12. For 2SLS, Wooldridge (1995) suggested a robust score test, but Baum, Schaffer and Stillman (2003) point out that it is numerically equivalent to the GMM overidentification statistic. Hence the bottom line appears to be that to allow heteroskedasticity or clustering, it is best to use a GMM approach.

In overidentified applications, it is always prudent to report an overidentification test. If the test is insignificant it means that the overidentifying restrictions are not rejected, supporting the estimated model. If the overidentifying test statistic is highly significant (if the p-value is very small) this is evidence that the overidentifying restrictions are violated. In this case we should be concerned that the model is misspecified and interpreting the parameter estimates should be done cautiously.

When reporting the results of an overidentification test, it seems reasonable to focus on very small significance levels, such as 1%. This means that we should only treat a model as “rejected” if the Sargan p-value is very small, e.g. less than 0.01. The reason to focus on very small significance levels is because it is very difficult to interpret the result “The model is rejected”. Stepping back a bit, it does not seem credible that any overidentified model is literally true, rather what seems potentially credible is that an overidentified model is a reasonable approximation. A test is asking the question “Is there evidence that a model is not true” when we really want to know the answer to “Is there evidence that the model is a poor approximation”. Consequently it seems reasonable to require strong evidence to lead to the conclusion “Let's reject this model”. The recommendation is that mild rejections (p-values between 1% and 5%) should be viewed as mildly worrisome, but

not critical evidence against a model. The results of an overidentification test should be integrated with other information before making a strong decision.

We illustrate the methods with the Card college proximity example. We have estimated two overidentified models by 2SLS, in columns 4 & 5 of Table 11.1. In each case, the number of overidentifying restrictions is 1. We report the Sargan statistic and its asymptotic p-value (calculated using the  $\chi_1^2$  distribution) in the table. Both p-values (0.36 and 0.52) are far from significant, indicating that there is no evidence that the models are misspecified.

We now prove Theorem 11.27.1. The statistic  $S$  is invariant to rotations of  $\mathbf{Z}$  (replacing  $\mathbf{Z}$  with  $\mathbf{ZC}$ ) so without loss of generality we assume  $\mathbb{E}(\mathbf{z}_i \mathbf{z}_i') = \mathbf{I}_\ell$ . As  $n \rightarrow \infty$ ,  $n^{-1/2} \mathbf{Z}' \mathbf{e} \xrightarrow{d} \sigma \mathbf{Z}$  where  $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I}_\ell)$ . Also  $\frac{1}{n} \mathbf{Z}' \mathbf{Z} \xrightarrow{p} \mathbf{I}_\ell$  and  $\frac{1}{n} \mathbf{Z}' \mathbf{X} \xrightarrow{p} \mathbf{Q}$ , say. Then

$$\begin{aligned} n^{-1/2} \mathbf{Z}' \hat{\mathbf{e}} &= \left( \mathbf{I}_\ell - \left( \frac{1}{n} \mathbf{Z}' \mathbf{X} \right) \left( \frac{1}{n} \mathbf{X}' \mathbf{P}_Z \mathbf{X} \right)^{-1} \left( \frac{1}{n} \mathbf{X}' \mathbf{Z} \right) \left( \frac{1}{n} \mathbf{Z}' \mathbf{Z} \right)^{-1} \right) n^{-1/2} \mathbf{Z}' \mathbf{e} \\ &\xrightarrow{d} \sigma \left( \mathbf{I}_\ell - \mathbf{Q} (\mathbf{Q}' \mathbf{Q})^{-1} \mathbf{Q}' \right) \mathbf{Z}. \end{aligned}$$

Since  $\hat{\sigma}^2 \xrightarrow{p} \sigma^2$  it follows that

$$S \xrightarrow{d} \mathbf{Z}' \left( \mathbf{I}_\ell - \mathbf{Q} (\mathbf{Q}' \mathbf{Q})^{-1} \mathbf{Q}' \right) \mathbf{Z} \sim \chi_{\ell-k}^2.$$

The distribution is  $\chi_{\ell-k}^2$  since  $\mathbf{I}_\ell - \mathbf{Q} (\mathbf{Q}' \mathbf{Q})^{-1} \mathbf{Q}'$  is idempotent with rank  $\ell - k$ .

The Sargan statistic test can be implemented in Stata using the command `estat overid` after `ivregress 2sls` or `ivregress liml` if a standard (non-robust) covariance matrix has been specified (that is, without the `,r` option), or by the command `estat overid, forcenonrobust` otherwise.

## 11.28 Subset OverIdentification Tests

Tests of  $\mathbb{H}_0 : \mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0}$  are typically interpreted as tests of model specification. The alternative  $\mathbb{H}_1 : \mathbb{E}(\mathbf{z}_i e_i) \neq \mathbf{0}$  means that at least one element of  $\mathbf{z}_i$  is correlated with the error  $e_i$  and is thus an invalid instrumental variable. In some cases it may be reasonable to test only a subset of the moment conditions.

As in the previous section we restrict attention to the homoskedasticity case  $\mathbb{E}(e_i^2 | \mathbf{z}_i) = \sigma^2$ .

Partition  $\mathbf{z}_i = (\mathbf{z}_{ai}, \mathbf{z}_{bi})$  with dimensions  $\ell_a$  and  $\ell_b$ , respectively, where  $\mathbf{z}_{ai}$  contains the instruments which are believed to be uncorrelated with  $e_i$ , and  $\mathbf{z}_{bi}$  contains the instruments which may be correlated with  $e_i$ . It is necessary to select this partition so that  $\ell_a > k$ , or equivalently  $\ell_b < \ell - k$ . This means that the model with just the instruments  $\mathbf{z}_{ai}$  is over-identified, or that  $\ell_b$  is smaller than the number of overidentifying restrictions. (If  $\ell_a = k$  then the tests described here exist but reduce to the Sargan test so are not interesting.) Hence the tests require that  $\ell - k > 1$ , that the number of overidentifying restrictions exceeds one.

Given this partition, the maintained hypothesis is that  $\mathbb{E}(\mathbf{z}_{ai} e_i) = \mathbf{0}$ . The null and alternative hypotheses are

$$\begin{aligned} \mathbb{H}_0 : \mathbb{E}(\mathbf{z}_{bi} e_i) &= \mathbf{0} \\ \mathbb{H}_1 : \mathbb{E}(\mathbf{z}_{bi} e_i) &\neq \mathbf{0}. \end{aligned}$$

That is, the null hypothesis is that the full set of moment conditions are valid, while the alternative hypothesis is that the instrument subset  $\mathbf{z}_{bi}$  is correlated with  $e_i$  and thus an invalid instrument. Rejection of  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  is then interpreted as evidence that  $\mathbf{z}_{bi}$  is misspecified as an instrument.

Based on the same reasoning as described in the previous section, to test  $\mathbb{H}_0$  against  $\mathbb{H}_1$  we consider a partitioned version of the regression (11.69)

$$e_i = \mathbf{z}_{ai}' \boldsymbol{\alpha}_a + \mathbf{z}_{bi}' \boldsymbol{\alpha}_b + \varepsilon_i$$

but now focus on the coefficient  $\alpha_b$ . Given  $\mathbb{E}(z_{ai}e_i) = \mathbf{0}$ ,  $\mathbb{H}_0$  is equivalent to  $\alpha_b = \mathbf{0}$ . The equation is estimated by least-squares, replacing the unobserved  $e_i$  with the 2SLS residual  $\hat{e}_i$ . The estimate of  $\alpha_b$  is

$$\hat{\alpha}_b = (\mathbf{Z}'_b \mathbf{M}_a \mathbf{Z}_b)^{-1} \mathbf{Z}'_b \mathbf{M}_a \hat{\mathbf{e}}$$

where  $\mathbf{M}_a = \mathbf{I}_n - \mathbf{Z}_a (\mathbf{Z}'_a \mathbf{Z}_a)^{-1} \mathbf{Z}'_a$ . Newey (1985) showed that an optimal (asymptotically most powerful) test of  $\mathbb{H}_0$  against  $\mathbb{H}_1$  is to reject for large values of the score statistic

$$\begin{aligned} N &= \hat{\alpha}'_b \left( \widehat{\text{var}}(\hat{\alpha}_b) \right)^{-1} \hat{\alpha}_b \\ &= \frac{\tilde{\mathbf{e}}' \mathbf{R} \left( \mathbf{R}' \mathbf{R} - \mathbf{R}' \widehat{\mathbf{X}} \left( \widehat{\mathbf{X}}' \widehat{\mathbf{X}} \right)^{-1} \widehat{\mathbf{X}}' \mathbf{R} \right)^{-1} \mathbf{R}' \tilde{\mathbf{e}}}{\hat{\sigma}^2} \end{aligned}$$

where  $\widehat{\mathbf{X}} = \mathbf{P} \mathbf{X}$ ,  $\mathbf{P} = \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}'$ ,  $\mathbf{R} = \mathbf{M}_a \mathbf{Z}_b$ , and  $\hat{\sigma}^2 = \frac{1}{n} \tilde{\mathbf{e}}' \tilde{\mathbf{e}}$ .

Independently from Newey (1985), Eichenbaum, Hansen, and Singleton (1988) proposed a test based on the difference of Sargan statistics. Letting  $S$  be the Sargan test statistic (11.70) based on the full instrument set and  $S_a$  be the Sargan test based on the instrument set  $z_{ai}$ , the Sargan difference statistic is

$$C = S - S_a.$$

Specifically, let  $\tilde{\beta}_{2\text{sls}}$  be the 2SLS estimator using the instruments  $z_{ai}$  only, set  $\tilde{e}_i = y_i - \mathbf{x}'_i \tilde{\beta}_{2\text{sls}}$ , and set  $\tilde{\sigma}^2 = \frac{1}{n} \tilde{\mathbf{e}}' \tilde{\mathbf{e}}$ . Then

$$S_a = \frac{\tilde{\mathbf{e}}' \mathbf{Z}_a (\mathbf{Z}'_a \mathbf{Z}_a)^{-1} \mathbf{Z}'_a \tilde{\mathbf{e}}}{\tilde{\sigma}^2}.$$

An advantage of the  $C$  statistic is that it is quite simple to calculate from the standard regression output.

At this point it is useful to reflect on our stated requirement that  $\ell_a > k$ . Indeed, if  $\ell_a < k$  then  $z_{ai}$  fails the order condition for identification and  $\tilde{\beta}_{2\text{sls}}$  cannot be calculated. Thus  $\ell_a \geq k$  is necessary to compute  $S_a$  and hence  $S$ . Furthermore, if  $\ell_a = k$  then  $z_{ai}$  is just identified so while  $\tilde{\beta}_{2\text{sls}}$  can be calculated, the statistic  $S_a = 0$  so  $C = S$ . Thus when  $\ell_a = k$  the subset test equals the full overidentification test so there is no gain from considering subset tests.

The  $C$  statistic  $S_a$  is asymptotically equivalent to replacing  $\tilde{\sigma}^2$  in  $S_a$  with  $\hat{\sigma}^2$ , yielding the statistic

$$C^* = \frac{\tilde{\mathbf{e}}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \tilde{\mathbf{e}}}{\hat{\sigma}^2} - \frac{\tilde{\mathbf{e}}' \mathbf{Z}_a (\mathbf{Z}'_a \mathbf{Z}_a)^{-1} \mathbf{Z}'_a \tilde{\mathbf{e}}}{\hat{\sigma}^2}.$$

It turns out that this is Newey's statistic  $N$ . These tests have chi-square asymptotic distributions.

Let  $c$  satisfy  $\alpha = 1 - G_{\ell_b}(c)$ .

**Theorem 11.28.1** *Algebraically,  $N = C^*$ . Under Assumption 11.14.1 and  $\mathbb{E}(e_i^2 | \mathbf{z}_i) = \sigma^2$ , as  $n \rightarrow \infty$ ,  $N \xrightarrow{d} \chi_{\ell_b}^2$  and  $C \xrightarrow{d} \chi_{\ell_b}^2$ . Thus the tests “Reject  $\mathbb{H}_0$  if  $N > c$ ” and “Reject  $\mathbb{H}_0$  if  $C > c$ ” are asymptotically equivalent and have asymptotic size  $\alpha$ .*

Theorem 11.28.1 shows that  $N$  and  $C^*$  are identical, and are near equivalents to the convenient statistic  $C^*$ , and the appropriate asymptotic distribution is  $\chi_{\ell_b}^2$ . Computationally, the easiest method to implement a subset overidentification test is to estimate the model twice by 2SLS, first using the full instrument set  $\mathbf{z}_i$  and the second using the partial instrument set  $\mathbf{z}_{ai}$ . Compute the Sargan statistics for both 2SLS regressions, and compute  $C$  as the difference in the Sargan statistics. In Stata, for example, this is simple to implement with a few lines of code.

We illustrate using the Card college proximity example. Our reported 2SLS estimates have  $\ell - k = 1$  so there is no role for a subset overidentification test. (Recall, the number of overidentifying restrictions must exceed one.) To illustrate we consider adding extra instruments to the estimates in column 5 of Table 1.1 (the 2SLS estimates using *public*, *private*, *age*, and *age*<sup>2</sup> as instruments for *education*, *experience*, and *experience*<sup>2</sup>/100). We add two instruments: the years of education of the *father* and the *mother* of the worker. These variables had been used in the earlier labor economics literature as instruments, but Card did not. (He used them as regression controls in some specifications.) The motivation for using parent's education as instruments is the hypothesis that parental education influences children's educational attainment, but does not directly influence their ability. The more modern labor economics literature has disputed this idea, arguing that children are educated in part at home, and thus parent's education has a direct impact on the skill attainment of children (and not just an indirect impact via educational attainment). The older view was that parent's education is a valid instrument, the modern view is that it is not valid. We can test this dispute using a overidentification subset test.

We do this by estimating the wage equation by 2SLS using *public*, *private*, *age*, *age*<sup>2</sup>, *father*, and *mother*, as instruments for *education*, *experience*, and *experience*<sup>2</sup>/100). We do not report the parameter estimates here, but observe that this model is overidentified with 3 overidentifying restrictions. We calculate the Sargan overidentification statistic. It is 7.9 with an asymptotic p-value (calculated using  $\chi_3^2$ ) of 0.048. This is a mild rejection of the null hypothesis of correct specification. As we argued in the previous section, this by itself is not reason to reject the model. Now we consider a subset overidentification test. We are interested in testing the validity of the two instruments *father* and *mother*, not the instruments *public*, *private*, *age*, *age*<sup>2</sup>. To test the hypothesis that these two instruments are uncorrelated with the structural error, we compute the difference in Sargan statistic,  $C = 7.9 - 0.5 = 7.4$ , which has a p-value (calculated using  $\chi_2^2$ ) of 0.025. This is marginally statistically significant, meaning that there is evidence that *father* and *mother* are not valid instruments for the wage equation. Since the p-value is not smaller than 1%, it is not overwhelming evidence, but it still supports Card's decision to not use parental education as instruments for the wage equation.

We now prove the results in Theorem 11.28.1.

We first show that  $N = C^*$ . Define  $\mathbf{P}_a = \mathbf{Z}_a (\mathbf{Z}_a' \mathbf{Z}_a)^{-1} \mathbf{Z}_a'$  and  $\mathbf{P}_R = \mathbf{R} (\mathbf{R}' \mathbf{R})^{-1} \mathbf{R}'$ . Since  $[\mathbf{Z}_a, \mathbf{R}]$  span  $\mathbf{Z}$  we find  $\mathbf{P} = \mathbf{P}_R + \mathbf{P}_a$  and  $\mathbf{P}_R \mathbf{P}_a = \mathbf{0}$ . It will be useful to note that

$$\begin{aligned} \mathbf{P}_R \widehat{\mathbf{X}} &= \mathbf{P}_R \mathbf{P} \mathbf{X} = \mathbf{P}_R \mathbf{X} \\ \widehat{\mathbf{X}}' \widehat{\mathbf{X}} - \widehat{\mathbf{X}}' \mathbf{P}_R \widehat{\mathbf{X}} &= \mathbf{X}' (\mathbf{P} - \mathbf{P}_R) \mathbf{X} = \mathbf{X}' \mathbf{P}_a \mathbf{X}. \end{aligned}$$

The fact that  $\mathbf{X}' \mathbf{P} \widehat{\mathbf{e}} = \widehat{\mathbf{X}}' \widehat{\mathbf{e}} = \mathbf{0}$  implies  $\mathbf{X}' \mathbf{P}_R \widehat{\mathbf{e}} = -\mathbf{X}' \mathbf{P}_a \widehat{\mathbf{e}}$ . Finally, since  $\mathbf{y} = \mathbf{X} \widehat{\boldsymbol{\beta}} + \widehat{\mathbf{e}}$ ,

$$\widetilde{\mathbf{e}} = \left( \mathbf{I}_n - \mathbf{X} (\mathbf{X}' \mathbf{P}_a \mathbf{X})^{-1} \mathbf{X}' \mathbf{P}_a \right) \widehat{\mathbf{e}}$$

so

$$\widetilde{\mathbf{e}}' \mathbf{P}_a \widetilde{\mathbf{e}} = \widetilde{\mathbf{e}}' \left( \mathbf{P}_a - \mathbf{P}_a \mathbf{X} (\mathbf{X}' \mathbf{P}_a \mathbf{X})^{-1} \mathbf{X}' \mathbf{P}_a \right) \widehat{\mathbf{e}}.$$

Applying the Woodbury matrix equality to the definition of  $N$ , and the above algebraic relationships,

$$\begin{aligned} N &= \frac{\widetilde{\mathbf{e}}' \mathbf{P}_R \widehat{\mathbf{e}} + \widetilde{\mathbf{e}}' \mathbf{P}_R \widehat{\mathbf{X}} \left( \widehat{\mathbf{X}}' \widehat{\mathbf{X}} - \widehat{\mathbf{X}}' \mathbf{P}_R \widehat{\mathbf{X}} \right)^{-1} \widehat{\mathbf{X}}' \mathbf{P}_R \widehat{\mathbf{e}}}{\widehat{\sigma}^2} \\ &= \frac{\widetilde{\mathbf{e}}' \mathbf{P} \widehat{\mathbf{e}} - \widetilde{\mathbf{e}}' \mathbf{P}_a \widehat{\mathbf{e}} + \widetilde{\mathbf{e}}' \mathbf{P}_a \mathbf{X} (\mathbf{X}' \mathbf{P}_a \mathbf{X})^{-1} \mathbf{X}' \mathbf{P}_a \widehat{\mathbf{e}}}{\widehat{\sigma}^2} \\ &= \frac{\widetilde{\mathbf{e}}' \mathbf{P} \widehat{\mathbf{e}} - \widetilde{\mathbf{e}}' \mathbf{P}_a \widetilde{\mathbf{e}}}{\widehat{\sigma}^2} \\ &= C^* \end{aligned}$$

as claimed.

We next establish the asymptotic distribution. Since  $\mathbf{Z}_a$  is a subset of  $\mathbf{Z}$ ,  $\mathbf{P}\mathbf{M}_a = \mathbf{M}_a\mathbf{P}$ , thus  $\mathbf{P}\mathbf{R} = \mathbf{R}$  and  $\mathbf{R}'\mathbf{X} = \mathbf{R}'\widehat{\mathbf{X}}$ . Consequently

$$\begin{aligned} \frac{1}{\sqrt{n}}\mathbf{R}'\widehat{\mathbf{e}} &= \frac{1}{\sqrt{n}}\mathbf{R}'(\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}) \\ &= \frac{1}{\sqrt{n}}\mathbf{R}'\left(\mathbf{I}_n - \mathbf{X}\left(\widehat{\mathbf{X}}'\widehat{\mathbf{X}}\right)^{-1}\widehat{\mathbf{X}}'\right)\mathbf{e} \\ &= \frac{1}{\sqrt{n}}\mathbf{R}'\left(\mathbf{I}_n - \widehat{\mathbf{X}}\left(\widehat{\mathbf{X}}'\widehat{\mathbf{X}}\right)^{-1}\widehat{\mathbf{X}}'\right)\mathbf{e}. \\ &\xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{V}_2) \end{aligned}$$

where

$$\mathbf{V}_2 = \text{plim}_{n \rightarrow \infty} \left( \frac{1}{n}\mathbf{R}'\mathbf{R} - \frac{1}{n}\mathbf{R}'\widehat{\mathbf{X}}\left(\frac{1}{n}\widehat{\mathbf{X}}'\widehat{\mathbf{X}}\right)^{-1}\frac{1}{n}\widehat{\mathbf{X}}'\mathbf{R} \right).$$

It follows that  $N = C^* \xrightarrow{d} \chi_{\ell_b}^2$  as claimed. Since  $C = C^* + o_p(1)$  it has the same limiting distribution.

## 11.29 Local Average Treatment Effects

In a pair of influential papers, Imbens and Angrist (1994) and Angrist, Imbens and Rubin (1996) proposed an new interpretation of the instrumental variables estimator using the potential outcomes model introduced in Section 2.29.

We will restrict attention to the case that the endogenous regressor  $x$  and excluded instrument  $z$  are binary variables. We write the model as a pair of potential outcome functions. The dependent variable  $y$  is a function of the regressor and an unobservable vector  $\mathbf{u}$

$$y = h(x, \mathbf{u})$$

and the endogenous regressor  $x$  is a function of the instrument  $z$  and  $\mathbf{u}$

$$x = g(z, \mathbf{u}).$$

By specifying  $\mathbf{u}$  as a vector there is no loss of generality in letting both equations depend on  $\mathbf{u}$ .

In this framework, the outcomes are determined by the random vector  $\mathbf{u}$  and the exogenous instrument  $z$ . This determines  $x$ , which determines  $y$ . To put this in the context of the college proximity example, the variable  $\mathbf{u}$  is everything specific about an individual. Given college proximity  $z$ , the person decides to attend college or not. The person's wage is determined by the individual attributes  $\mathbf{u}$  as well as college attendance  $x$ , but is not directly affected by college proximity  $z$ .

We can omit the random variable  $\mathbf{u}$  from the notation as follows. An individual  $i$  has a realization  $\mathbf{u}_i$ . We then set  $y_i(x) = h(x, \mathbf{u}_i)$  and  $x_i(z) = g(z, \mathbf{u}_i)$ . Also, given a realization  $z_i$  the observables are  $x_i = x_i(z_i)$  and  $y_i = y_i(x_i)$ .

In this model the causal effect of college is for individual  $i$  is

$$C_i = y_i(1) - y_i(0).$$

As discussed in Section 2.29, in general this is individual-specific.

We would like to learn about the distribution of the causal effects, or at least features of the distribution. A common feature of interest is the average treatment effect (ATE)

$$ATE = \mathbb{E}(C_i) = \mathbb{E}(y_i(1) - y_i(0)).$$

This, however, is typically not feasible to estimate allowing for endogenous  $x$  without strong assumptions (such as that the causal effect  $C_i$  is constant across individuals). The treatment effect literature has explored what features of the distribution of  $C_i$  can be estimated.

One particular feature of interest, and emphasized by Imbens and Angrist (1994), is known as the local average treatment effect (LATE), and is roughly the average effect upon those effected by the instrumental variable. To understand LATE, it is helpful to consider the college proximity example using the potential outcomes framework. In this framework, each person is fully characterized by their individual unobservable  $\mathbf{u}_i$ . Given  $\mathbf{u}_i$ , their decision to attend college is a function of the proximity indicator  $z_i$ . For some students, proximity has no effect on their decision. For other students, it has an effect in the specific sense that given  $z_i = 1$  they choose to attend college while if  $z_i = 0$  they choose to not attend. We can summarize the possibilities with the following chart, which is based on labels developed by Angrist, Imbens and Rubin (1996).

	$x(0) = 0$	$x(0) = 1$
$x(1) = 0$	Never Takers	Deniers
$x(1) = 1$	Compliers	Always Takers

The columns indicate the college attendance decision given  $z = 0$ . The rows indicate the college attendance decision given  $z = 1$ . The four entries are labels given four types of individuals based on these decisions. The upper-left entry are the individuals who do not attend college regardless of  $z$ . They are called “Never Takers”. The lower-right entry are the individuals who conversely attend college regardless of  $z$ . They are called “Always Takers”. The bottom left are the individuals who only attend college if they live close to one. They are called “Compliers”. The upper right entry is a bit of a challenge. These are individuals who attend college only if they do not live close to one. They are called “Deniers”. Imbens and Angrist discovered that to identify the parameters of interest we need to assume that there are no Deniers, or equivalently that  $x(1) \geq x(0)$ , which they label as a “monotonicity” condition – that increasing the instrument cannot decrease  $x$  for any individual.

We can distinguish the types in the table by the relative values of  $x(1) - x(0)$ . For Never-Takers and Always-Takers,  $x(1) - x(0) = 0$ , while for Deniers,  $x(1) - x(0) = 1$ .

We are interested in the causal effect  $C_i = h(1, \mathbf{u}) - h(0, \mathbf{u})$  of college attendance on wages. Consider the average causal effect among the different types. Among Never-Takers and Always-Takers,  $x(1) = x(0)$  so

$$\mathbb{E}(y_i(1) - y_i(0) | x(1) = x(0))$$

Suppose we try and estimate its average value, conditional for each the three types of individuals: Never-Takers, Always-Takers, and Compliers. It would be impossible for the Never-Takers and Always-Takers. For the former, none attend college so it would be impossible to ascertain the effect of college attendance, and similarly for the latter since they all attend college. Thus the only group for which we can estimate the average causal effect are the Compliers. This is

$$\text{LATE} = \mathbb{E}(y_i(1) - y_i(0) | x_i(1) > x_i(0)).$$

Imbens and Angrist called this the **local average treatment effect (LATE)** as it is the average treatment effect for the sub-population whose endogenous regressor is affected by changes in the instrumental variable.

Interestingly, we show below that

$$\text{LATE} = \frac{\mathbb{E}(y_i | z_i = 1) - \mathbb{E}(y_i | z_i = 0)}{\mathbb{E}(x_i | z_i = 1) - \mathbb{E}(x_i | z_i = 0)}. \quad (11.71)$$

That is, LATE equals the Wald expression (11.32) for the slope coefficient in the IV regression model. This means that the standard IV estimator is an estimator of LATE. Thus when treatment effects are potentially heterogeneous, we can interpret IV as an estimator of LATE. The equality (11.71) occurs under the following conditions.



**Assumption 11.29.1**  $u_i$  and  $z_i$  are independent; and  $\Pr(x_i(1) - x_i(0) < 0) = 0$ .

One interesting feature about LATE is that its value can depend on the instrument  $z_i$  and the distribution of causal effects  $C_i$  in the population. To make this concrete, suppose that instead of the Card proximity instrument, we consider an instrument based on the financial cost of local college attendance. It is reasonable to expect that while the set of students affected by these two instruments are similar, the two sets of students will not be the same. That is, some students may be responsive to proximity but not finances, and conversely. If the causal effect  $C_i$  has a different average in these two groups of students, then LATE will be different when calculated with these two instruments. Thus LATE can vary by the choice of instrument.

How can that be? How can a well-defined parameter depend on the choice of instrument? Doesn't this contradict the basic IV regression model? The answer is that the basic IV regression model is more restrictive – it specifies that the causal effect  $\beta$  is common across all individuals. Thus its value is the same regardless of the choice of specific instrument (so long as it satisfies the instrumental variables assumptions). In contrast, the potential outcomes framework is more general, allowing for the causal effect to vary across individuals. What this analysis shows us is that in this context is quite possible for the LATE coefficient to vary by instrument. This occurs when causal effects are heterogeneous.

One implication of the LATE framework is that IV estimates should be interpreted as causal effects only for the population of compliers. Interpretation should focus on the population of potential compliers and extension to other populations should be done with caution. For example, in the Card proximity model, the IV estimates of the causal return to schooling presented in Table 11.1 should be interpreted as applying to the population of students who are incentivized to attend college by the presence of a college within their home county. The estimates should not be applied to other students.

Formally, the analysis of this section examined the case of a binary instrument and endogenous regressor. How does this generalize? Suppose that the regressor  $x$  is discrete, taking  $J + 1$  discrete values. We can then rewrite the model as one with  $J$  binary endogenous regressors. If we then have  $J$  binary instruments, we are back in the Imbens-Angrist framework (assuming the instruments have a monotonic impact on the endogenous regressors). A benefit is that with a larger set of instruments it is plausible that the set of compliers in the population is expanded.

We close this section by showing (11.71) under Assumption 11.29.1. The realized value of  $x_i$  can be written as

$$x_i = (1 - z_i)x_i(0) + z_i x_i(1) = x_i(0) + z_i(x_i(1) - x_i(0)).$$

Similarly

$$y_i = y_i(0) + x_i(y_i(1) - y_i(0)) = y_i(0) + x_i C_i.$$

Combining,

$$y_i = y_i(0) + x_i(0)C_i + z_i(x_i(1) - x_i(0))C_i.$$

The independence of  $u_i$  and  $z_i$  implies independence of  $(y_i(0), y_i(1), x_i(0), x_i(1), C_i)$  and  $z_i$ . Thus

$$\mathbb{E}(y_i | z_i = 1) = \mathbb{E}(y_i(0)) + \mathbb{E}(x_i(0)C_i) + \mathbb{E}((x_i(1) - x_i(0))C_i)$$

and

$$\mathbb{E}(y_i | z_i = 0) = \mathbb{E}(y_i(0)) + \mathbb{E}(x_i(0)C_i).$$

Subtracting we obtain

$$\begin{aligned} \mathbb{E}(y_i | z_i = 1) - \mathbb{E}(y_i | z_i = 0) &= \mathbb{E}((x_i(1) - x_i(0))C_i) \\ &= 1 \cdot \mathbb{E}(C_i | x_i(1) - x_i(0) = 1) \Pr(x_i(1) - x_i(0) = 1) \\ &\quad + 0 \cdot \mathbb{E}(C_i | x_i(1) - x_i(0) = 0) \Pr(x_i(1) - x_i(0) = 0) \\ &\quad + (-1) \cdot \mathbb{E}(C_i | x_i(1) - x_i(0) = -1) \Pr(x_i(1) - x_i(0) = -1) \\ &= \mathbb{E}(C_i | x_i(1) - x_i(0) = 1) (\mathbb{E}(x_i | z_i = 1) - \mathbb{E}(x_i | z_i = 0)) \end{aligned}$$

where the final equality uses  $\Pr(x_i(1) - x_i(0) < 0) = 0$  and

$$\Pr(x_i(1) - x_i(0) = 1) = \mathbb{E}(x_i(1) - x_i(0)) = \mathbb{E}(x_i | z_i = 1) - \mathbb{E}(x_i | z_i = 0).$$

Rearranging

$$\text{LATE} = \mathbb{E}(C_i | x_i(1) - x_i(0) = 1) = \frac{\mathbb{E}(y_i | z_i = 1) - \mathbb{E}(y_i | z_i = 0)}{\mathbb{E}(x_i | z_i = 1) - \mathbb{E}(x_i | z_i = 0)}$$

as claimed.

### 11.30 Identification Failure

Recall the reduced form equation

$$\mathbf{x}_{2i} = \mathbf{\Gamma}'_{12} \mathbf{z}_{1i} + \mathbf{\Gamma}'_{22} \mathbf{z}_{2i} + \mathbf{u}_{2i}.$$

The parameter  $\beta$  fails to be identified if  $\mathbf{\Gamma}_{22}$  has deficient rank. The consequences of identification failure for inference are quite severe.

Take the simplest case where  $k_1 = 0$  and  $k_2 = \ell_2 = 1$ . Then the model may be written as

$$\begin{aligned} y_i &= x_i \beta + e_i \\ x_i &= z_i \gamma + u_i \end{aligned} \tag{11.72}$$

and  $\Gamma_{22} = \gamma = \mathbb{E}(z_i x_i) / \mathbb{E}(z_i^2)$ . We see that  $\beta$  is identified if and only if  $\gamma \neq 0$ , which occurs when  $\mathbb{E}(x_i z_i) \neq 0$ . Thus identification hinges on the existence of correlation between the excluded exogenous variable and the included endogenous variable.

Suppose this condition fails. In this case  $\gamma = 0$  and  $\mathbb{E}(x_i z_i) = 0$ . We now analyze the distribution of the least-squares and IV estimators of  $\beta$ . For simplicity we assume conditional homoskedasticity and normalize the variances to unity. Thus

$$\begin{aligned} \text{var} \left( \begin{pmatrix} e_i \\ u_i \end{pmatrix} \mid z_i \right) &= \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \\ \mathbb{E}(z_i^2) &= 1. \end{aligned} \tag{11.73}$$

The errors have non-zero correlation  $\rho \neq 0$  which occurs when the variables are endogenous.

By the CLT we have the joint convergence

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \begin{pmatrix} z_i e_i \\ z_i u_i \end{pmatrix} \xrightarrow{d} \begin{pmatrix} \xi_1 \\ \xi_2 \end{pmatrix} \sim N \left( 0, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right). \tag{11.74}$$

It is convenient to define  $\xi_0 = \xi_1 - \rho \xi_2$  which is normal and independent of  $\xi_2$ .

As a benchmark, it is useful to observe that the least-squares estimator of  $\beta$  satisfies

$$\hat{\beta}_{\text{ols}} - \beta = \frac{n^{-1} \sum_{i=1}^n u_i e_i}{n^{-1} \sum_{i=1}^n u_i^2} \xrightarrow{p} \rho \neq 0 \tag{11.75}$$

so endogeneity causes  $\hat{\beta}_{\text{ols}}$  to be inconsistent for  $\beta$ .

Under identification failure  $\gamma = 0$  the asymptotic distribution of the IV estimator is

$$\hat{\beta}_{\text{iv}} - \beta = \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n z_i e_i}{\frac{1}{\sqrt{n}} \sum_{i=1}^n z_i x_i} \xrightarrow{d} \frac{\xi_1}{\xi_2} = \rho + \frac{\xi_0}{\xi_2}.$$

This asymptotic convergence result uses the continuous mapping theorem, which applies since the function  $\xi_1/\xi_2$  is continuous everywhere except at  $\xi_2 = 0$ , which occurs with probability equal to zero.

This limiting distribution has several notable features.

First,  $\hat{\beta}_{iv}$  does not converge in probability to a limit, rather it converges in distribution to a random variable. Thus the IV estimator is inconsistent. Indeed, it is not possible to consistently estimate an unidentified parameter and  $\beta$  is not identified when  $\gamma = 0$ .

Second, the ratio  $\xi_0/\xi_2$  is symmetrically distributed about zero, so the median of the limiting distribution of  $\hat{\beta}_{iv}$  is  $\beta + \rho$ . This means that the IV estimator is median biased under endogeneity. Thus under identification failure the IV estimator does not correct the centering (median bias) of least-squares.

Third, the ratio  $\xi_0/\xi_2$  of two independent normal random variables is Cauchy distributed. This is particularly nasty, as the Cauchy distribution does not have a finite mean. The distribution has thick tails meaning that extreme values occur with higher frequency than the normal, and inferences based on the normal distribution can be quite incorrect.

Together, these results show that  $\gamma = 0$  renders the IV estimator particularly poorly behaved – it is inconsistent, median biased, and non-normally distributed.

We can also examine the behavior of the t-statistic. For simplicity consider the classical (homoskedastic) t-statistic. The error variance estimate has the asymptotic distribution

$$\begin{aligned}\hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n \left( y_i - x_i \hat{\beta}_{iv} \right)^2 \\ &= \frac{1}{n} \sum_{i=1}^n e_i^2 - \frac{2}{n} \sum_{i=1}^n e_i x_i \left( \hat{\beta}_{iv} - \beta \right) + \frac{1}{n} \sum_{i=1}^n x_i^2 \left( \hat{\beta}_{iv} - \beta \right)^2 \\ &\xrightarrow{d} 1 - 2\rho \frac{\xi_1}{\xi_2} + \left( \frac{\xi_1}{\xi_2} \right)^2.\end{aligned}$$

Thus the t-statistic has the asymptotic distribution

$$T = \frac{\hat{\beta}_{iv} - \beta}{\sqrt{\hat{\sigma}^2 \sum_{i=1}^n z_i^2 / |\sum_{i=1}^n z_i x_i|}} \xrightarrow{d} \frac{\xi_1/\xi_2}{\sqrt{1 - 2\rho \frac{\xi_1}{\xi_2} + \left( \frac{\xi_1}{\xi_2} \right)^2}}.$$

The limiting distribution is non-normal, meaning that inference using the normal distribution will be (considerably) incorrect. This distribution depends on the correlation  $\rho$ . The distortion from the normal is increasing in  $\rho$ . Indeed as  $\rho \rightarrow 1$  we have  $\xi_1/\xi_2 \rightarrow_p 1$  and the unexpected finding  $\hat{\sigma}^2 \rightarrow_p 0$ . The latter means that the conventional standard error  $s(\hat{\beta}_{iv})$  for  $\hat{\beta}_{iv}$  also converges in probability to zero. This implies that the t-statistic diverges in the sense  $|T| \rightarrow_p \infty$ . In this situations users may incorrectly interpret estimates as precise, despite the fact that they are useless.

## 11.31 Weak Instruments

In the previous section we examined the extreme consequences of full identification failure. Unfortunately many of the same problems extend to the context where identification is weak in the sense that the reduced form coefficient matrix  $\mathbf{\Gamma}_{22}$  is full rank but small.

A rich asymptotic distribution theory has been developed to understand this setting by modeling  $\mathbf{\Gamma}_{22}$  as “local-to-zero”. The seminal contributions are Staiger and Stock (1997) and Stock and Yogo (2005). The theory was extended to nonlinear GMM estimation by Stock and Wright (2000).

In this section we focus exclusively on the case of one right-hand-side endogenous variable ( $k_2 = 1$ ). We consider the case of multiple endogenous variables in the next section. Our general theory will allow for any arbitrary number of instruments and regressors, but for the sake of clear

exposition we will focus on the very simple case of no included exogenous variables ( $k_1 = 0$ ) and just one exogenous instrument ( $\ell_2 = 1$ ), which is model (11.72) from the previous section

$$\begin{aligned} y_i &= x_i\beta + e_i \\ x_i &= z_i\gamma + u_i. \end{aligned}$$

Furthermore, as in Section 11.30 we assume conditional homoskedasticity and normalize the variances as in (11.73).

The question of primary interest is to determine conditions on the reduced form under which the IV estimator of the structural equation is well behaved, and secondly, what statistical tests can be used to learn if these conditions are satisfied.

In Section 11.30 we assumed complete identification failure in the sense that  $\gamma = 0$ . We now want to assume that identification does not completely fail, but is weak in the sense that  $\gamma$  is small. The technical device which yields a useful distributional theory is to assume that the reduced form parameter is **local-to-zero**, specifically

$$\gamma = n^{-1/2}\mu \quad (11.76)$$

where  $\mu$  is a free parameter. The  $n^{-1/2}$  scaling is picked because it provides just the right balance to allow a useful distribution theory. The local-to-zero assumption (11.76) is not meant to be taken literally but rather is meant to be a useful distributional approximation. The parameter  $\mu$  indexes the degree of identification. Larger  $|\mu|$  implies stronger identification; smaller  $|\mu|$  implies weaker identification.

We now derive the asymptotic distribution of the least-squares and IV estimators under the local-to-unity assumption (11.76).

First, the least-squares estimator satisfies

$$\hat{\beta}_{\text{ols}} - \beta = \frac{n^{-1} \sum_{i=1}^n x_i e_i}{n^{-1} \sum_{i=1}^n x_i^2} = \frac{n^{-1} \sum_{i=1}^n u_i e_i}{n^{-1} \sum_{i=1}^n u_i^2} + o_p(1) \xrightarrow{p} \rho \neq 0$$

which is the same as in (11.75). Thus the least-squares estimator is inconsistent for  $\beta$  under endogeneity.

Second, we derive the distribution of the IV estimator. The joint convergence (11.74) holds, and the local-to-zero assumption implies

$$\begin{aligned} \frac{1}{\sqrt{n}} \sum_{i=1}^n z_i x_i &= \frac{1}{\sqrt{n}} \sum_{i=1}^n z_i^2 \gamma + \frac{1}{\sqrt{n}} \sum_{i=1}^n z_i u_i \\ &= \frac{1}{n} \sum_{i=1}^n z_i^2 \mu + \frac{1}{\sqrt{n}} \sum_{i=1}^n z_i u_i \\ &\xrightarrow{d} \mu + \xi_2. \end{aligned}$$

This allows us to calculate the asymptotic distribution of the IV estimator.

$$\hat{\beta}_{\text{ols}} - \beta = \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n z_i e_i}{\frac{1}{\sqrt{n}} \sum_{i=1}^n z_i x_i} \xrightarrow{d} \frac{\xi_1}{\mu + \xi_2}.$$

This asymptotic convergence result uses the continuous mapping theorem, which applies since the function  $\xi_1/(\mu + \xi_2)$  is a continuous function everywhere except at  $\xi_2 = -\mu$ , which occurs with probability equal to zero.

As in the case of complete identification failure, we find that  $\hat{\beta}_{\text{iv}}$  is inconsistent for  $\beta$  and its asymptotic distribution is non-normal. The distortion is affected by the coefficient  $\mu$ . As  $\mu \rightarrow \infty$

the distribution converges in probability to zero, meaning that  $\hat{\beta}_{iv}$  is consistent for  $\beta$ . This is the classic “strong identification” context.

We also examine the behavior of the classical (homoskedastic) t-statistic for the IV estimator. Note

$$\begin{aligned}\hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n \left( y_i - x_i \hat{\beta}_{iv} \right)^2 \\ &= \frac{1}{n} \sum_{i=1}^n e_i^2 - \frac{2}{n} \sum_{i=1}^n e_i x_i \left( \hat{\beta}_{iv} - \beta \right) + \frac{1}{n} \sum_{i=1}^n x_i^2 \left( \hat{\beta}_{iv} - \beta \right)^2 \\ &\xrightarrow{d} 1 - 2\rho \frac{\xi_1}{\mu + \xi_2} + \left( \frac{\xi_1}{\mu + \xi_2} \right)^2.\end{aligned}$$

Thus

$$T = \frac{\hat{\beta}_{iv} - \beta}{\sqrt{\hat{\sigma}^2 \sum_{i=1}^n z_i^2 / \left| \sum_{i=1}^n z_i x_i \right|}} \xrightarrow{d} \frac{\xi_1}{\sqrt{1 - 2\rho \frac{\xi_1}{\mu + \xi_2} + \left( \frac{\xi_1}{\mu + \xi_2} \right)^2}} \stackrel{def}{=} S. \quad (11.77)$$

In general,  $S$  is non-normal, and its distribution depends on the parameters  $\rho$  and  $\mu$ .

Can we use the distribution  $S$  for inference on  $\beta$ ? The distribution depends on two unknown parameters, and neither is consistently estimable. (Thus we cannot simply use the distribution in (11.77) with  $\rho$  and  $\mu$  replaced with estimates.) To eliminate the dependence on  $\rho$  one possibility is to use the “worst case” value, which turns out to be  $\rho = 1$ . By worst-case we mean that value which causes the greatest distortion away from normal critical values. Setting  $\rho = 1$  we have the considerable simplification

$$S = S_1 = \xi \left| 1 + \frac{\xi}{\mu} \right| \quad (11.78)$$

where  $\xi \sim N(0, 1)$ . When the model is strongly identified (so  $|\mu|$  is very large) then  $S_1 \approx \xi$  is standard normal, consistent with classical theory. However when  $|\mu|$  is very small (but non-zero)  $|S_1| \approx \xi^2/\mu$  (in the sense that this term dominates), which is a scaled  $\chi_1^2$  and quite far from normal. As  $|\mu| \rightarrow 0$  we find the extreme case  $|S_1| \rightarrow_p \infty$ .

While (11.78) is a convenient simplification it does not yield a useful approximation for inference since the distribution in (11.78) is highly dependent on the unknown  $\mu$ . If we try to take the worst-case value of  $\mu$ , which is  $\mu = 0$ , we find that  $|S_1|$  diverges and all distributional approximations fail.

To break this impasse, Stock and Yogo (2005) recommended a constructive alternative. Rather than using the worst-case  $\mu$ , they suggested finding a threshold such that if  $\mu$  exceeds this threshold then the distribution (11.78) is not “too badly” distorted from the normal distribution.

Specifically, the Stock-Yogo recommendation can be summarized by two steps. First, the distribution result (11.78) can be used to find a threshold value  $\tau^2$  such that if  $\mu^2 \geq \tau^2$  then the size of the nominal<sup>1</sup> 5% test “Reject if  $|T| \geq 1.96$ ” has asymptotic size  $\Pr(|S_1| \geq 1.96) \leq 0.15$ . This means that while the goal is to obtain a test with size 5%, we recognize that there may be size distortion due to weak instruments and are willing to tolerate a specific size distortion, for example 10% distortion (allow for actual size up to 15%, or more generally  $r$ ). Second, they use the asymptotic distribution of the reduced-form (first stage)  $F$  statistic to test if the actual unknown value of  $\mu^2$  exceeds the threshold  $\tau^2$ . These two steps together give rise to the rule-of-thumb that the first-stage  $F$  statistic should exceed 10 in order to achieve reliable IV inference. (This is for the case of one instrumental variable. If there is more than one instrument then the rule-of-thumb changes.) We now describe the steps behind this reasoning in more detail.

<sup>1</sup>The term “nominal size” of a test is the official intended size – the size which would obtain under ideal circumstances. In this context the test “Reject if  $|T| \geq 1.96$ ” has nominal size 0.05 as this would be the asymptotic rejection probability in the ideal context of strong instruments.

The first step is to use the distribution (11.77) to determine the threshold  $\tau^2$ . Formally, the goal is to find the value of  $\tau^2 = \mu^2$  at which the asymptotic size of a nominal 5% test is actually  $r$  (e.g.  $r = 0.15$ )

$$\Pr(|S_1| \geq 1.96) \leq r.$$

By some algebra and using the quadratic formula the event  $|\xi(1 + \xi/\mu)| < x$  is the same as

$$\frac{\mu^2}{4} - x\mu < \left(\xi + \frac{\mu}{2}\right)^2 < \frac{\mu^2}{4} + x\mu.$$

The random variable between the inequalities is distributed  $\chi_1^2(\mu^2/4)$ , a noncentral chi-square with one degree of freedom and noncentrality parameter  $\mu^2/4$ . Thus

$$\begin{aligned} \Pr(|S_1| \geq x) &= \Pr\left(\chi_1^2\left(\frac{\mu^2}{4}\right) \geq \frac{\mu^2}{4} + x\mu\right) + \Pr\left(\chi_1^2\left(\frac{\mu^2}{4}\right) \leq \frac{\mu^2}{4} - x\mu\right) \\ &= 1 - G\left(\frac{\mu^2}{4} + x\mu, \frac{\mu^2}{4}\right) + G\left(\frac{\mu^2}{4} - x\mu, \frac{\mu^2}{4}\right) \end{aligned} \quad (11.79)$$

where  $G(u, \lambda)$  is the distribution function of  $\chi_1^2(\lambda)$ . Hence the desired threshold  $\tau^2$  solves

$$1 - G\left(\frac{\tau^2}{4} + 1.96\tau, \frac{\tau^2}{4}\right) + G\left(\frac{\tau^2}{4} - 1.96\tau, \frac{\tau^2}{4}\right) = r$$

or effectively

$$G\left(\frac{\tau^2}{4} + 1.96\tau, \frac{\tau^2}{4}\right) = 1 - r$$

since  $\tau^2/4 - 1.96\tau < 0$  for relevant values of  $\tau$ . The numerical solution (computed with the non-central chi-square distribution function, e.g. `ncx2cdf` in MATLAB) is  $\tau^2 = 1.70$  when  $r = 0.15$ . (That is, the command `ncx2cdf(1.7/4+1.96*sqrt(1.7),1,1.7/4)` yields the answer 0.8500. Stock and Yogo (2005) approximate the same calculation using simulation methods and report  $\tau^2 = 1.82$ .)

This calculation means that if the true reduced form coefficient satisfies  $\mu^2 \geq 1.7$ , or equivalently if  $\gamma^2 \geq 1.7/n$ , then the (asymptotic) size of a nominal 5% test on the structural parameter is no larger than 15%.

To summarize the Stock-Yogo first step, we calculate the minimum value  $\tau^2$  for  $\mu^2$  sufficient to ensure that the asymptotic size of a nominal 5% t-test does not exceed  $r$ , and find that  $\tau^2 = 1.70$  for  $r = 0.15$ .

The Stock-Yogo second step is to find a critical value for the first-stage  $F$  statistic sufficient to reject the hypothesis that  $\mathbb{H}_0 : \mu^2 = \tau^2$  against  $\mathbb{H}_1 : \mu^2 > \tau^2$ . We now describe this procedure.

They suggest testing  $\mathbb{H}_0 : \mu^2 = \tau^2$  at the 5% size using the first stage  $F$  statistic. If the  $F$  statistic is small so that the test does not reject then we should be worried that the true value of  $\mu^2$  is small and there is a weak instrument problem. On the other hand if the  $F$  statistic is large so that the test rejects then we can have some confidence that the true value of  $\mu^2$  is sufficiently large that the weak instrument problem is not too severe.

To implement the test we need to calculate an appropriate critical value. It should be calculated under the null hypothesis  $\mathbb{H}_0 : \mu^2 = \tau^2$ . This is different from a conventional  $F$  test (which has the null hypothesis  $\mathbb{H}_0 : \mu^2 = 0$ ).

We start by calculating the asymptotic distribution of  $F$ . Since there is just one regressor and one instrument in our simplified setting, the first-stage  $F$  statistic is the squared t-statistic from the reduced form, and given our previous calculations has the asymptotic distribution

$$F = \frac{\hat{\gamma}^2}{s(\hat{\gamma})^2} = \frac{(\sum_{i=1}^n z_i x_i)^2}{(\sum_{i=1}^n x_i^2) \hat{\sigma}_u^2} \xrightarrow{d} (\mu + \xi_2)^2 \sim \chi_1^2(\mu^2).$$

This is a non-central chi-square distribution with one degree of freedom and non-centrality parameter  $\mu^2$ . The distribution function of the latter is  $G(u, \mu^2)$ .

To test  $\mathbb{H}_0 : \mu^2 = \tau^2$  against  $\mathbb{H}_1 : \mu^2 > \tau^2$  we reject for  $F \geq c$  where  $c$  is selected so that the asymptotic rejection probability

$$\Pr(F \geq c) \rightarrow \Pr(\chi_1^2(\mu^2) \geq c) = 1 - G(c, \mu^2)$$

equals 0.05 under  $\mathbb{H}_0 : \mu^2 = \tau^2$ , or equivalently

$$G(c, \tau^2) = G(c, 1.7) = 0.95.$$

This can be found using the non-central chi-square quantile function, e.g. the function  $Q(p, d)$  which solves  $G(Q(p, d), d) = p$ . We find that

$$c = Q(0.95, 1.7) = 8.7.$$

In MATLAB, this can be computed by `ncx2inv(.95,1.7)`. (Stock and Yogo (2005) report  $c = 9.0$  since they used  $\tau^2 = 1.82$ .)

This means that if  $F > 8.7$  we can reject  $\mathbb{H}_0 : \mu^2 = 1.7$  against  $\mathbb{H}_1 : \mu^2 > 1.7$  with an asymptotic 5% test. In this context we should expect the IV estimate and tests to be reasonably well behaved. However, if  $F < 8.7$  then we should be cautious about the IV estimator, confidence intervals, and tests. This finding led Staiger and Stock (1997) to propose the informal “rule of thumb” that the first stage  $F$  statistic should exceed 10. Notice that  $F$  exceeding 8.7 (or 10) is equivalent to the reduced form t-statistic exceeding 2.94 (or 3.16), which is considerably larger than a conventional check if the t-statistic is “significant”. Equivalently, the recommended rule-of-thumb for the case of a single instrument is to estimate the reduced form and verify that the t-statistic for exclusion of the instrumental variable exceeds 3 in absolute value.

Does the proposed procedure control the asymptotic size of a 2SLS test? The first step has asymptotic size bounded below  $r$  (e.g. 15%). The second step has asymptotic size 5%. By the Bonferroni bound (see Section 9.20) the two steps together have asymptotic size bounded below  $r + 0.05$  (e.g. 20%). We can thus call the Stock-Yogo procedure a rigorous test with asymptotic size  $r + 0.05$  (or 20%).

Our analysis has been confined to the case  $k_2 = \ell_2 = 1$ . Stock and Yogo (2005) also examine the case of  $\ell_2 > 1$  (which requires numerical simulation to solve), and both the 2SLS and LIML estimators. They show that the  $F$  statistic critical values depend on the number of instruments  $\ell_2$  as well as the estimator. We report their calculations here.

F Statistic 5% Critical Value for Weak Instruments,  $k_2 = 1$

	Maximal Size $r$							
	2SLS				LIML			
$\ell_2$	0.10	0.15	0.20	0.25	0.10	0.15	0.20	0.25
1	16.4	9.0	6.7	5.5	16.4	9.0	6.7	5.5
2	19.9	11.6	8.7	7.2	8.7	5.3	4.4	3.9
3	22.3	12.8	9.5	7.8	6.5	4.4	3.7	3.3
4	24.6	14.0	10.3	8.3	5.4	3.9	3.3	3.0
5	26.9	15.1	11.0	8.8	4.8	3.6	3.0	2.8
6	29.2	16.2	11.7	9.4	4.4	3.3	2.9	2.6
7	31.5	17.4	12.5	9.9	4.2	3.2	2.7	2.5
8	33.8	18.5	13.2	10.5	4.0	3.0	2.6	2.4
9	36.2	19.7	14.0	11.1	3.8	2.9	2.5	2.3
10	38.5	20.9	14.8	11.6	3.7	2.8	2.5	2.2
15	50.4	26.8	18.7	12.2	3.3	2.5	2.2	2.0
20	62.3	32.8	22.7	17.6	3.2	2.3	2.1	1.9
25	74.2	38.8	26.7	20.6	3.8	2.2	2.0	1.8
30	86.2	44.8	30.7	23.6	3.9	2.2	1.9	1.7

One striking feature about these critical values is that those for the 2SLS estimator are strongly increasing in  $\ell_2$  while those for the LIML estimator are decreasing in  $\ell_2$ . This means that when the number of instruments  $\ell_2$  is large, 2SLS requires a much stronger reduced form (larger  $\mu^2$ ) in order for inference to be reliable, but this is not the case for LIML. This is direct evidence that inference is less sensitive to weak instruments when estimation is by LIML rather than 2SLS. This makes a strong case for using LIML rather than 2SLS, especially when  $\ell_2$  is large or the instruments are potentially weak.

We now summarize the recommended Staiger-Stock/Stock-Yogo procedure for  $k_1 \geq 1$ ,  $k_2 = 1$ , and  $\ell_2 \geq 1$ . The structural equation and reduced form equations are

$$\begin{aligned} y_i &= \mathbf{x}'_{1i}\beta_1 + x_{2i}\beta_2 + e_i \\ x_{2i} &= \mathbf{x}'_{1i}\gamma_1 + \mathbf{z}'_{2i}\gamma_2 + u_i. \end{aligned}$$

The reduced form is estimated by least-squares

$$x_{2i} = \mathbf{x}'_{1i}\hat{\gamma}_1 + \mathbf{z}'_{2i}\hat{\gamma}_2 + \hat{u}_i$$

and the structural equation by either 2SLS or LIML:

$$y_i = \mathbf{x}'_{1i}\hat{\beta}_1 + x_{2i}\hat{\beta}_2 + \hat{e}_i.$$

Let  $F$  be the  $F$  statistic for  $\mathbb{H}_0 : \gamma_2 = 0$  in the reduced form equation. Let  $s(\hat{\beta}_2)$  be a standard error for  $\beta_2$  in the structural equation. The procedure is:

1. Compare  $F$  with the critical values  $c$  in the above table, with the row selected to match the number of excluded instruments  $\ell_2$ , and the columns to match the estimation method (2SLS or LIML) and the desired size  $r$ .
2. If  $F > c$  then report the 2SLS or LIML estimates with conventional inference.

The Stock-Yogo test can be implemented in Stata using the command `estat firststage` after `ivregress 2sls` or `ivregress liml` if a standard (non-robust) covariance matrix has been specified (that is, without the `,r` option).

There are possible extensions to the Stock-Yogo procedure.

One modest extension is to use the information to convey the degree of confidence in the accuracy of a confidence interval. Suppose in an application you have  $\ell_2 = 5$  excluded instruments and have estimated your equation by 2SLS. Now suppose that your reduced form  $F$  statistic equals 12. You check the Stock-Yogo table, and find that  $F = 12$  is significant with  $r = 0.20$ . Thus we can interpret the conventional 2SLS confidence interval as having coverage of 80% (or 75% if we make the Bonferroni correction). On the other hand if  $F = 27$  we would conclude that the test for weak instruments is significant with  $r = 0.10$ , meaning that the conventional 2SLS confidence interval can be interpreted as having coverage of 90% (or 85% after Bonferroni correction).

A more substantive extension, which we now discuss, reverses the steps. Unfortunately this discussion will be limited to the case  $\ell_2 = 1$ , where 2SLS and LIML are equivalent. First, use the reduced form  $F$  statistic to find a one-sided confidence interval for  $\mu^2$  of the form  $[\mu_L^2, \infty)$ . Second, use the lower bound  $\mu_L^2$  to calculate a critical value  $C$  for  $S_1$  such that the 2SLS test has asymptotic size bounded below 0.05. This produces better size control than the Stock-Yogo procedure and produces more informative confidence intervals for  $\beta_2$ . We now describe the steps in detail.

The first goal is to find a one-sided confidence interval for  $\mu^2$ . This is found by test inversion. As we described earlier, for any  $\tau^2$  we reject  $\mathbb{H}_0 : \mu^2 = \tau^2$  in favor of  $\mathbb{H}_1 : \mu^2 > \tau^2$  if  $F > c$  where  $G(c, \tau^2) = 0.95$ . Equivalently, we reject if  $G(F, \tau^2) > 0.95$ . By the test inversion principle, an asymptotic 95% confidence interval  $[\mu_L^2, \infty)$  can be formed as the set of all values of  $\tau^2$  which



are not rejected by this test. Since  $G(F, \tau^2) \geq 0.95$  for all  $\tau^2$  in this set, the lower bound  $\mu_L^2$  satisfies  $G(F, \mu_L^2) = 0.95$ . The lower bound is found from this equation. Since this solution is not generally programmed, it needs to be found numerically. In MATLAB, the solution is `mu2` when `ncx2cdf(F,1,mu2)` returns `0.95`.

The second goal is to find the critical value  $C$  such that  $\Pr(|S_1| \geq C) = 0.05$  when  $\mu^2 = \mu_L^2$ . From (11.79), this is achieved when

$$1 - G\left(\frac{\mu_L^2}{4} + C\mu_L, \frac{\mu_L^2}{4}\right) + G\left(\frac{\mu_L^2}{4} - C\mu_L, \frac{\mu_L^2}{4}\right) = 0.05. \quad (11.80)$$

This can be solved as

$$G\left(\frac{\mu_L^2}{4} + C\mu_L, \frac{\mu_L^2}{4}\right) = 0.95.$$

(The third term on the left-hand-side of (11.80) is zero for all solutions so can be ignored.) Using the non-central chi-square quantile function  $Q(p, d)$ , this  $C$  equals

$$C = \frac{Q\left(0.95, \frac{\mu_L^2}{4}\right) - \frac{\mu_L^2}{4}}{\mu_L}.$$

For example, in MATLAB this is found as `C=(ncx2inv(.95,1,mu2/4)-mu2/4)/sqrt(mu2)`. 95% confidence intervals for  $\beta_2$  are then calculated as

$$\hat{\beta}_{IV} \pm Cs(\hat{\beta}_{IV}).$$

We can also calculate a p-value for the t-statistic  $T$  for  $\beta_2$ . These are

$$p = 1 - G\left(\frac{\mu_L^2}{4} + |T|\mu_L, \frac{\mu_L^2}{4}\right) + G\left(\frac{\mu_L^2}{4} - |T|\mu_L, \frac{\mu_L^2}{4}\right)$$

where the third term equals zero if  $|T| \geq \mu_L/4$ . In MATLAB, for example, this can be calculated by the commands

```
T1 = mu2/4 + abs(T) * sqrt(mu2);
T2 = mu2/4 - abs(T) * sqrt(mu2);
p = -ncx2cdf(T1, 1, mu2/4) + ncx2cdf(T2, 1, mu2/4);
```

These confidence intervals and p-values will be larger than the conventional intervals and p-values, reflecting the incorporation of information about the strength of the instruments through the first-stage  $F$  statistic. Also, by the Bonferroni bound these tests have asymptotic size bounded below 10% and the confidence intervals have asymptotic coverage exceeding 90%, unlike the Stock-Yogo method which has size of 20% and coverage of 80%.

The augmented procedure suggested here, only for the  $\ell_2 = 1$  case, is

1. Find  $\mu_L^2$  which solves  $G(F, \mu_L^2) = 0.95$ . In MATLAB, the solution is `mu2` when `ncx2cdf(F,1,mu2)` returns `0.95`.
2. Find  $C$  which solves  $G(\mu_L^2/4 + C\mu_L, \mu_L^2/4) = 0.95$ . In MATLAB, the command is `C=(ncx2inv(.95,1,mu2/4)-mu2/4)/sqrt(mu2)`
3. Report the confidence interval  $\hat{\beta}_2 \pm Cs(\hat{\beta}_2)$  for  $\beta_2$ .
4. For the t statistic  $T = (\hat{\beta}_2 - \beta_2) / s(\hat{\beta}_2)$  the asymptotic p-value is

$$p = 1 - G\left(\frac{\mu_L^2}{4} + |T|\mu_L, \frac{\mu_L^2}{4}\right) + G\left(\frac{\mu_L^2}{4} - |T|\mu_L, \frac{\mu_L^2}{4}\right)$$

which is computed in MATLAB by `T1=mu2/4+abs(T)*sqrt(mu2); T2=mu2/4-abs(T)*sqrt(mu2);` and `p=1-ncx2cdf(T1,1,mu2/4)+ncx2cdf(T2,1,mu2/4)`.

We have described an extension to the Stock-Yogo procedure for the case of one instrumental variable  $\ell_2 = 1$ . This restriction was due to the use of the analytic formula (11.80) for the asymptotic distribution, which is only available when  $\ell_2 > 0$ . In principle the procedure could be extended using simulation or bootstrap methods, but this has not been done to my knowledge.

To illustrate the Stock-Yogo and extended procedures, let us return to the Card proximity example. First, let's take the IV estimates reported in the second column of Table 11.1 which used *college* proximity as a single instrument. The reduced form estimates for the endogenous variable *education* is reported in the second column of Table 11.2. The excluded instrument *college* has a t-ratio of 4.2 which implies an  $F$  statistic of 17.8. The  $F$  statistic exceeds the rule-of thumb of 10, so the structural estimates pass the Stock-Yogo threshold. Based on the Stock-Yogo recommendation, this means that we can interpret the estimates conventionally. However, the conventional confidence interval, e.g. for the returns to education,  $0.132 \pm 0.049 * 1.96 = [0.04, 0.23]$  has an asymptotic coverage of 80%, rather than the nominal 95% rate.

Now consider the extended procedure. Given  $F = 17.8$  we can calculate the lower bound  $\mu_L^2 = 6.6$ . This implies a critical value of  $C = 2.7$ . Hence an improved confidence interval for the returns to education in this equation is  $0.132 \pm 0.049 * 2.7 = [0.01, 0.26]$ . This is a wider confidence interval, but has improved asymptotic coverage of 90%. The p-value for  $\beta_2 = 0$  is  $p = 0.012$ .

Next, let's take the 2SLS estimates reported in the fourth column of Table 11.1 which use the two instruments *public* and *private*. The reduced form equation is reported in column six of Table 11.2. An  $F$  statistic for exclusion of the two instruments is  $F = 13.9$ , which exceeds the 15% size threshold for 2SLS and all thresholds for LIML, indicating that the structural estimates pass the Stock-Yogo threshold test and can be interpreted conventionally.

The weak instrument methods described here are important for applied econometrics as they discipline researchers to assess the quality of their reduced form relationships before reporting structural estimates. The theory, however, has limitations and shortcomings. A major limitation is that the theory requires the strong assumption of conditional homoskedasticity. Despite this theoretical limitation, in practice researchers apply the Stock-Yogo recommendations to estimates computed with heteroskedasticity-robust standard errors as it is the currently the best known approach. This is an active area of research so the recommended methods may change in the years ahead.

### James Stock

James Stock (1955-) is a American econometrician and empirical macro-economist who has made several important contributions, most notably his work on weak instruments, unit root testing, cointegration, and forecasting. He is also well-known for his undergraduate textbook *Introduction to Econometrics* (2014) co-authored with Mark Watson

## 11.32 Weak Instruments with $k_2 > 1$

When there are more than one endogenous regressor ( $k_2 > 1$ ) it is better to examine the reduced form as a system. Staiger and Stock (1997) and Stock and Yogo (2005) provided an analysis of this case and constructed a test for weak instruments. The theory is considerably more involved than the  $k_2 = 1$  case, so we briefly summarize it here excluding many details, emphasizing their suggested methods.

The structural equation and reduced form equations are

$$\begin{aligned} y_i &= \mathbf{x}'_{1i}\boldsymbol{\beta}_1 + \mathbf{x}'_{2i}\boldsymbol{\beta}_2 + e_i \\ \mathbf{z}_{2i} &= \boldsymbol{\Gamma}'_{12}\mathbf{z}_{1i} + \boldsymbol{\Gamma}'_{22}\mathbf{z}_{2i} + \mathbf{u}_{2i}. \end{aligned}$$

As in the previous section we assume that the errors are conditionally homoskedastic.

Identification of  $\boldsymbol{\beta}_2$  requires the matrix  $\boldsymbol{\Gamma}_{22}$  to be full rank. A necessary condition is that each row of  $\boldsymbol{\Gamma}'_{22}$  is non-zero, but this is not sufficient.

We focus on the size performance of the homoskedastic Wald statistic for the 2SLS estimator of  $\boldsymbol{\beta}_2$ . For simplicity assume that the variance of  $e_i$  is known and normalized to one. Using representation (11.37), the Wald statistic can be written as

$$W = e'\tilde{\mathbf{Z}}_2 \left( \tilde{\mathbf{Z}}_2' \tilde{\mathbf{Z}}_2 \right)^{-1} \tilde{\mathbf{Z}}_2' \mathbf{X}_2 \left( \mathbf{X}_2' \tilde{\mathbf{Z}}_2 \left( \tilde{\mathbf{Z}}_2' \tilde{\mathbf{Z}}_2 \right)^{-1} \tilde{\mathbf{Z}}_2' \mathbf{X}_2 \right)^{-1} \left( \mathbf{X}_2' \tilde{\mathbf{Z}}_2 \left( \tilde{\mathbf{Z}}_2' \tilde{\mathbf{Z}}_2 \right)^{-1} \tilde{\mathbf{Z}}_2' e \right)$$

where  $\tilde{\mathbf{Z}}_2 = (\mathbf{I}_n - \mathbf{P}_1) \mathbf{Z}_2$  and  $\mathbf{P}_1 = \mathbf{X}_1 (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1'$ .

Stock and Staiger model the excluded instruments  $\mathbf{z}_{2i}$  as weak by setting  $\boldsymbol{\Gamma}_{22} = n^{-1/2} \mathbf{C}$  for some matrix  $\mathbf{C}$ . This is the multivariate analog of the simple case examined in the previous section. In this framework we have the asymptotic distribution results

$$\begin{aligned} \frac{1}{n} \tilde{\mathbf{Z}}_2' \tilde{\mathbf{Z}}_2 &\xrightarrow{p} \mathbf{Q} = \mathbb{E}(\mathbf{z}_{2i} \mathbf{z}_{2i}') - \mathbb{E}(\mathbf{z}_{2i} \mathbf{z}_{1i}') \left( \mathbb{E}(\mathbf{z}_{1i} \mathbf{z}_{1i}') \right)^{-1} \mathbb{E}(\mathbf{z}_{1i} \mathbf{z}_{2i}') \\ \frac{1}{\sqrt{n}} \tilde{\mathbf{Z}}_2' e &\xrightarrow{d} \mathbf{Q}^{1/2} \boldsymbol{\xi}_0 \end{aligned}$$

where  $\boldsymbol{\xi}_0$  is a matrix normal variate whose columns are independent  $N(\mathbf{0}, \mathbf{I})$ . Furthermore, setting  $\boldsymbol{\Sigma} = \mathbb{E}(\mathbf{u}_{2i} \mathbf{u}_{2i}')$  and  $\overline{\mathbf{C}} = \mathbf{Q}^{1/2} \mathbf{C} \boldsymbol{\Sigma}^{-1/2}$ ,

$$\frac{1}{\sqrt{n}} \tilde{\mathbf{Z}}_2' \mathbf{X}_2 = \frac{1}{n} \tilde{\mathbf{Z}}_2' \tilde{\mathbf{Z}}_2 \mathbf{C} + \frac{1}{\sqrt{n}} \tilde{\mathbf{Z}}_2' \mathbf{U}_2 \xrightarrow{d} \mathbf{Q}^{1/2} \overline{\mathbf{C}} \boldsymbol{\Sigma}^{1/2} + \mathbf{Q}^{1/2} \boldsymbol{\xi}_2 \boldsymbol{\Sigma}^{1/2}$$

where  $\boldsymbol{\xi}_2$  is a matrix normal variates whose columns are independent  $N(\mathbf{0}, \mathbf{I})$ . The variables  $\boldsymbol{\xi}_0$  and  $\boldsymbol{\xi}_2$  are correlated. Together we obtain the asymptotic distribution of the Wald statistic

$$W \xrightarrow{d} S = \boldsymbol{\xi}_0' (\overline{\mathbf{C}} + \boldsymbol{\xi}_2) \left( \overline{\mathbf{C}}' \overline{\mathbf{C}} \right)^{-1} (\overline{\mathbf{C}} + \boldsymbol{\xi}_2)' \boldsymbol{\xi}_0.$$

Using the spectral decomposition,  $\overline{\mathbf{C}}' \overline{\mathbf{C}} = \mathbf{H}' \boldsymbol{\Lambda} \mathbf{H}$  where  $\mathbf{H}' \mathbf{H} = \mathbf{I}$  and  $\boldsymbol{\Lambda}$  is diagonal. Thus we can write

$$S = \boldsymbol{\xi}_0' \bar{\boldsymbol{\xi}}_2 \boldsymbol{\Lambda}^{-1} \bar{\boldsymbol{\xi}}_2' \boldsymbol{\xi}_0$$

where  $\bar{\boldsymbol{\xi}}_2 = \overline{\mathbf{C}} \mathbf{H}' + \boldsymbol{\xi}_2 \mathbf{H}'$ . The matrix  $\boldsymbol{\xi}^* = (\boldsymbol{\xi}_0, \bar{\boldsymbol{\xi}}_2)$  is multivariate normal, so  $\boldsymbol{\xi}^{*'} \boldsymbol{\xi}^*$  has what is called a non-central Wishart distribution. It only depends on the matrix  $\overline{\mathbf{C}}$  through  $\mathbf{H} \overline{\mathbf{C}}' \overline{\mathbf{C}} \mathbf{H}' = \boldsymbol{\Lambda}$ , which are the eigenvalues of  $\overline{\mathbf{C}}' \overline{\mathbf{C}}$ . Since  $S$  is a function of  $\boldsymbol{\xi}^*$  only through  $\bar{\boldsymbol{\xi}}_2' \boldsymbol{\xi}_0$  we conclude that  $S$  is a function of  $\overline{\mathbf{C}}$  only through these eigenvalues.

This is a very quick derivation of a rather involved derivation, but the conclusion drawn by Stock and Yogo is that the asymptotic distribution of the Wald statistic is non-standard, and a function of the model parameters only through the eigenvalues of  $\overline{\mathbf{C}}' \overline{\mathbf{C}}$  and the correlations between the normal variates  $\boldsymbol{\xi}_0$  and  $\bar{\boldsymbol{\xi}}_2$ . The worst-case can be summarized by the maximal correlation between  $\boldsymbol{\xi}_0$  and  $\bar{\boldsymbol{\xi}}_2$  and the smallest eigenvalue of  $\overline{\mathbf{C}}' \overline{\mathbf{C}}$ . For convenience, they rescale the latter by dividing by the number of endogenous variables. Define

$$\mathbf{G} = \overline{\mathbf{C}}' \overline{\mathbf{C}} / k_2 = \boldsymbol{\Sigma}^{-1/2} \mathbf{C}' \mathbf{Q} \mathbf{C} \boldsymbol{\Sigma}^{-1/2} / k_2$$

and

$$g = \lambda_{\min}(\mathbf{G}) = \lambda_{\min} \left( \boldsymbol{\Sigma}^{-1/2} \mathbf{C}' \mathbf{Q} \mathbf{C} \boldsymbol{\Sigma}^{-1/2} \right) / k_2.$$

This can be estimated from the reduced-form regression

$$\mathbf{x}_{2i} = \hat{\mathbf{\Gamma}}'_{12} \mathbf{z}_{1i} + \hat{\mathbf{\Gamma}}'_{22} \mathbf{z}_{2i} + \hat{\mathbf{u}}_{2i}.$$

The estimator is

$$\begin{aligned} \hat{\mathbf{G}} &= \hat{\mathbf{\Sigma}}^{-1/2} \hat{\mathbf{\Gamma}}'_{22} \left( \tilde{\mathbf{Z}}'_2 \tilde{\mathbf{Z}}_2 \right) \hat{\mathbf{\Gamma}}_{22} \hat{\mathbf{\Sigma}}^{-1/2} / k_2 \\ &= \hat{\mathbf{\Sigma}}^{-1/2} \left( \mathbf{X}'_2 \tilde{\mathbf{Z}}_2 \left( \tilde{\mathbf{Z}}'_2 \tilde{\mathbf{Z}}_2 \right)^{-1} \tilde{\mathbf{Z}}'_2 \mathbf{X}_2 \right) \hat{\mathbf{\Sigma}}^{-1/2} / k_2 \\ \hat{\mathbf{\Sigma}} &= \frac{1}{n-k} \sum_{i=1}^n \hat{\mathbf{u}}_{2i} \hat{\mathbf{u}}'_{2i} \\ \hat{g} &= \lambda_{\min} \left( \hat{\mathbf{G}} \right). \end{aligned}$$

$\hat{\mathbf{G}}$  is a matrix  $F$ -type statistic for the coefficient matrix  $\hat{\mathbf{\Gamma}}_{22}$ .

The statistic  $\hat{g}$  was proposed by Craig and Donald (1993) as a test for underidentification. Stock and Yogo (2005) use it as a test for weak instruments. Using simulation methods, they determined critical values for  $\hat{g}$  similar to those for the  $k_2 = 1$  case. For given size  $r > 0.05$ , there is a critical value  $c$  (reported in the table below) such that if  $\hat{g} > c$ , then the 2SLS (or LIML) Wald statistic  $W$  for  $\hat{\beta}_2$  has asymptotic size bounded below  $r$ . On the other hand, if  $\hat{g} \leq c$  then we cannot bound the asymptotic size below  $r$  and we cannot reject the hypothesis of weak instruments.

The Stock-Yogo critical values for  $k_2 = 2$  are presented in the following table. The methods and theory applies to the cases  $k_2 > 2$  as well, but those critical values have not been calculated. As for the  $k_2 = 1$  case, the critical values for 2SLS are dramatically increasing in  $\ell_2$ . Thus when the model is over-identified, we need quite a large value of  $\hat{g}$  to reject the hypothesis of weak instruments. This is a strong cautionary message to check the  $\hat{g}$  statistic in applications. Furthermore, the critical values for LIML are generally decreasing in  $\ell_2$  (except for  $r = 0.10$ , where the critical values are increasing for large  $\ell_2$ ). This means that for over-identified models, LIML inference is much less sensitive to weak instruments than 2SLS, and may be the preferred estimation method.

The Stock-Yogo test can be implemented in Stata for  $k_2 \leq 2$  using the command `estat firststage` after `ivregress 2sls` or `ivregress liml` if a standard (non-robust) covariance matrix has been specified (that is, without the ‘,r’ option).

$\hat{g}$  5% Critical Value for Weak Instruments,  $k_2 = 2$

	Maximal Size $r$							
	2SLS				LIML			
$\ell_2$	0.10	0.15	0.20	0.25	0.10	0.15	0.20	0.25
2	7.0	4.6	3.9	3.6	7.0	4.6	3.9	3.6
3	13.4	8.2	6.4	5.4	5.4	3.8	3.3	3.1
4	16.9	9.9	7.5	6.3	4.7	3.4	3.0	2.8
5	19.4	11.2	8.4	6.9	4.3	3.1	2.8	2.6
6	21.7	12.3	9.1	7.4	4.1	2.9	2.6	2.5
7	23.7	13.3	9.8	7.9	3.9	2.8	2.5	2.4
8	25.6	14.3	10.4	8.4	3.8	2.7	2.4	2.3
9	27.5	15.2	11.0	8.8	3.7	2.7	2.4	2.2
10	29.3	16.2	11.6	9.3	3.6	2.6	2.3	2.1
15	38.0	20.6	14.6	11.6	3.5	2.4	2.1	2.0
20	46.6	25.0	17.6	13.8	3.6	2.4	2.0	1.9
25	55.1	29.3	20.6	16.1	3.6	2.4	1.97	1.8
30	63.5	33.6	23.5	18.3	4.1	2.4	1.95	1.7

### 11.33 Many Instruments

Some applications have available a large number  $\ell$  of instruments. If they are all valid, using a large number should reduce the asymptotic variance relative to estimation with a smaller number of instruments. Is it then good practice to use many instruments? Or is there a cost to this practice? Bekker (1994) initiated a large literature investigating this question by formalizing the idea of “many instruments”. Bekker proposed an asymptotic approximation which treats the number of instruments  $\ell$  as proportional to the sample size, that is  $\ell = \alpha n$ , or equivalently that  $\ell/N \rightarrow \alpha \in [0, 1)$ .

We examine this idea in the simplified setting of one endogenous regressor and no included exogenous regressors

$$\begin{aligned} y_i &= \beta x_i + e_i \\ x_i &= \mathbf{z}_i' \boldsymbol{\gamma} + u_i \end{aligned} \tag{11.81}$$

with  $\mathbf{z}_i$   $\ell \times 1$ . As in the previous two sections we make the simplifying assumption that the errors are conditionally homoskedastic and unit variance

$$\text{var} \left( \begin{pmatrix} e_i \\ u_i \end{pmatrix} \mid \mathbf{z}_i \right) = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}. \tag{11.82}$$

In addition we assume that the conditional fourth moments are bounded

$$\mathbb{E}(e_i^4 \mid \mathbf{z}_i) \leq C < \infty \quad \mathbb{E}(u_i^4 \mid \mathbf{z}_i) \leq C < \infty. \tag{11.83}$$

The idea that there are “many instruments” is formalized by the assumption that the number of instruments is increasing proportionately with the sample size

$$\frac{\ell}{n} \longrightarrow \alpha. \tag{11.84}$$

The best way to think about this is to view  $\alpha$  as the ratio of  $\ell$  to  $n$  in a given sample. Thus if an application has  $n = 100$  observations and  $\ell = 10$  instruments, then we should treat  $\alpha = 0.10$ .

Consider the variance of the endogenous regressor  $x_i$  from the reduced form:  $\text{var}(x_i) = \text{var}(\mathbf{z}_i' \boldsymbol{\gamma}) + \text{var}(u_i)$ . Suppose that  $\text{var}(x_i)$  and  $\text{var}(u_i)$  are unchanging as  $\ell$  increases. This implies that  $\text{var}(\mathbf{z}_i' \boldsymbol{\gamma})$  is unchanging as well. This will be a useful assumption, as it implies that the population  $R^2$  of the reduced form is not changing with  $\ell$ . We don't need this exact condition, rather we simply assume that the sample version converges in probability to a fixed constant

$$\frac{1}{n} \sum_{i=1}^n \boldsymbol{\gamma}' \mathbf{z}_i \mathbf{z}_i' \boldsymbol{\gamma} \xrightarrow{p} c \tag{11.85}$$

for  $0 < c < \infty$ . Again, this essentially implies that the  $R^2$  of the reduced form regression for  $x_i$  converges to a constant.

As a baseline it is useful to examine the behavior of the least-squares estimator of  $\beta$ . First, observe that the variances of  $n^{-1} \sum_{i=1}^n \boldsymbol{\gamma}' \mathbf{z}_i e_i$  and  $n^{-1} \sum_{i=1}^n \boldsymbol{\gamma}' \mathbf{z}_i u_i$ , conditional on  $\mathbf{Z}$ , are both equal to

$$n^{-2} \sum_{i=1}^n \boldsymbol{\gamma}' \mathbf{z}_i \mathbf{z}_i' \boldsymbol{\gamma} \xrightarrow{p} 0$$

by (11.85). Thus they converge in probability to zero:

$$n^{-1} \sum_{i=1}^n \boldsymbol{\gamma}' \mathbf{z}_i e_i \xrightarrow{p} 0 \tag{11.86}$$

and

$$n^{-1} \sum_{i=1}^n \gamma' \mathbf{z}_i \mathbf{u}_i \xrightarrow{p} 0. \quad (11.87)$$

Combined with (11.85) and the WLLN we find

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n x_i e_i &= \frac{1}{n} \sum_{i=1}^n \gamma' \mathbf{z}_i e_i + \frac{1}{n} \sum_{i=1}^n u_i e_i \xrightarrow{p} \rho \\ \frac{1}{n} \sum_{i=1}^n x_i^2 &= \frac{1}{n} \sum_{i=1}^n \gamma' \mathbf{z}_i \mathbf{z}_i' \gamma + \frac{2}{n} \sum_{i=1}^n \gamma' \mathbf{z}_i u_i + \frac{1}{n} \sum_{i=1}^n u_i^2 \xrightarrow{p} c + 1. \end{aligned}$$

Hence

$$\hat{\beta}_{\text{ols}} = \beta + \frac{\frac{1}{n} \sum_{i=1}^n x_i e_i}{\frac{1}{n} \sum_{i=1}^n x_i^2} \xrightarrow{p} \beta + \frac{\rho}{c + 1}.$$

Thus least-squares is inconsistent for  $\beta$  under endogeneity.

Now consider the 2SLS estimator. In matrix notation, setting  $\mathbf{P} = \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}'$ ,

$$\hat{\beta}_{2\text{sls}} - \beta = \frac{\frac{1}{n} \mathbf{X}' \mathbf{P} \mathbf{e}}{\frac{1}{n} \mathbf{X}' \mathbf{P} \mathbf{X}} = \frac{\frac{1}{n} \gamma' \mathbf{Z}' \mathbf{e} + \frac{1}{n} \mathbf{u}' \mathbf{P} \mathbf{e}}{\frac{1}{n} \gamma' \mathbf{Z}' \mathbf{Z} \gamma + \frac{2}{n} \gamma' \mathbf{Z}' \mathbf{u} + \frac{1}{n} \mathbf{u}' \mathbf{P} \mathbf{u}}. \quad (11.88)$$

In the expression on the right-side of (11.88), three of the components have been examined in (11.85), (11.86), and (11.87). We now examine the remaining components  $\frac{1}{n} \mathbf{u}' \mathbf{P} \mathbf{e}$  and  $\frac{1}{n} \mathbf{u}' \mathbf{P} \mathbf{e} = \mathbf{u}$ .

First, it is simple to take their expectations under the conditional homoskedasticity assumption. We have

$$\mathbb{E} \left( \frac{1}{n} \mathbf{u}' \mathbf{P} \mathbf{e} \right) = \frac{1}{n} \text{tr} \mathbb{E} (\mathbf{P} \mathbf{e} \mathbf{u}') = \frac{1}{n} \text{tr} (\mathbf{P}) \rho = \frac{\ell}{n} \rho \quad (11.89)$$

since  $\text{tr} (\mathbf{P}) = \ell$ . Similarly

$$\mathbb{E} \left( \frac{1}{n} \mathbf{u}' \mathbf{P} \mathbf{u} \right) = \frac{1}{n} \text{tr} \mathbb{E} (\mathbf{P} \mathbf{u} \mathbf{u}') = \frac{1}{n} \text{tr} (\mathbf{P}) = \frac{\ell}{n}.$$

Second, we examine their variances, which is a more cumbersome exercise. Let  $P_{ij} = \mathbf{z}_i' (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{z}_j$  be the  $ij^{\text{th}}$  element of  $\mathbf{P}$ . Then  $\mathbf{u}' \mathbf{P} \mathbf{e} = \sum_{i=1}^n \sum_{j=1}^n u_i e_j P_{ij}$  and  $\mathbf{u}' \mathbf{P} \mathbf{u} = \sum_{i=1}^n \sum_{j=1}^n u_i u_j P_{ij}$ .

The matrix  $\mathbf{P}$  is idempotent. It therefore has the properties  $\sum_{i=1}^n P_{ii} = \text{tr} (\mathbf{P}) = \ell$  and  $0 \leq P_{ii} \leq 1$ . The property  $\mathbf{P} \mathbf{P} = \mathbf{P}$  also implies  $\sum_{j=1}^n P_{ij}^2 = P_{ii}$ . Then

$$\begin{aligned} \text{var} \left( \frac{1}{n} \mathbf{u}' \mathbf{P} \mathbf{e} \right) &= \frac{1}{n^2} \mathbb{E} \left( \sum_{i=1}^n \sum_{j=1}^n (u_i e_j - \rho 1(i=j)) P_{ij} \right)^2 \\ &= \frac{1}{n^2} \mathbb{E} \left( \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n (u_i e_j - \rho 1(i=j)) P_{ij} (u_k e_l - \rho 1(k=l)) P_{kl} \right) \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left( (u_i e_i - \rho)^2 P_{ii}^2 \right) \end{aligned} \quad (11.90)$$

$$+ \frac{1}{n^2} \sum_{i=1}^n \sum_{j \neq i} \mathbb{E} (u_i^2 e_j^2 P_{ij}^2) \quad (11.91)$$

$$+ \frac{1}{n^2} \sum_{i=1}^n \sum_{j \neq i} \mathbb{E} (u_i e_j e_i u_j P_{ij}^2) \quad (11.92)$$

$$= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}(u_i^2 e_i^2 P_{ii}^2) - 2 \frac{\rho}{n^2} \sum_{i=1}^n \mathbb{E}(u_i e_i P_{ii}^2) + \rho^2 \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}(P_{ii}^2).$$

The third equality holds because the remaining cross-products have zero expectation since the observations are independent and the errors have zero mean. We then calculate that (11.90) is bounded by

$$(C - \rho^2) \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} P_{ii}^2 \leq (C - \rho^2) \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}(P_{ii}) = (C - \rho^2) \frac{\ell}{n^2} \rightarrow 0$$

under (11.84). The first inequality is  $P_{ii} \leq 1$  and the equality is  $\sum_{i=1}^n P_{ii} = \ell$ . Next, the conditional homoskedasticity assumption implies that (11.91) plus (11.92) equals  $(1 + \rho^2)$  times

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j \neq i} \mathbb{E}(P_{ij}^2) \leq \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}(P_{ij}^2) = \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}(P_{ii}) = \frac{\ell}{n^2} \rightarrow 0$$

under (11.84). The first equality is  $\sum_{j=1}^n P_{ij}^2 = P_{ii}$ . Together, we have shown that

$$\text{var} \left( \frac{1}{n} \mathbf{u}' \mathbf{P} \mathbf{e} \right) \rightarrow 0.$$

Using (11.89) and Markov's inequality

$$\frac{1}{n} \mathbf{u}' \mathbf{P} \mathbf{e} - \frac{\ell}{n} \rho \xrightarrow{p} 0.$$

Combined with (11.84) we find

$$\frac{1}{n} \mathbf{u}' \mathbf{P} \mathbf{e} \xrightarrow{p} \alpha \rho. \quad (11.93)$$

The analysis for  $\frac{1}{n} \mathbf{u}' \mathbf{P} \mathbf{u}$  is quite similar. We deduce that

$$\frac{1}{n} \mathbf{u}' \mathbf{P} \mathbf{u} \xrightarrow{p} \alpha. \quad (11.94)$$

Returning to the 2SLS estimator (11.88) and combining (11.85), (11.86), (11.87), (11.93) and (11.94), we find

$$\widehat{\beta}_{2\text{sls}} \xrightarrow{p} \beta + \frac{\alpha \rho}{c + \alpha}.$$

We can state this formally.

**Theorem 11.33.1** *In model (11.81), under assumptions (11.82), (11.83) and (11.84), then as  $n \rightarrow \infty$ .*

$$\widehat{\beta}_{\text{ols}} \xrightarrow{p} \beta + \frac{\rho}{c + 1}$$

$$\widehat{\beta}_{2\text{sls}} \xrightarrow{p} \beta + \frac{\alpha \rho}{c + \alpha}.$$

This result is quite insightful. It shows that while endogeneity ( $\rho \neq 0$ ) renders the least-squares estimator inconsistent, the 2SLS estimator is also inconsistent if the number of instruments diverges proportionately with  $n$ . The limit in Theorem 11.33.1 shows a continuity between least-squares and 2SLS. The probability limit of the 2SLS estimator is continuous in  $\alpha$ , with the extreme case ( $\alpha = 1$ )

implying that 2SLS and least-squares have the same probability limit. The general implication is that the inconsistency of 2SLS is increasing in  $\alpha$ .

Hence using a large number of instruments in an application comes at a cost.

In an application, users should calculate the “many instrument ratio”  $\alpha = \ell/n$ . Unfortunately there is no known rule-of-thumb for  $\alpha$  which should lead to acceptable inference, but a minimum criterion is that if  $\alpha \geq 0.05$  you should be seriously concerned about the many-instrument problem. In general, if it is desired to use a large number of instruments then it is recommended to use an estimation method other than 2SLS such as LIML.

### 11.34 Example: Acemoglu, Johnson and Robinson (2001)

One particularly well-cited instrument variable regression is in Acemoglu, Johnson and Robinson (2001) with additional details published in (2012). They are interested in the effect of political institutions on economic performance. The theory is that good institutions (rule-of-law, property rights) should result in a country having higher long-term economic output than if the same country had poor institutions. To investigate this question, they focus on a sample of 64 former European colonies. Their data is in the file `AJR2001` on the textbook website.

The authors’ premise is that modern political institutions will have been influenced by the colonizing country. In particular, they argue that colonizing countries tended to set up colonies as either an “extractive state” or as a “migrant colony”. An extractive state was used by the colonizer to extract resources for the colonizing country, but was not largely settled by the European colonists. In this case the colonists would have had no incentive to set up good political institutions. In contrast, if a colony was set up as a “migrant colony”, then large numbers of European settlers migrated to the colony to live. These settlers would have desired institutions similar to those in their home country, and hence would have had a positive incentive to set up good political institutions. The nature of institutions is quite persistent over time, so these 19<sup>th</sup>-century foundations would affect the nature of modern institutions. The authors conclude that the 19<sup>th</sup>-century nature of the colony should be predictive of the nature of modern institutions, and hence modern economic growth.

To start the investigation they report an OLS regression of log GDP per capita in 1995 on a measure of political institutions they call “risk”, which is a measure of the protection against expropriation risk. This variable ranges from 0 to 10, with 0 the lowest protection against appropriation, and 10 the highest. For each country the authors take the average value of the index over 1985 to 1995 (the mean is 6.5 with a standard deviation of 1.5). Their reported OLS estimates (intercept omitted) are

$$\log(\widehat{GDP \text{ per Capita}}) = \begin{matrix} 0.52 \\ (0.06) \end{matrix} \text{ risk.} \quad (11.95)$$

These estimates imply a 52% difference in GDP between countries with a 1-unit difference in *risk*.

The authors argue that the *risk* is likely endogenous, since economic output influences political institutions, and because the variable *risk* is undoubtedly measured with error. These issues induce least-square bias in different directions and thus the overall bias effect is unclear.

To correct for the endogeneity bias the authors argue the need for an instrumental variable which does not directly affect economic performance yet is associated with political institutions. Their innovative suggestion was to use the mortality rate which faced potential European settlers in the 19<sup>th</sup> century. Colonies with high expected mortality would have been less attractive to European settlers, resulting in lower levels of European migrants. As a consequence the authors expect such colonies to have been more likely structured as an extractive state rather than a migrant colony. To measure the expected mortality rate the authors use estimates provided by historical research of the annualized deaths per 1000 soldiers, labeled *mortality*. (They used military mortality rates



as the military maintained high-quality records.) The first-stage regression is

$$\begin{aligned} risk = & -0.61 \log(mortality) + \hat{u}. \\ & (0.13) \end{aligned} \tag{11.96}$$

These estimates confirm that 19<sup>th</sup>-century high settler mortality rates are associated with countries with lower quality modern institutions. Using  $\log(mortality)$  as an instrument for  $risk$ , they estimate the structural equation using 2SLS and report

$$\begin{aligned} \log(\widehat{GDP \text{ per Capita}}) = & 0.94 \ risk. \\ & (0.16) \end{aligned} \tag{11.97}$$

This estimate is much higher than the OLS estimate from (11.95). The estimate is consistent with a near doubling of GDP due to a 1-unit difference in the risk index.

These are simple regressions involving just one right-hand-side variable. The authors considered a range of other models. Included in these results are a reversal of a traditional finding. In a conventional (least-squares) regression two relevant variables for output are *latitude* (distance from the equator) and *africa* (a dummy variable for countries from Africa), both of which are difficult to interpret causally. But in the proposed instrumental variables regression the variables *latitude* and *africa* have much smaller – and statistically insignificant – coefficients.

To assess the specification, we can use the Stock-Yogo and endogeneity tests. The Stock-Yogo test is from the reduced form (11.96). The instrument has a t-ratio of 4.8 (or  $F = 23$ ) which exceeds the Stock-Yogo critical value and hence can be treated as strong. For an endogeneity test, we take the least-squares residual  $\hat{u}$  from this equation and include it in the structural equation and estimate by least-squares. We find a coefficient on  $\hat{u}$  of  $-0.57$  with a t-ratio of 4.7, which is highly significant. We conclude that the least-squares and 2SLS estimates are statistically different, and reject the hypothesis that the variable  $risk$  is exogenous for the GDP structural equation.

In Exercise 11.23 you will replicate and extend these results using the authors' data.

This paper is a creative and careful use of the instrumental variables method. The creativity stems from the careful historical analysis which lead to the focus on mortality as a potential predictor of migration choices. The care comes in the implementation, as the authors needed to gather country-level data on political institutions and mortality from distinct sources. Putting these pieces together is the art of the project.

### 11.35 Example: Angrist and Krueger (1991)

Another influential instrument variable regression is in Angrist and Krueger (1991). Their concern, similar to Card (1995), is estimation of the structural returns to education while treating educational attainment as endogenous. Like Card, their goal is to find an instrument which is exogenous for wages yet has an impact on educational attainment. A subset of their data in the file **AK1991** on the textbook website.

Their creative suggestion was to focus on compulsory school attendance policies and their interaction with birthdates. Compulsory schooling laws vary across states in the United States, but typically require that youth remain in school until their sixteenth or seventeenth birthday. Angrist and Krueger argue that compulsory schooling has a causal effect on wages – youth who would have chosen to drop out of school stay in school for more years – and thus have more education which causally impacts their earnings as adults.

Angrist and Krueger next observe that these policies have differential impact on youth who are born early or late in the school year. Students who are born early in the calendar year are typically older when they enter school. Consequently when they attain the legal dropout age they

have attended less school than those born near the end of the year. This means that birthdate (early in the calendar year versus late) exogenously impacts educational attainment, and thus wages through education. Yet birthdate must be exogenous for the structural wage equation, as there is no reason to believe that birthdate itself has a causal impact on a person's ability or wages. These considerations together suggest that birthdate is a valid instrumental variable for education in a causal wage equation.

Typical wage datasets include age, but not birthdates. To obtain information on birthdate, Angrist and Krueger used a U.S. Census data which includes an individual's quarter of birth (January-March, April-June, etc.). They use this variable to construct 2SLS estimates of the return to education.

Their paper carefully documents that educational attainment varies by quarter of birth (as predicted by the above discussion), and reports a large set of least-squares and 2SLS estimates. We focus on two estimates at the core of their analysis, reported in column (6) of their Tables V and VII. This involves data from the 1980 census with men born in 1930-1939, with 329,509 observations. The first equation is

$$\widehat{\log(wage)} = 0.080 \text{ edu} - 0.230 \text{ black} + 0.158 \text{ urban} + 0.244 \text{ married} \quad (11.98)$$

(0.016)                      (0.026)                      (0.017)                      (0.005)

where *edu* years of education, and *black*, *urban*, and *married* are dummy variables indicating race (1 if black, 0 otherwise), lives in a metropolitan area, and if married. In addition to the reported coefficients, the equation also includes as regressors nine year-of-birth dummies and eight region-of-residence dummies. The equation is estimated by 2SLS. The instrumental variables are the 30 interactions of three quarter-of-birth times ten year-of-birth dummy variables.

This equation indicates an 8% increase in wages due to each year of education.

Angrist and Krueger observe that the effect of compulsory education laws are likely to vary across states, so expand the instrument set to include interactions with state-of-birth. They estimate the following equation by 2SLS

$$\widehat{\log(wage)} = 0.083 \text{ edu} - 0.233 \text{ black} + 0.151 \text{ urban} + 0.244 \text{ married}. \quad (11.99)$$

(0.010)                      (0.011)                      (0.010)                      (0.003)

This equation also adds fifty state-of-birth dummy variables as regressors. The instrumental variables are the 180 interactions of quarter-of-birth times year-of-birth dummy variables, plus quarter-of-birth times state-of-birth interactions.

This equation shows a similar estimated causal effect of education on wages as in (11.98). More notably, the standard error is smaller in (11.99), suggesting improved precision by the expanded instrumental variable set.

However, these estimates seem excellent candidates for weak instruments and many instruments. Indeed, this paper (published in 1991) helped sparked these two literatures. We can use the Stock-Yogo tools to explore the instrument strength and the implications for the Angrist-Krueger estimates.

We first take equation (11.98). Using the original Angrist-Krueger data, we estimate the corresponding reduced form, and calculate the  $F$  statistic for the 30 excluded instruments. We find  $F = 4.7$ . It has an asymptotic p-value of 0.000, suggesting that we can reject (at any significance level) the hypothesis that the coefficients on the excluded instruments are zero. Thus Angrist and Krueger appear to be correct that quarter of birth helps to explain educational attainment and are thus a valid instrumental variable set. However, using the Stock-Yogo test,  $F = 4.7$  is not high enough to reject the hypothesis that the instruments are weak. Specifically, for  $\ell_2 = 30$  the critical value for the  $F$  statistic is 45 (if we want to bound size below 15%). The actual value of 4.7 is

far below 45. Since we cannot reject that the instruments are weak, this indicates that we cannot interpret the 2SLS estimates and test statistics in (11.98) as reliable.

Second, take (11.99) with the expanded regressor and instrument set. Estimating the corresponding reduced form, we find the  $F$  statistic for the 180 excluded instruments is  $F = 2.15$  which also has an asymptotic p-value of 0.000 indicating that we can reject at any significance level the hypothesis that the excluded instruments have no effect on educational attainment. However, using the Stock-Yogo test we also cannot reject the hypothesis that the instruments are weak. While Stock and Yogo did not calculate the critical values for  $\ell_2 = 180$ , the 2SLS critical values are increasing in  $\ell_2$  so we can use those for  $\ell_2 = 30$  as a lower bound. Hence the observed value of  $F = 2.15$  is far below the level needed for significance. Consequently the results in (11.99) cannot be viewed as reliable. In particular, the observation that the standard errors in (11.99) are smaller than those in (11.98) should not be interpreted as evidence of greater precision. Rather, they should be viewed as evidence of unreliability due to weak instruments.

When instruments are weak, one constructive suggestion is to use LIML estimation rather than 2SLS. Another constructive suggestion is to alter the instrument set. While Angrist and Krueger used a large number of instrumental variables, we can consider using a smaller set. Take equation (11.98). Rather than estimating it using the 30 interaction instruments, consider using only the three quarter-of-birth dummy variables. We report the reduced form estimates here:

$$\widehat{edu} = -1.57 \quad black + 1.05 \quad urban + 0.225 \quad married + 0.050 \quad Q_2 + 0.101 \quad Q_3 + 0.142 \quad Q_4 \\ (0.02) \quad (0.01) \quad (0.016) \quad (0.016) \quad (0.016) \quad (0.016) \quad (0.016) \quad (0.016) \quad (11.100)$$

where  $Q_2$ ,  $Q_3$  and  $Q_4$  are dummy variables for birth in the 2<sup>nd</sup>, 3<sup>rd</sup>, and 4<sup>th</sup> quarter. The regression also includes nine year-of-birth and eight region-of-residence dummy variables.

The reduced form coefficients in (11.100) on the quarter-of-birth dummies are quite instructive. The coefficients are positive and increasing, consistent with the Angrist-Krueger hypothesis that individuals born later in the year achieve higher average education. Focusing on the weak instrument problem, the  $F$  test for exclusion of these three variables is  $F = 30$ . The Stock-Yogo critical value is 12.8 for  $\ell_2 = 3$  and a size of 15%, and is 22.3 for a size of 10%. Since  $F = 30$  exceeds both these thresholds we can reject the hypothesis that this reduced form is weak. Estimating the model by 2SLS with these three instruments we find

$$\log(\widehat{wage}) = 0.098 \quad edu - 0.217 \quad black + 0.137 \quad urban + 0.240 \quad married. \quad (11.101) \\ (0.020) \quad (0.022) \quad (0.017) \quad (0.006)$$

These estimates indicate a slightly larger (10%) causal impact of education on wages, but with a larger standard error. The Stock-Yogo analysis indicates that we can interpret the confidence intervals from these estimates as having asymptotic coverage 85%.

While the original Angrist-Krueger estimates suffer due to weak instruments, their paper is a very creative and thoughtful application of the **natural experiment** methodology. They discovered a completely exogenous variation present in the world – birthdate – and showed how this has a small but measurable effect on educational attainment, and thereby on earnings. Their crafting of this natural experiment regression is extremely clever and demonstrates a style of analysis which can successfully underlie an effective instrumental variables empirical analysis.

### Joshua Angrist

Joshua Angrist (1960-) is an Israeli-American econometrician and labor economist who is known for his advocacy of natural experiments to motivate instrumental variables estimation. He is also well-known for his book *Mostly Harmless Econometrics* (2009) co-authored with Jörn-Steffen Pischke.

## 11.36 Programming

We now present Stata code for some of the empirical work reported in this chapter.

### Stata do File for Card Example

```
use Card1995.dta, clear
set more off
gen exp = age76 - ed76 - 6
gen exp2 = (exp^2)/100
* Drop observations with missing wage
drop if lwage76==.
* Least squares baseline
reg lwage76 ed76 exp exp2 smsa76r reg76r, r
* Reduced form estimates using college as instrument
reg lwage76 nearc4 exp exp2 smsa76r reg76r, r
reg ed76 nearc4 exp exp2 smsa76r reg76r, r
* IV estimates
ivregress 2sls lwage76 exp exp2 smsa76r reg76r (ed76=nearc4), r
* Reduced form using public and private as instruments
reg ed76 nearc4a nearc4b exp exp2 smsa76r reg76r, r
* F test for excluded instruments
testparm nearc4a nearc4b
predict u2, residual
* 2SLS estimates using both instruments
ivregress 2sls lwage76 exp exp2 smsa76r reg76r (ed76=nearc4a nearc4b), r
* Control function regressions
reg lwage76 ed76 exp exp2 smsa76r reg76r u2
reg lwage76 ed76 exp exp2 smsa76r reg76r u2, r
* LIML estimates
ivregress liml lwage76 exp exp2 smsa76r reg76r (ed76=nearc4a nearc4b), r
```

### Stata do File for Acemoglu-Johnson-Robinson Example

```
use AJR2001.dta, clear
reg loggdp risk
reg risk logmort0
predict u, residual
ivregress 2sls loggdp (risk=logmort0)
reg loggdp risk u
```

**Stata do File for Angrist-Krueger Example**

```
use AK1991.dta, clear
ivregress 2sls logwage black smsa married i.yob i.region (edu = i.qob#i.yob)
reg edu black smsa married i.yob i.region i.qob#i.yob
testparm i.qob#i.yob
ivregress 2sls logwage black smsa married i.yob i.region i.state (edu =
i.qob#i.yob i.qob#i.state)
reg edu black smsa married i.yob i.region i.state i.qob#i.yob i.qob#i.state
testparm i.qob#i.yob i.qob#i.state
reg edu black smsa married i.yob i.region i.qob
testparm i.qob
ivregress 2sls logwage black smsa married i.yob i.region (edu = i.qob)
```

## Exercises

**Exercise 11.1** Consider the single equation model

$$y_i = z_i\beta + e_i,$$

where  $y_i$  and  $z_i$  are both real-valued ( $1 \times 1$ ). Let  $\hat{\beta}$  denote the IV estimator of  $\beta$  using as an instrument a dummy variable  $d_i$  (takes only the values 0 and 1). Find a simple expression for the IV estimator in this context.

**Exercise 11.2** In the linear model

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(e_i \mid \mathbf{x}_i) &= 0 \end{aligned}$$

suppose  $\sigma_i^2 = \mathbb{E}(e_i^2 \mid \mathbf{x}_i)$  is known. Show that the GLS estimator of  $\boldsymbol{\beta}$  can be written as an IV estimator using some instrument  $\mathbf{z}_i$ . (Find an expression for  $\mathbf{z}_i$ .)

**Exercise 11.3** Take the linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}.$$

Let the OLS estimator for  $\boldsymbol{\beta}$  be  $\hat{\boldsymbol{\beta}}$  and the OLS residual be  $\hat{\mathbf{e}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ .

Let the IV estimator for  $\boldsymbol{\beta}$  using some instrument  $\mathbf{Z}$  be  $\tilde{\boldsymbol{\beta}}$  and the IV residual be  $\tilde{\mathbf{e}} = \mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}$ . If  $\mathbf{X}$  is indeed endogenous, will IV “fit” better than OLS, in the sense that  $\tilde{\mathbf{e}}'\tilde{\mathbf{e}} < \hat{\mathbf{e}}'\hat{\mathbf{e}}$ , at least in large samples?

**Exercise 11.4** The reduced form between the regressors  $\mathbf{x}_i$  and instruments  $\mathbf{z}_i$  takes the form

$$\mathbf{x}_i = \boldsymbol{\Gamma}' \mathbf{z}_i + \mathbf{u}_i$$

or

$$\mathbf{X} = \mathbf{Z}\boldsymbol{\Gamma} + \mathbf{U}$$

where  $\mathbf{x}_i$  is  $k \times 1$ ,  $\mathbf{z}_i$  is  $l \times 1$ ,  $\mathbf{X}$  is  $n \times k$ ,  $\mathbf{Z}$  is  $n \times l$ ,  $\mathbf{U}$  is  $n \times k$ , and  $\boldsymbol{\Gamma}$  is  $l \times k$ . The parameter  $\boldsymbol{\Gamma}$  is defined by the population moment condition

$$\mathbb{E}(\mathbf{z}_i \mathbf{u}_i') = \mathbf{0}$$

Show that the method of moments estimator for  $\boldsymbol{\Gamma}$  is  $\hat{\boldsymbol{\Gamma}} = (\mathbf{Z}'\mathbf{Z})^{-1}(\mathbf{Z}'\mathbf{X})$ .

**Exercise 11.5** In the structural model

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{e} \\ \mathbf{X} &= \mathbf{Z}\boldsymbol{\Gamma} + \mathbf{U} \end{aligned}$$

with  $\boldsymbol{\Gamma}$   $l \times k$ ,  $l \geq k$ , we claim that  $\boldsymbol{\beta}$  is identified (can be recovered from the reduced form) if  $\text{rank}(\boldsymbol{\Gamma}) = k$ . Explain why this is true. That is, show that if  $\text{rank}(\boldsymbol{\Gamma}) < k$  then  $\boldsymbol{\beta}$  cannot be identified.

**Exercise 11.6** For Theorem 11.16.1, establish that  $\hat{\mathbf{V}}_{\boldsymbol{\beta}} \xrightarrow{p} \mathbf{V}_{\boldsymbol{\beta}}$ .

**Exercise 11.7** Take the linear model

$$\begin{aligned} y_i &= x_i\beta + e_i \\ \mathbb{E}(e_i \mid x_i) &= 0. \end{aligned}$$

where  $x_i$  and  $\beta$  are  $1 \times 1$ .

- (a) Show that  $\mathbb{E}(x_i e_i) = 0$  and  $\mathbb{E}(x_i^2 e_i) = 0$ . Is  $\mathbf{z}_i = (x_i \ x_i^2)'$  a valid instrumental variable for estimation of  $\beta$ ?
- (b) Define the 2SLS estimator of  $\beta$ , using  $\mathbf{z}_i$  as an instrument for  $x_i$ . How does this differ from OLS?

**Exercise 11.8** Suppose that price and quantity are determined by the intersection of the linear demand and supply curves

$$\text{Demand : } Q = a_0 + a_1 P + a_2 Y + e_1$$

$$\text{Supply : } Q = b_0 + b_1 P + b_2 W + e_2$$

where income ( $Y$ ) and wage ( $W$ ) are determined outside the market. In this model, are the parameters identified?

**Exercise 11.9** Consider the model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(e_i | \mathbf{z}_i) = 0$$

with  $y_i$  scalar and  $\mathbf{x}_i$  and  $\mathbf{z}_i$  each a  $k$  vector. You have a random sample  $(y_i, \mathbf{x}_i, \mathbf{z}_i : i = 1, \dots, n)$ .

- (a) Suppose that  $\mathbf{x}_i$  is exogenous in the sense that  $E(e_i | \mathbf{z}_i, \mathbf{x}_i) = 0$ . Is the IV estimator  $\hat{\boldsymbol{\beta}}_{\text{iv}}$  unbiased for  $\boldsymbol{\beta}$ ?
- (b) Continuing to assume that  $\mathbf{x}_i$  is exogenous, find the variance matrix for  $\hat{\boldsymbol{\beta}}_{\text{iv}}$ ,  $\text{var}(\hat{\boldsymbol{\beta}}_{\text{iv}} | \mathbf{X}, \mathbf{Z})$ .

**Exercise 11.10** Consider the model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbf{x}_i = \boldsymbol{\Gamma}' \mathbf{z}_i + \mathbf{u}_i$$

$$\mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0}$$

$$\mathbb{E}(\mathbf{z}_i \mathbf{u}_i') = \mathbf{0}$$

with  $y_i$  scalar and  $\mathbf{x}_i$  and  $\mathbf{z}_i$  each a  $k$  vector. You have a random sample  $(y_i, \mathbf{x}_i, \mathbf{z}_i : i = 1, \dots, n)$ . Take the control function equation

$$e_i = \mathbf{u}_i' \boldsymbol{\gamma} + \varepsilon_i$$

$$\mathbb{E}(\mathbf{u}_i \varepsilon_i) = \mathbf{0}$$

and assume for simplicity that  $\mathbf{u}_i$  is observed. Inserting into the structural equation we find

$$y_i = \mathbf{z}_i' \boldsymbol{\beta} + \mathbf{u}_i' \boldsymbol{\gamma} + \varepsilon_i$$

The control function estimator  $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\gamma}})$  is OLS estimation of this equation.

- (a) Show that  $\mathbb{E}(\mathbf{x}_i \varepsilon_i) = \mathbf{0}$  (algebraically)
- (b) Derive the asymptotic distribution of  $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\gamma}})$ .

**Exercise 11.11** Consider the structural equation

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + e_i \tag{11.102}$$

with  $x_i$  treated as endogenous so that  $E(x_i e_i) \neq 0$ . Assume  $y_i$  and  $x_i$  are scalar. Suppose we also have a scalar instrument  $z_i$  which satisfies

$$\mathbb{E}(e_i | z_i) = 0$$

so in particular  $\mathbb{E}(e_i) = 0$ ,  $\mathbb{E}(z_i e_i) = 0$  and  $\mathbb{E}(z_i^2 e_i) = 0$ .

- (a) Should  $x_i^2$  be treated as endogenous or exogenous?
- (b) Suppose we have a scalar instrument  $z_i$  which satisfies

$$x_i = \gamma_0 + \gamma_1 z_i + u_i \quad (11.103)$$

with  $u_i$  independent of  $z_i$  and mean zero.

Consider using  $(1, z_i, z_i^2)$  as instruments. Is this a sufficient number of instruments? (Would this be just-identified, over-identified, or under-identified)?

- (c) Write out the reduced form equation for  $x_i^2$ . Under what condition on the reduced form parameters (11.103) are the parameters in (11.102) identified?

**Exercise 11.12** Consider the structural equation and reduced form

$$\begin{aligned} y_i &= \beta x_i^2 + e_i \\ x_i &= \gamma z_i + u_i \\ \mathbb{E}(z_i e_i) &= 0 \\ \mathbb{E}(z_i u_i) &= 0 \end{aligned}$$

with  $x_i^2$  treated as endogenous so that  $\mathbb{E}(x_i^2 e_i) \neq 0$ . For simplicity assume no intercepts.  $y_i$ ,  $z_i$ , and  $x_i$  are scalar. Assume  $\gamma \neq 0$ . Consider the following estimator. First, estimate  $\gamma$  by OLS of  $x_i$  on  $z_i$  and construct the fitted values  $\hat{x}_i = \hat{\gamma} z_i$ . Second, estimate  $\beta$  by OLS of  $y_i$  on  $\hat{x}_i^2$ .

- (a) Write out this estimator  $\hat{\beta}$  explicitly as a function of the sample
- (b) Find its probability limit as  $n \rightarrow \infty$
- (c) In general, is  $\hat{\beta}$  consistent for  $\beta$ ? Is there a reasonable condition under which  $\hat{\beta}$  is consistent?

**Exercise 11.13** Consider the structural equation

$$\begin{aligned} y_i &= \mathbf{x}'_{1i} \boldsymbol{\beta}_1 + \mathbf{x}'_{2i} \boldsymbol{\beta}_2 + e_i \\ \mathbb{E}(\mathbf{z}_i e_i) &= 0 \end{aligned}$$

where  $\mathbf{x}_{2i}$  is  $k_2 \times 1$  and treated as endogenous. The variables  $\mathbf{z}_i = (\mathbf{x}_{1i}, \mathbf{z}_{2i})$  are treated as exogenous, where  $\mathbf{z}_{2i}$  is  $\ell_2 \times 1$  and  $\ell_2 \geq k_2$ . You are interested in testing the hypothesis

$$\mathbb{H}_0 : \boldsymbol{\beta}_2 = 0.$$

Consider the reduced form equation for  $y_i$

$$y_i = \mathbf{x}'_{1i} \boldsymbol{\lambda}_1 + \mathbf{z}'_{2i} \boldsymbol{\lambda}_2 + v_i. \quad (11.104)$$

Show how to test  $\mathbb{H}_0$  using only the OLS estimates of (11.104).

Hint: This will require an analysis of the reduced form equations and their relation to the structural equation.

**Exercise 11.14** Take the linear instrumental variables equation

$$\begin{aligned} y_i &= \mathbf{x}'_{1i} \boldsymbol{\beta}_1 + \mathbf{x}'_{2i} \boldsymbol{\beta}_2 + e_i \\ \mathbb{E}(\mathbf{z}_i e_i) &= 0 \end{aligned}$$

where  $\mathbf{x}_{1i}$  is  $k_1 \times 1$ ,  $\mathbf{x}_{2i}$  is  $k_2 \times 1$ , and  $\mathbf{z}_i$  is  $\ell \times 1$ , with  $\ell \geq k = k_1 + k_2$ . The sample size is  $n$ . Assume that  $\mathbf{Q}_{zz} = \mathbb{E}(\mathbf{z}_i \mathbf{z}_i') > 0$  and  $\mathbf{Q}_{zx} = \mathbb{E}(\mathbf{z}_i \mathbf{x}_i')$  has full rank  $k$ .

Suppose that only  $(y_i, \mathbf{x}_{1i}, \mathbf{z}_i)$  are available, and  $\mathbf{x}_{2i}$  is missing from the dataset.

Consider the 2SLS estimator  $\hat{\boldsymbol{\beta}}_1$  of  $\boldsymbol{\beta}_1$  obtained from the misspecified IV regression, by regressing  $y_i$  on  $\mathbf{x}_{1i}$  only, using  $\mathbf{z}_i$  as an instrument for  $\mathbf{x}_{1i}$ .



- (a) Find a stochastic decomposition  $\widehat{\beta}_1 = \beta_1 + \mathbf{b}_{1n} + \mathbf{r}_{1n}$  where  $\mathbf{r}_{1n}$  depends on the error  $e_i$ , and  $\mathbf{b}_{1n}$  does not depend on the error  $e_i$ .
- (b) Show that  $\mathbf{r}_{1n} \rightarrow_p 0$  as  $n \rightarrow \infty$ .
- (c) Find the probability limit of  $\mathbf{b}_{1n}$  and  $\widehat{\beta}_1$  as  $n \rightarrow \infty$ .
- (d) Does  $\widehat{\beta}_1$  suffer from “omitted variables bias”? Explain. Under what conditions is there no omitted variables bias?
- (e) Find the asymptotic distribution as  $n \rightarrow \infty$  of

$$\sqrt{n} \left( \widehat{\beta}_1 - \beta_1 - \mathbf{b}_{1n} \right).$$

**Exercise 11.15** Take the linear instrumental variables equation

$$\begin{aligned} y_i &= x_i \beta_1 + z_i \beta_2 + e_i \\ \mathbb{E}(e_i | z_i) &= 0 \end{aligned}$$

where for simplicity both  $x_i$  and  $z_i$  are scalar  $1 \times 1$ .

- (a) Can the coefficients  $(\beta_1, \beta_2)$  be estimated by 2SLS using  $z_i$  as an instrument for  $x_i$ ? Why or why not?
- (b) Can the coefficients  $(\beta_1, \beta_2)$  be estimated by 2SLS using  $z_i$  and  $z_i^2$  as instruments?
- (c) For the 2SLS estimator suggested in (b), what is the implicit exclusion restriction?
- (d) In (b), what is the implicit assumption about instrument relevance?  
[Hint: Write down the implied reduced form equation for  $x_i$ .]
- (e) In a generic application, would you be comfortable with the assumptions in (c) and (d)?

**Exercise 11.16** Take a linear equation with endogeneity and a just-identified linear reduced form

$$\begin{aligned} y_i &= x_i \beta + e_i \\ x_i &= \gamma z_i + u_i \end{aligned}$$

where both  $x_i$  and  $z_i$  are scalar  $1 \times 1$ . Assume that

$$\begin{aligned} \mathbb{E}(z_i e_i) &= 0 \\ \mathbb{E}(z_i u_i) &= 0 \end{aligned}$$

- (a) Derive the reduced form equation

$$y_i = z_i \lambda + v_i.$$

Show that  $\beta = \lambda/\gamma$  if  $\gamma \neq 0$ , and that  $\mathbb{E}(z_i v_i) = 0$

- (b) Let  $\widehat{\lambda}$  denote the OLS estimate from linear regression of  $Y$  on  $Z$ , and let  $\widehat{\gamma}$  denote the OLS estimate from linear regression of  $X$  on  $Z$ . Write  $\theta = (\lambda, \gamma)'$  and let  $\widehat{\theta} = (\widehat{\lambda}, \widehat{\gamma})'$ . Define the error vector  $\boldsymbol{\xi}_i = \begin{pmatrix} v_i \\ u_i \end{pmatrix}$ . Write  $\sqrt{n}(\widehat{\theta} - \theta)$  using a single expression as a function of the error  $\boldsymbol{\xi}_i$ .
- (c) Show that  $\mathbb{E}(z_i \boldsymbol{\xi}_i) = 0$

- (d) Derive the joint asymptotic distribution of  $\sqrt{n}(\hat{\theta} - \theta)$  as  $n \rightarrow \infty$ . Hint: Define  $\mathbf{\Omega}_{\xi} = \mathbb{E}(z_i^2 \xi_i \xi_i')$
- (e) Using the previous result and the Delta Method, find the asymptotic distribution of the Indirect Least Squares estimator  $\hat{\beta} = \hat{\lambda}/\hat{\gamma}$
- (f) Is the answer in (e) the same as the asymptotic distribution of the 2SLS estimator in Theorem 11.14.1?

Hint: Show that  $\begin{pmatrix} 1 & -\beta \end{pmatrix} \xi_i = e_i$  and  $\begin{pmatrix} 1 & -\beta \end{pmatrix} \mathbf{\Omega}_{\xi} \begin{pmatrix} 1 \\ -\beta \end{pmatrix} = \mathbb{E}(z_i^2 e_i^2)$ .

**Exercise 11.17** Take the model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(z_i e_i) = 0$$

and consider the two-stage least-squares estimator. The first-stage estimate is

$$\widehat{\mathbf{X}} = \mathbf{Z} \widehat{\Gamma}$$

$$\widehat{\Gamma} = (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X}$$

and the second-stage is least-squares of  $y_i$  on  $\widehat{\mathbf{x}}_i$ :

$$\widehat{\beta} = \left( \widehat{\mathbf{X}}' \widehat{\mathbf{X}} \right)^{-1} \widehat{\mathbf{X}}' \mathbf{y}$$

with least-squares residuals

$$\widehat{\mathbf{e}} = \mathbf{y} - \widehat{\mathbf{X}} \widehat{\beta}.$$

Consider  $\widehat{\sigma}^2 = \frac{1}{n} \widehat{\mathbf{e}}' \widehat{\mathbf{e}}$  as an estimator for  $\sigma^2 = \mathbb{E}(e_i^2)$ . Is this appropriate? If not, propose an alternative estimator.

**Exercise 11.18** You have two independent iid samples  $(y_{1i}, \mathbf{x}_{1i}, \mathbf{z}_{1i} : i = 1, \dots, n)$  and  $(y_{2i}, \mathbf{x}_{2i}, \mathbf{z}_{2i} : i = 1, \dots, n)$ . The dependent variables  $y_{1i}$  and  $y_{2i}$  are real-valued. The regressors  $\mathbf{x}_{1i}$  and  $\mathbf{x}_{2i}$  and instruments  $\mathbf{z}_{1i}$  and  $\mathbf{z}_{2i}$  are  $k$ -vectors. The model is standard just-identified linear instrumental variables

$$y_{1i} = \mathbf{x}_{1i}' \boldsymbol{\beta}_1 + e_{1i}$$

$$\mathbb{E}(\mathbf{z}_{1i} e_{1i}) = \mathbf{0}$$

$$y_{2i} = \mathbf{x}_{2i}' \boldsymbol{\beta}_2 + e_{2i}$$

$$\mathbb{E}(\mathbf{z}_{2i} e_{2i}) = \mathbf{0}$$

For concreteness, sample 1 are women and sample 2 are men. You want to test  $\mathbb{H}_0 : \boldsymbol{\beta}_1 = \boldsymbol{\beta}_2$ , that the two samples have the same coefficients.

- (a) Develop a test statistic for  $\mathbb{H}_0$ .
- (b) Derive the asymptotic distribution of the test statistic.
- (c) Describe (in brief) the testing procedure.

**Exercise 11.19** To estimate  $\beta$  in the model  $y_i = x_i \beta + e_i$  with  $x_i$  scalar and endogenous, with household level data, you want to use as an the instrument the state of residence.

- (a) What are the assumptions needed to justify this choice of instrument?

(b) Is the model just identified or overidentified?

**Exercise 11.20** The model is

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0}$$

An economist wants to obtain the 2SLS estimates and standard errors for  $\boldsymbol{\beta}$ . He uses the following steps

- Regresses  $\mathbf{x}_i$  on  $\mathbf{z}_i$ , obtains the predicted values  $\hat{\mathbf{x}}_i$ .
- Regresses  $y_i$  on  $\hat{\mathbf{x}}_i$ , obtains the coefficient estimate  $\hat{\boldsymbol{\beta}}$  and standard error  $s(\hat{\boldsymbol{\beta}})$  from this regression.

Is this correct? Does this produce the 2SLS estimates and standard errors?

**Exercise 11.21** Let

$$y_i = \mathbf{x}_{1i}' \boldsymbol{\beta}_1 + \mathbf{x}_{2i}' \boldsymbol{\beta}_2 + e_i$$

Let  $(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\beta}}_2)$  denote the 2SLS estimates of  $(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2)$  when  $\mathbf{z}_{2i}$  is used as an instrument for  $\mathbf{x}_{2i}$  and they are the same dimension (so the model is just identified). Let  $(\hat{\boldsymbol{\lambda}}_1, \hat{\boldsymbol{\lambda}}_2)$  be the OLS estimates from the regression

$$y_i = \mathbf{x}_{1i}' \hat{\boldsymbol{\lambda}}_1 + \mathbf{z}_{2i}' \hat{\boldsymbol{\lambda}}_2 + e_i$$

Show that  $\hat{\boldsymbol{\beta}}_1 = \hat{\boldsymbol{\lambda}}_1$ .

**Exercise 11.22** In the linear model

$$y_i = x_i \beta + e_i$$

suppose  $\sigma_i^2 = E(e_i^2 | x_i)$  is known. Show that the GLS estimator of  $\beta$  can be written as an instrumental variables estimator using some instrument  $z_i$ . (Find an expression for  $z_i$ .)

**Exercise 11.23** You will replicate and extend the work reported in Acemoglu, Johnson and Robinson (2001). The authors provided an expanded set of controls when they published their 2012 extension and posted the data on the AER website. This dataset is **AJR2001** on the textbook website..

- Estimate the OLS regression (11.95), the reduced form regression (11.96) and the 2SLS regression (11.97). (Which point estimate is different by 0.01 from the reported values? This is a common phenomenon in empirical replication).
- For the above estimates, calculate both homoskedastic and heteroskedastic-robust standard errors. Which were used by the authors (as reported in (11.95)-(11.96)-(11.97)?)
- Calculate the 2SLS estimates by the Indirect Least Squares formula. Are they the same?
- Calculate the 2SLS estimates by the two-stage approach. Are they the same?
- Calculate the 2SLS estimates by the control variable approach. Are they the same?
- Acemoglu, Johnson and Robinson (2001) reported many specifications including alternative regressor controls, for example *latitude* and *africa*. Estimate by least-squares the equation for logGDP adding *latitude* and *africa* as regressors. Does this regression suggest that *latitude* and *africa* are predictive of the level of GDP?

- (g) Now estimate the same equation as in (f) but by 2SLS using log mortality as an instrument for *risk*. How does the interpretation of the effect of *latitude* and *africa* change?
- (h) Return to our baseline model (without including *latitude* and *africa*). The authors' reduced form equation uses log(mortality) as the instrument, rather than, say, the level of mortality. Estimate the reduced form for risk with *mortality* as the instrument. (This variable is not provided in the dataset, so you need to take the exponential of the mortality variable.) Can you explain why the authors preferred the equation with log(mortality)?
- (i) Try an alternative reduced form, including both log(mortality) and the square of log(mortality). Interpret the results. Re-estimate the structural equation by 2SLS using both log(mortality) and its square as instruments. How do the results change?
- (j) For the estimates in (i), are the instruments strong or weak using the Stock-Yogo test?
- (k) Calculate and interpret a test for exogeneity of the instruments.
- (l) Estimate the equation by LIML, using the instruments log(mortality) and the square of log(mortality).

**Exercise 11.24** You will replicate and extend the work reported in the chapter relating to Card (1995). The data is from the author's website, and is posted as `Card1995`. The model we focus on is labeled 2SLS(a) in Table 11.1, which uses *public* and *private* as instruments for *Edu*. The variables you will need for this exercise include *lwage76*, *ed76*, *age76*, *smsa76r*, *reg76r*, *nearc2*, *nearc4*, *nearc4a*, *nearc4b*. See the description file for definitions.

$$\log(Wage) = \beta_0 + \beta_1 Edu + \beta_2 Exp + \beta_3 Exp^2/100 + \beta_4 South + \beta_5 Black + e$$

where *Edu* = *Education* (Years), *Exp* = *Experience* (Years), and *South* and *Black* are regional and racial dummy variables. The variables *Exp* = *Age* - *Edu* - 6 and *Exp*<sup>2</sup>/100 are not in the dataset, they need to be generated.

- (a) First, replicate the reduced form regression presented in the final column of Table 11.2, and the 2SLS regression described above (using *public* and *private* as instruments for *Edu*) to verify that you have the same variable definitions.
- (b) Now try a different reduced form model. The variable *nearc2* means "grew up near a 2-year college". See if adding it to the reduced form equation is useful.
- (c) Now try more interactions in the reduced form. Create the interactions *nearc4a\*age76* and *nearc4a\*age76*<sup>2</sup>/100, and add them to the reduced form equation. Estimate this by least-squares. Interpret the coefficients on the two new variables.
- (d) Estimate the structural equation by 2SLS using the expanded instrument set  $\{nearc4a, nearc4b, nearc4a*age76, nearc4a*age76^2/100\}$ .  
What is the impact on the structural estimate of the return to schooling?
- (e) Using the Stock-Yogo test, are the instruments strong or weak?
- (f) Test the hypothesis that *Edu* is exogenous for the structural return to schooling.
- (g) Re-estimate the last equation by LIML. Do the results change meaningfully?

**Exercise 11.25** You will extend Angrist and Krueger (1991). In their Table VIII, they report their estimates of an analog of (11.99) for the subsample of 26,913 black men. Use this sub-sample for the following analysis.

- (a) Start by considering estimation of an equation which is identical in form to (11.99), with the same additional regressors (year-of-birth, region-of-residence, and state-of-birth dummy variables) and 180 excluded instrumental variables (the interactions of quarter-of-birth times year-of-birth dummy variables, and quarter-of-birth times state-of-birth interactions). But now, it is estimated on the subsample of black men. One regressor must be omitted to achieve identification. Which variable is this?
- (b) Estimate the reduced form for the above equation by least-squares. Calculate the  $F$  statistic for the excluded instruments. What do you conclude about the strength of the instruments?
- (c) Repeat, now estimating the reduced form for the analog of (11.98) which has 30 excluded instrumental variables, and does not include the state-of-birth dummy variables in the regression. What do you conclude about the strength of the instruments?
- (d) Repeat, now estimating the reduced form for the analog of (11.101) which has only 3 excluded instrumental variables. Are the instruments sufficiently strong for 2SLS estimation? For LIML estimation?
- (e) Estimate the structural wage equation using what you believe is the most appropriate set of regressors, instruments, and the most appropriate estimation method. What is the estimated return to education (for the subsample of black men) and its standard error? Without doing a formal hypothesis test, do these results (or in which way?) appear meaningfully different from the results for the full sample?

# Chapter 12

## Generalized Method of Moments

### 12.1 Moment Equation Models

All of the models that have been introduced so far can be written as **moment equation models**, where the population parameters solve a system of moment equations. Moment equation models are much broader than the models so far considered, and understanding their common structure opens up straightforward techniques to handle new econometric models.

Moment equation models take the following form. Let  $\mathbf{g}_i(\boldsymbol{\beta})$  be a known  $\ell \times 1$  function of the  $i^{\text{th}}$  observation and a  $k \times 1$  parameter  $\boldsymbol{\beta}$ . A moment equation model is summarized by the moment equations

$$\mathbb{E}(\mathbf{g}_i(\boldsymbol{\beta})) = \mathbf{0} \quad (12.1)$$

and a parameter space  $\boldsymbol{\beta} \in \mathbf{B}$ . For example, in the instrumental variables model  $\mathbf{g}_i(\boldsymbol{\beta}) = \mathbf{z}_i(y_i - \mathbf{x}_i'\boldsymbol{\beta})$ .

In general, we say that a parameter  $\boldsymbol{\beta}$  is **identified** if there is a unique mapping from the data distribution to  $\boldsymbol{\beta}$ . In the context of the model (12.1) this means that there is a unique  $\boldsymbol{\beta}$  satisfying (12.1). Since (12.1) is a system of  $\ell$  equations with  $k$  unknowns, then it is necessary that  $\ell \geq k$  for there to be a unique solution. If  $\ell = k$  we say that the model is **just identified**, meaning that there is just enough information to identify the parameters. If  $\ell > k$  we say that the model is **overidentified**, meaning that there is excess information (which can improve estimation efficiency). If  $\ell < k$  we say that the model is **underidentified**, meaning that there is insufficient information to identify the parameters. In general, we assume that  $\ell \geq k$  so the model is either just identified or overidentified.

### 12.2 Method of Moments Estimators

In this section we consider the just-identified case  $\ell = k$ .

Define the sample analog of (12.1)

$$\bar{\mathbf{g}}_n(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \mathbf{g}_i(\boldsymbol{\beta}). \quad (12.2)$$

The **method of moments estimator (MME)**  $\hat{\boldsymbol{\beta}}_{\text{mm}}$  for  $\boldsymbol{\beta}$  is defined as the parameter value which sets  $\bar{\mathbf{g}}_n(\boldsymbol{\beta}) = \mathbf{0}$ . Thus

$$\bar{\mathbf{g}}_n(\hat{\boldsymbol{\beta}}_{\text{mm}}) = \frac{1}{n} \sum_{i=1}^n \mathbf{g}_i(\hat{\boldsymbol{\beta}}_{\text{mm}}) = \mathbf{0}. \quad (12.3)$$

The equations (12.3) are known as the **estimating equations** as they are the equations which determine the estimator  $\hat{\boldsymbol{\beta}}_{\text{mm}}$ .

In some contexts (such as those discussed in the examples below), there is an explicit solution for  $\hat{\beta}_{\text{mm}}$ . In other cases the solution must be found numerically.

We now show how most of the estimators discussed so far in the textbook can be written as method of moments estimators.

**Mean:** Set  $g_i(\mu) = y_i - \mu$ . The MME is  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n y_i$ .

**Mean and Variance:** Set

$$g_i(\mu, \sigma^2) = \begin{pmatrix} y_i - \mu \\ (y_i - \mu)^2 - \sigma^2 \end{pmatrix}.$$

The MME are  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n y_i$  and  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\mu})^2$ .

**OLS:** Set  $g_i(\beta) = x_i(y_i - x_i'\beta)$ . The MME is  $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{y})$ .

**OLS and Variance:** Set

$$g_i(\beta, \sigma^2) = \begin{pmatrix} x_i(y_i - x_i'\beta) \\ (y_i - x_i'\beta)^2 - \sigma^2 \end{pmatrix}.$$

The MME is  $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{y})$  and  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - x_i'\hat{\beta})^2$ .

**Multivariate Least Squares, vector form:** Set  $g_i(\beta) = \mathbf{X}_i(\mathbf{y}_i - \mathbf{X}_i'\beta)$ . The MME is  $\hat{\beta} = (\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i')^{-1} (\sum_{i=1}^n \mathbf{X}_i \mathbf{y}_i)$  which is (10.3).

**Multivariate Least Squares, matrix form:** Set  $g_i(\mathbf{B}) = \text{vec}(\mathbf{x}_i(\mathbf{y}_i' - \mathbf{x}_i'\mathbf{B}))$ . The MME is  $\hat{\mathbf{B}} = (\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i')^{-1} (\sum_{i=1}^n \mathbf{x}_i \mathbf{y}_i')$  which is (10.5).

**Seemingly Unrelated Regression:** Set

$$g_i(\beta, \Sigma) = \begin{pmatrix} \mathbf{X}_i \Sigma^{-1} (\mathbf{y}_i - \mathbf{X}_i' \beta) \\ \text{vec}(\Sigma - (\mathbf{y}_i - \mathbf{X}_i' \beta)(\mathbf{y}_i - \mathbf{X}_i' \beta)') \end{pmatrix}.$$

The MME is  $\hat{\beta} = (\sum_{i=1}^n \mathbf{X}_i \hat{\Sigma}^{-1} \mathbf{X}_i')^{-1} (\sum_{i=1}^n \mathbf{X}_i \hat{\Sigma}^{-1} \mathbf{y}_i)$  and  $\hat{\Sigma} = n^{-1} \sum_{i=1}^n (\mathbf{y}_i - \mathbf{X}_i' \hat{\beta})(\mathbf{y}_i - \mathbf{X}_i' \hat{\beta})'$ .

**IV:** Set  $g_i(\beta) = \mathbf{z}_i(y_i - \mathbf{x}_i'\beta)$ . The MME is  $\hat{\beta} = (\sum_{i=1}^n \mathbf{z}_i \mathbf{x}_i')^{-1} (\sum_{i=1}^n \mathbf{z}_i y_i)$ .

**Generated Regressors:** Set

$$g_i(\beta, \mathbf{A}) = \begin{pmatrix} \mathbf{A}' \mathbf{z}_i (y_i - \mathbf{z}_i' \mathbf{A} \beta) \\ \text{vec}(\mathbf{z}_i (\mathbf{x}_i' - \mathbf{z}_i' \mathbf{A})) \end{pmatrix}.$$

The MME is  $\hat{\mathbf{A}} = (\sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i')^{-1} (\sum_{i=1}^n \mathbf{z}_i \mathbf{x}_i')$  and  $\hat{\beta} = (\hat{\mathbf{A}}' \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i' \hat{\mathbf{A}})^{-1} (\hat{\mathbf{A}}' \sum_{i=1}^n \mathbf{z}_i y_i)$ .

A common feature unifying these examples is that the estimator can be written as the solution to a set of estimating equations (12.3). This provides a common framework which enables a convenient development of a unified distribution theory.

### 12.3 Overidentified Moment Equations

In the instrumental variables model  $g_i(\beta) = z_i(y_i - x_i'\beta)$ . Thus (12.2) is

$$\bar{g}_n(\beta) = \frac{1}{n} \sum_{i=1}^n g_i(\beta) = \frac{1}{n} \sum_{i=1}^n z_i(y_i - x_i'\beta) = \frac{1}{n} (\mathbf{Z}'\mathbf{y} - \mathbf{Z}'\mathbf{X}\beta). \quad (12.4)$$

We have defined the method of moments estimator for  $\beta$  as the parameter value which sets  $\bar{g}_n(\beta) = \mathbf{0}$ . However, when the model is overidentified (if  $\ell > k$ ) then this is generally impossible as there are more equations than free parameters. Equivalently, there is no choice of  $\beta$  which sets (12.4) to zero. Thus the method of moments estimator is not defined for the overidentified case.

While we cannot find an estimator which sets  $\bar{g}_n(\beta)$  equal to zero, we can try to find an estimator which makes  $\bar{g}_n(\beta)$  as close to zero as possible. Let's think what that means. Since  $\bar{g}_n(\beta)$  is an  $\ell \times 1$  vector, this means we are trying to find a value for  $\beta$  which sets  $\bar{g}_n(\beta)$  as close as possible to the zero vector.

One way to think about this is to define the vector  $\mu = \mathbf{Z}'\mathbf{y}$ , the matrix  $\mathbf{G} = \mathbf{Z}'\mathbf{X}$  and the “error”  $\eta = \mu - \mathbf{G}\beta$ . Then we can write (12.4) as

$$\mu = \mathbf{G}\beta + \eta.$$

This looks like a regression equation with the  $\ell \times 1$  dependent variable  $\mu$ , the  $\ell \times k$  regressor matrix  $\mathbf{G}$ , and the  $\ell \times 1$  error vector  $\eta$ . Recall, the goal is to make the error vector  $\eta$  as small as possible. Recalling our knowledge about least-squares, we know that a simple method is to use least-squares regression of  $\mu$  on  $\mathbf{G}$ , which minimizes the sum-of-squares  $\eta'\eta$ . This is certainly one way to make  $\eta$  “small”. This least-squares solution is  $\hat{\beta} = (\mathbf{G}'\mathbf{G})^{-1}(\mathbf{G}'\mu)$ .

More generally, we know that when errors are non-homogeneous it can be more efficient to estimate by weighted least squares. Thus for some weight matrix  $\mathbf{W}$ , consider the estimator

$$\begin{aligned} \hat{\beta} &= (\mathbf{G}'\mathbf{W}\mathbf{G})^{-1}(\mathbf{G}'\mathbf{W}\mu) \\ &= (\mathbf{X}'\mathbf{Z}\mathbf{W}\mathbf{Z}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{Z}\mathbf{W}\mathbf{Z}'\mathbf{y}). \end{aligned}$$

This minimizes the weighted sum of squares  $\eta'\mathbf{W}\eta$ . This solution is known as the generalized method of moments (GMM).

The estimator is typically defined as follows. Given a set of moment equations (12.2) and an  $\ell \times \ell$  weight matrix  $\mathbf{W} > 0$ , the GMM criterion function is defined as

$$J(\beta) = n \cdot \bar{g}_n(\beta)'\mathbf{W}\bar{g}_n(\beta).$$

The factor “ $n$ ” is not important for the definition of the estimator, but is convenient for the distribution theory. The criterion  $J(\beta)$  is the weighted sum of squared moment equation errors. When  $\mathbf{W} = \mathbf{I}_\ell$ , then  $J(\beta) = n \cdot \bar{g}_n(\beta)'\bar{g}_n(\beta) = n \cdot \|\bar{g}_n(\beta)\|^2$ , the square of the Euclidean length. Since we restrict attention to positive definite weight matrices  $\mathbf{W}$ , the criterion  $J(\beta)$  is always non-negative.

The **Generalized Method of Moments (GMM)** estimator is defined as the minimizer of the GMM criterion  $J(\beta)$ .

**Definition 12.3.1** *The Generalized Method of Moments estimator is*  

$$\hat{\beta}_{\text{gmm}} = \underset{\beta}{\operatorname{argmin}} J_n(\beta).$$



Recall that in the just-identified case  $k = \ell$ , the method of moments estimator  $\hat{\beta}_{\text{mm}}$  solves  $\bar{g}_n(\hat{\beta}_{\text{mm}}) = \mathbf{0}$ . Hence in this case  $J_n(\hat{\beta}_{\text{mm}}) = 0$  which means that  $\hat{\beta}_{\text{mm}}$  minimizes  $J_n(\beta)$  and equals  $\hat{\beta}_{\text{gmm}} = \hat{\beta}_{\text{mm}}$ . This means that GMM includes MME as a special case. This implies that all of our results for GMM will apply to any method of moments estimators as a special case.

In the over-identified case the GMM estimator will depend on the choice of weight matrix  $\mathbf{W}$  and so this is an important focus of the theory. In the just-identified case, the GMM estimator simplifies to the MME which does not depend on  $\mathbf{W}$ .

The method and theory of the generalized method of moments was developed in an influential paper by Lars Hansen (1982). This paper introduced the method, its asymptotic distribution, the form of the efficient weight matrix, and tests for overidentification.

### Lars Peter Hansen

Lars Hansen (1952-) is an American econometrician and macroeconomist. In econometrics, he is famously known for the GMM estimator which has transformed theoretical and empirical economics. He was awarded the Nobel Memorial Prize in Economics in 2013.

## 12.4 Linear Moment Models

One of the great advantages of the moment equation framework is that it allows both linear and nonlinear models. However, when the moment equations are linear in the parameters then we have explicit solutions for the estimates and a straightforward asymptotic distribution theory. Hence we start by confining attention to linear moment equations, and return to nonlinear moment equations later. In the examples listed earlier, the estimators which have linear moment equations include the sample mean, OLS, multivariate least squares, IV, and 2SLS. The estimates which have non-linear moment equations include the sample variance, SUR, and generated regressors.

In particular, we focus on the overidentified IV model

$$\mathbf{g}_i(\beta) = \mathbf{z}_i(y_i - \mathbf{x}_i'\beta) \quad (12.5)$$

where  $\mathbf{z}_i$  is  $\ell \times 1$  and  $\mathbf{x}_i$  is  $k \times 1$ .

## 12.5 GMM Estimator

Given (12.5) the sample moment equations are (12.4). The GMM criterion can be written as

$$J(\beta) = n (\mathbf{Z}'\mathbf{y} - \mathbf{Z}'\mathbf{X}\beta)' \mathbf{W} (\mathbf{Z}'\mathbf{y} - \mathbf{Z}'\mathbf{X}\beta).$$

The GMM estimator minimizes  $J(\beta)$ . The first order conditions are

$$\begin{aligned} \mathbf{0} &= \frac{\partial}{\partial \beta} J(\hat{\beta}) \\ &= 2 \frac{\partial}{\partial \beta} \bar{\mathbf{g}}_n(\hat{\beta})' \mathbf{W} \bar{\mathbf{g}}_n(\hat{\beta}) \\ &= -2 \left( \frac{1}{n} \mathbf{X}'\mathbf{Z} \right) \mathbf{W} \left( \frac{1}{n} \mathbf{Z}' (\mathbf{y} - \mathbf{X}\hat{\beta}) \right). \end{aligned}$$

The solution is given as follows.

**Theorem 12.5.1** *For the overidentified IV model*

$$\hat{\beta}_{\text{gmm}} = (\mathbf{X}'\mathbf{Z}\mathbf{W}\mathbf{Z}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{Z}\mathbf{W}\mathbf{Z}'\mathbf{y}). \quad (12.6)$$

While the estimator depends on  $\mathbf{W}$ , the dependence is only up to scale. This is because if  $\mathbf{W}$  is replaced by  $c\mathbf{W}$  for some  $c > 0$ ,  $\hat{\beta}_{\text{gmm}}$  does not change.

When  $\mathbf{W}$  is fixed by the user, we call  $\hat{\beta}_{\text{gmm}}$  a **one-step GMM estimator**.

The GMM estimator (12.6) resembles the 2SLS estimator (11.34). In fact they are equal when  $\mathbf{W} = (\mathbf{Z}'\mathbf{Z})^{-1}$ . This means that the 2SLS estimator is a one-step GMM estimator for the linear model. In the just-identified case it also simplifies to the IV estimator (11.29).

**Theorem 12.5.2** *If  $\mathbf{W} = (\mathbf{Z}'\mathbf{Z})^{-1}$  then  $\hat{\beta}_{\text{gmm}} = \hat{\beta}_{\text{2sls}}$ . Furthermore, if  $k = \ell$  then  $\hat{\beta}_{\text{gmm}} = \hat{\beta}_{\text{iv}}$ .*

## 12.6 Distribution of GMM Estimator

Let

$$\mathbf{Q} = \mathbb{E}(\mathbf{z}_i \mathbf{x}_i')$$

and

$$\mathbf{\Omega} = \mathbb{E}(\mathbf{z}_i \mathbf{z}_i' e_i^2) = \mathbb{E}(\mathbf{g}_i \mathbf{g}_i')$$

where  $\mathbf{g}_i = \mathbf{z}_i e_i$ . Then

$$\left(\frac{1}{n}\mathbf{X}'\mathbf{Z}\right) \mathbf{W} \left(\frac{1}{n}\mathbf{Z}'\mathbf{X}\right) \xrightarrow{p} \mathbf{Q}'\mathbf{W}\mathbf{Q}$$

and

$$\left(\frac{1}{n}\mathbf{X}'\mathbf{Z}\right) \mathbf{W} \left(\frac{1}{\sqrt{n}}\mathbf{Z}'\mathbf{e}\right) \xrightarrow{d} \mathbf{Q}'\mathbf{W} \cdot \mathbf{N}(\mathbf{0}, \mathbf{\Omega}).$$

We conclude:

**Theorem 12.6.1** *Asymptotic Distribution of GMM Estimator.*  
Under Assumption 11.14.1, as  $n \rightarrow \infty$

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{V}_\beta)$$

where

$$\mathbf{V}_\beta = (\mathbf{Q}'\mathbf{W}\mathbf{Q})^{-1} (\mathbf{Q}'\mathbf{W}\mathbf{\Omega}\mathbf{W}\mathbf{Q}) (\mathbf{Q}'\mathbf{W}\mathbf{Q})^{-1}. \quad (12.7)$$

We find that the GMM estimator is asymptotically normal with a “sandwich form” asymptotic variance.

Our derivation treated the weight matrix  $\mathbf{W}$  as if it is non-random, but Theorem 12.6.1 carries over to the case where the weight matrix  $\hat{\mathbf{W}}$  is random so long as it converges in probability to some positive definite limit  $\mathbf{W}$ . This may require scaling the weight matrix, for example replacing  $\hat{\mathbf{W}} = (\mathbf{Z}'\mathbf{Z})^{-1}$  with  $\hat{\mathbf{W}} = (n^{-1}\mathbf{Z}'\mathbf{Z})^{-1}$ . Since rescaling the weight matrix does not affect the estimator this is ignored in implementation.

## 12.7 Efficient GMM

The asymptotic distribution of the GMM estimator  $\hat{\beta}_{\text{gmm}}$  depends on the weight matrix  $\mathbf{W}$  through the asymptotic variance  $\mathbf{V}_{\beta}$ . The asymptotically optimal weight matrix  $\mathbf{W}_0$  is one which minimizes  $\mathbf{V}_{\beta}$ . This turns out to be  $\mathbf{W}_0 = \mathbf{\Omega}^{-1}$ . The proof is left to Exercise 12.4.

When the GMM estimator  $\hat{\beta}$  is constructed with  $\mathbf{W} = \mathbf{W}_0 = \mathbf{\Omega}^{-1}$  (or a weight matrix which is a consistent estimator of  $\mathbf{W}_0$ ) we call it the **Efficient GMM** estimator:

$$\hat{\beta}_{\text{gmm}} = (\mathbf{X}'\mathbf{Z}\mathbf{\Omega}^{-1}\mathbf{Z}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{Z}\mathbf{\Omega}^{-1}\mathbf{Z}'\mathbf{y}).$$

Its asymptotic distribution takes a simpler form than in Theorem 12.6.1. By substituting  $\mathbf{W} = \mathbf{W}_0 = \mathbf{\Omega}^{-1}$  into (12.7) we find

$$\mathbf{V}_{\beta} = (\mathbf{Q}'\mathbf{\Omega}^{-1}\mathbf{Q})^{-1} (\mathbf{Q}'\mathbf{\Omega}^{-1}\mathbf{\Omega}\mathbf{\Omega}^{-1}\mathbf{Q}) (\mathbf{Q}'\mathbf{\Omega}^{-1}\mathbf{Q})^{-1} = (\mathbf{Q}'\mathbf{\Omega}^{-1}\mathbf{Q})^{-1}.$$

This is the asymptotic variance of the efficient GMM estimator.

**Theorem 12.7.1** *Asymptotic Distribution of GMM with Efficient Weight Matrix.* Under Assumption 11.14.1 and  $\mathbf{\Omega} > 0$ , as  $n \rightarrow \infty$

$$\sqrt{n}(\hat{\beta}_{\text{gmm}} - \beta) \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{V}_{\beta})$$

where

$$\mathbf{V}_{\beta} = (\mathbf{Q}'\mathbf{\Omega}^{-1}\mathbf{Q})^{-1}.$$

**Theorem 12.7.2** *Efficient GMM.* Under Assumption 11.14.1 and  $\mathbf{\Omega} > 0$ , for any  $\mathbf{W} > 0$ ,

$$(\mathbf{Q}'\mathbf{W}\mathbf{Q})^{-1} (\mathbf{Q}'\mathbf{W}\mathbf{\Omega}\mathbf{W}\mathbf{Q}) (\mathbf{Q}'\mathbf{W}\mathbf{Q})^{-1} - (\mathbf{Q}'\mathbf{\Omega}^{-1}\mathbf{Q})^{-1} > 0.$$

Thus if  $\hat{\beta}_{\text{gmm}}$  is the efficient GMM estimator and  $\tilde{\beta}_{\text{gmm}}$  is another GMM estimator, then

$$\text{avar}(\hat{\beta}_{\text{gmm}}) \leq \text{avar}(\tilde{\beta}_{\text{gmm}}).$$

For a proof, see Exercise 12.4.

This means that the smallest possible GMM covariance matrix (in the positive definite sense) is achieved by the efficient GMM weight matrix.

$\mathbf{W}_0 = \mathbf{\Omega}^{-1}$  is not known in practice but it can be estimated consistently as we discuss in Section 12.9. For any  $\hat{\mathbf{W}} \xrightarrow{p} \mathbf{W}_0$ , the asymptotic distribution in Theorem 12.7.1 is unaffected. Consequently we still call any  $\hat{\beta}_{\text{gmm}}$  constructed with an estimate of the efficient weight matrix an efficient GMM estimator.

By “efficient”, we mean that this estimator has the smallest asymptotic variance in the class of GMM estimators with this set of moment conditions. This is a weak concept of optimality, as we are only considering alternative weight matrices  $\hat{\mathbf{W}}$ . However, it turns out that the GMM estimator is semiparametrically efficient as shown by Gary Chamberlain (1987). If it is known that  $\mathbb{E}(g(\mathbf{w}_i, \beta)) = \mathbf{0}$ , and this is all that is known, this is a semi-parametric problem as the

distribution of the data is unknown. Chamberlain showed that in this context no semiparametric estimator (one which is consistent globally for the class of models considered) can have a smaller asymptotic variance than  $(\mathbf{G}'\boldsymbol{\Omega}^{-1}\mathbf{G})^{-1}$  where  $\mathbf{G} = \mathbb{E}\left(\frac{\partial}{\partial\boldsymbol{\beta}}\mathbf{g}_i(\boldsymbol{\beta})\right)$ . Since the GMM estimator has this asymptotic variance, it is semiparametrically efficient.

The results in this section show that in the linear model no estimator has better asymptotic efficiency than the efficient linear GMM estimator. No estimator can do better (in this first-order asymptotic sense), without imposing additional assumptions.

## 12.8 Efficient GMM versus 2SLS

For the linear model we introduced the 2SLS estimator as a standard estimator for  $\boldsymbol{\beta}$ . Now we have introduced the GMM estimator which includes 2SLS as a special case. Is there a context where 2SLS is efficient?

To answer this question, recall that the 2SLS estimator is GMM given the weight matrix  $\widehat{\mathbf{W}} = (\mathbf{Z}'\mathbf{Z})^{-1}$  or equivalently  $\widehat{\mathbf{W}} = (n^{-1}\mathbf{Z}'\mathbf{Z})^{-1}$  since scaling doesn't matter. Since  $\widehat{\mathbf{W}} \xrightarrow{p} (\mathbb{E}(\mathbf{z}_i\mathbf{z}_i'))^{-1}$ , this is asymptotically equivalent to using the weight matrix  $\mathbf{W} = (\mathbb{E}(\mathbf{z}_i\mathbf{z}_i'))^{-1}$ . In contrast, the efficient weight matrix takes the form  $(\mathbb{E}(\mathbf{z}_i\mathbf{z}_i'e_i^2))^{-1}$ . Now suppose that the structural equation error  $e_i$  is conditionally homoskedastic in the sense that  $\mathbb{E}(e_i^2 | \mathbf{z}_i) = \sigma^2$ . Then the efficient weight matrix equals  $\mathbf{W} = (\mathbb{E}(\mathbf{z}_i\mathbf{z}_i'))^{-1}\sigma^{-2}$ , or equivalently  $\mathbf{W} = (\mathbb{E}(\mathbf{z}_i\mathbf{z}_i'))^{-1}$  since scaling doesn't matter. The latter weight matrix is the same as the 2SLS asymptotic weight matrix. This shows that the 2SLS weight matrix is the efficient weight matrix under conditional homoskedasticity.

**Theorem 12.8.1** *Under Assumption 11.14.1 and  $\mathbb{E}(e_i^2 | \mathbf{z}_i) = \sigma^2$  then  $\widehat{\boldsymbol{\beta}}_{2\text{sls}}$  is efficient GMM.*

This shows that 2SLS is efficient under homoskedasticity. When homoskedasticity holds, there is no reason to use efficient GMM over 2SLS. More broadly, when homoskedasticity is a reasonable approximation then 2SLS will be a reasonable estimator. However, this result also shows that in the general case where the error is conditionally heteroskedastic, then 2SLS is generically inefficient relative to efficient GMM.

## 12.9 Estimation of the Efficient Weight Matrix

To construct the efficient GMM estimator we need a consistent estimator  $\widehat{\mathbf{W}}$  of  $\mathbf{W}_0 = \boldsymbol{\Omega}^{-1}$ . The convention is to form an estimate  $\widehat{\boldsymbol{\Omega}}$  of  $\boldsymbol{\Omega}$  and then set  $\widehat{\mathbf{W}} = \widehat{\boldsymbol{\Omega}}^{-1}$ .

The **two-step GMM estimator** proceeds by using a one-step consistent estimate of  $\boldsymbol{\beta}$  to construct the weight matrix estimator  $\widehat{\mathbf{W}}$ . In the linear model the natural one-step estimator for  $\boldsymbol{\beta}$  is the 2SLS estimator  $\widehat{\boldsymbol{\beta}}_{2\text{sls}}$ . Set  $\tilde{e}_i = y_i - \mathbf{x}_i'\widehat{\boldsymbol{\beta}}_{2\text{sls}}$ ,  $\tilde{\mathbf{g}}_i = \mathbf{g}_i(\boldsymbol{\beta}) = \mathbf{z}_i\tilde{e}_i$  and  $\bar{\mathbf{g}}_n = n^{-1}\sum_{i=1}^n \tilde{\mathbf{g}}_i$ . Two moment estimators of  $\boldsymbol{\Omega}$  are then

$$\widehat{\boldsymbol{\Omega}} = \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{g}}_i \tilde{\mathbf{g}}_i' \quad (12.8)$$

and

$$\widehat{\boldsymbol{\Omega}}^* = \frac{1}{n} \sum_{i=1}^n (\tilde{\mathbf{g}}_i - \bar{\mathbf{g}}_n)(\tilde{\mathbf{g}}_i - \bar{\mathbf{g}}_n)'. \quad (12.9)$$

The estimator (12.8) is an uncentered covariance matrix estimator while the estimator (12.9) is a centered version. Either estimator is consistent when  $\mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0}$  which holds under correct

specification. However under misspecification we may have  $\mathbb{E}(z_i e_i) \neq \mathbf{0}$ . In the latter context  $\widehat{\Omega}^*$  may be viewed as a robust estimator. For some testing problems it turns out to be preferable to use a covariance matrix estimator which is robust to the alternative hypothesis. For these reasons estimator (12.9) is generally preferred. Unfortunately, estimator (12.8) is more commonly seen in practice since it is the default choice by most packages. It is also worth observing that when the model is just identified then  $\bar{g}_n = \mathbf{0}$  so the two are algebraically identically.

Given the choice of covariance matrix estimator we set  $\widehat{W} = \widehat{\Omega}^{-1}$  or  $\widehat{W} = \widehat{\Omega}^{*-1}$ . Given this weight matrix, we then construct the **two-step GMM estimator** as (12.6) using the weight matrix  $\widehat{W}$ .

Since the 2SLS estimator is consistent for  $\beta$ , by arguments nearly identical to those used for covariance matrix estimation, we can show that  $\widehat{\Omega}$  and  $\widehat{\Omega}^*$  are consistent for  $\Omega$  and thus  $\widehat{W}$  is consistent for  $\Omega^{-1}$ . See Exercise 12.3.

This also means that the two-step GMM estimator satisfies the conditions for Theorem 12.7.1. We have established.

**Theorem 12.9.1** *Under Assumption 11.14.1 and  $\Omega > 0$ , if  $\widehat{W} = \widehat{\Omega}^{-1}$  or  $\widehat{W} = \widehat{\Omega}^{*-1}$  where the latter are defined in (12.8) and (12.9) then as  $n \rightarrow \infty$*

$$\sqrt{n}(\widehat{\beta}_{\text{gmm}} - \beta) \xrightarrow{d} N(\mathbf{0}, V_\beta)$$

where

$$V_\beta = (Q' \Omega^{-1} Q)^{-1}.$$

This shows that the two-step GMM estimator is asymptotically efficient.

The two-step GMM estimator of the IV regression equation can be computed in Stata using the `ivregress gmm` command. By default it uses formula (12.8). The centered version (12.9) may be selected using the `center` option.

## 12.10 Iterated GMM

The asymptotic distribution of the two-step GMM estimator does not depend on the choice of the preliminary one-step estimator. However, the actual value of the estimator depends on this choice, and so will the finite sample distribution. This is undesirable and likely inefficient. To remove this dependence we can iterate the estimation sequence. Specifically, given  $\widehat{\beta}_{\text{gmm}}$  we can construct an updated weight matrix estimate  $\widehat{W}$  and then re-estimate  $\widehat{\beta}_{\text{gmm}}$ . This updating can be iterated until convergence<sup>1</sup>. The result is called the **iterated GMM estimator** and is a common implementation of efficient GMM.

Interestingly, B. Hansen and Lee (2018) show that the iterated GMM estimator is unaffected if the weight matrix is computed with or without centering. Standard errors and test statistics, however, will be affected by the choice.

The iterated GMM estimator of the IV regression equation can be computed in Stata using the `ivregress gmm` command using the `igmm` option.

<sup>1</sup>In practice, “convergence” obtains when the difference between the estimates obtained at subsequent steps is smaller than a pre-specified tolerance. A sufficient condition for convergence is that the sequence is a contraction mapping. Indeed, B. Hansen and Lee (2018) have shown that the iterated GMM estimator generally satisfies this condition in large samples.

## 12.11 Covariance Matrix Estimation

An estimator of the asymptotic variance of  $\hat{\beta}_{\text{gmm}}$  can be obtained by replacing the matrices in the asymptotic variance formula by consistent estimates.

For the one-step GMM estimator the covariance matrix estimator is

$$\hat{\mathbf{V}}_{\beta} = \left( \hat{\mathbf{Q}}' \hat{\mathbf{W}} \hat{\mathbf{Q}} \right)^{-1} \left( \hat{\mathbf{Q}}' \hat{\mathbf{W}} \hat{\Omega} \hat{\mathbf{W}} \hat{\mathbf{Q}} \right) \left( \hat{\mathbf{Q}}' \hat{\mathbf{W}} \hat{\mathbf{Q}} \right)^{-1}$$

where

$$\hat{\mathbf{Q}} = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{x}_i'$$

and using either the uncentered estimator (12.8) or centered estimator (12.9) with the residuals  $\hat{e}_i = y_i - \mathbf{x}_i' \hat{\beta}_{\text{gmm}}$ .

For the two-step or iterated gmm estimator the covariance matrix estimator is

$$\hat{\mathbf{V}}_{\beta} = \left( \hat{\mathbf{Q}}' \hat{\Omega}^{-1} \hat{\mathbf{Q}} \right)^{-1} = \left( \left( \frac{1}{n} \mathbf{X}' \mathbf{Z} \right) \hat{\Omega}^{-1} \left( \frac{1}{n} \mathbf{Z}' \mathbf{X} \right) \right)^{-1}. \quad (12.10)$$

Again,  $\hat{\Omega}$  can be computed using either the uncentered estimator (12.8) or centered estimator (12.9), but should use the final residuals  $\hat{e}_i = y_i - \mathbf{x}_i' \hat{\beta}_{\text{gmm}}$ .

Asymptotic standard errors are given by the square roots of the diagonal elements of  $n^{-1} \hat{\mathbf{V}}_{\beta}$ .

In Stata, the default covariance matrix estimation method is determined by the choice of weight matrix. Thus if the centered estimator (12.9) is used for the weight matrix, it is also used for the covariance matrix estimator.

## 12.12 Clustered Dependence

In Section 4.20 we introduced clustered dependence and in Section 11.21 described covariance matrix estimation for 2SLS. The methods extend naturally to GMM, but with the additional complication of potentially altering weight matrix calculation.

As before, the structural equation for the  $g^{\text{th}}$  cluster can be written as the matrix system

$$\mathbf{y}_g = \mathbf{X}_g \beta + \mathbf{e}_g.$$

Using this notation the centered GMM estimator with weight matrix  $\mathbf{W}$  can be written as

$$\hat{\beta}_{\text{gmm}} = (\mathbf{X}' \mathbf{Z} \mathbf{W} \mathbf{Z}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Z} \mathbf{W} \left( \sum_{g=1}^G \mathbf{Z}_g' \mathbf{e}_g \right).$$

The cluster-robust covariance matrix estimator for  $\hat{\beta}_{\text{gmm}}$  is then

$$\hat{\mathbf{V}}_{\beta} = (\mathbf{X}' \mathbf{Z} \mathbf{W} \mathbf{Z}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Z} \mathbf{W} \hat{\mathbf{S}} \mathbf{W} \mathbf{Z}' \mathbf{X} (\mathbf{X}' \mathbf{Z} \mathbf{W} \mathbf{Z}' \mathbf{X})^{-1} \quad (12.11)$$

with

$$\hat{\mathbf{S}} = \sum_{g=1}^G \mathbf{Z}_g' \hat{\mathbf{e}}_g \hat{\mathbf{e}}_g' \mathbf{Z}_g \quad (12.12)$$

and the clustered residuals

$$\hat{\mathbf{e}}_g = \mathbf{y}_g - \mathbf{X}_g \hat{\beta}_{\text{gmm}}. \quad (12.13)$$

The cluster-robust estimator (12.11) is appropriate for the one-step GMM estimator. It is also appropriate for the two-step and iterated estimators when the latter use a conventional (non-clustered) efficient weight matrix. However in the clustering context it is more natural to use a

cluster-robust weight matrix such as  $\mathbf{W} = \hat{\mathbf{S}}^{-1}$  where  $\hat{\mathbf{S}}$  is a cluster-robust covariance estimator as in (12.12) based on a one-step or iterated residual. This gives rise to the cluster-robust GMM estimator

$$\hat{\beta}_{\text{gmm}} = \left( \mathbf{X}' \mathbf{Z} \hat{\mathbf{S}}^{-1} \mathbf{Z}' \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{Z} \hat{\mathbf{S}}^{-1} \mathbf{Z}' \mathbf{y}. \quad (12.14)$$

For this estimator an appropriate cluster-robust covariance matrix estimator is

$$\hat{\mathbf{V}}_{\beta} = \left( \mathbf{X}' \mathbf{Z} \hat{\mathbf{S}}^{-1} \mathbf{Z}' \mathbf{X} \right)^{-1}$$

where  $\hat{\mathbf{S}}$  is calculated using the final residuals.

To implement a cluster-robust weight matrix, use the 2SLS estimator for first step estimator. Compute the cluster residuals (12.13) and covariance matrix (12.12). Then (12.14) is the two-step GMM estimator. Updating the residuals and covariance matrix, we can iterate the sequence to obtain the iterated GMM estimator.

In Stata, using the `ivregress gmm` command with the `cluster` option implements the two-step GMM estimator using the cluster-robust weight matrix and cluster-robust covariance matrix estimator. To use the centered covariance matrix use the `center` option, and to implement the iterated GMM estimator use the `igmm` option. Alternatively, you can use the `wmatrix` and `vce` options to separately specify the weight matrix and covariance matrix estimation methods.

## 12.13 Wald Test

For a given function  $\mathbf{r}(\beta) : \mathbb{R}^k \rightarrow \Theta \subset \mathbb{R}^q$  we define the parameter  $\theta = \mathbf{r}(\beta)$ . The GMM estimator of  $\theta$  is  $\hat{\theta}_{\text{gmm}} = \mathbf{r}(\hat{\beta}_{\text{gmm}})$ . By the delta method it is asymptotically normal with covariance matrix

$$\begin{aligned} \mathbf{V}_{\theta} &= \mathbf{R}' \mathbf{V}_{\beta} \mathbf{R} \\ \mathbf{R} &= \frac{\partial}{\partial \beta} \mathbf{r}(\beta)'. \end{aligned}$$

An estimator of the asymptotic covariance matrix is

$$\begin{aligned} \hat{\mathbf{V}}_{\theta} &= \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\beta} \hat{\mathbf{R}} \\ \hat{\mathbf{R}} &= \frac{\partial}{\partial \beta} \mathbf{r}(\hat{\beta}_{\text{gmm}})'. \end{aligned}$$

When  $\theta$  is scalar then an asymptotic standard error for  $\hat{\theta}_{\text{gmm}}$  is formed as  $\sqrt{n^{-1} \hat{\mathbf{V}}_{\theta}}$ .

A standard test of the hypothesis

$$\mathbb{H}_0 : \theta = \theta_0$$

against

$$\mathbb{H}_1 : \theta \neq \theta_0$$

is based on the Wald statistic

$$W = n \left( \hat{\theta} - \theta_0 \right)' \hat{\mathbf{V}}_{\hat{\theta}}^{-1} \left( \hat{\theta} - \theta_0 \right).$$

Let  $G_q(u)$  denote the  $\chi_q^2$  distribution function.

**Theorem 12.13.1** *Under Assumption 11.14.1 and  $\Omega > 0$ , if  $\mathbf{r}(\boldsymbol{\beta})$  is continuously differentiable at  $\boldsymbol{\beta}$ , and  $\mathbb{H}_0$  holds, then as  $n \rightarrow \infty$ ,*

$$W \xrightarrow{d} \chi_q^2.$$

*For  $c$  satisfying  $\alpha = 1 - G_q(c)$ ,*

$$\Pr(W > c \mid \mathbb{H}_0) \longrightarrow \alpha$$

*so the test “Reject  $\mathbb{H}_0$  if  $W > c$ ” has asymptotic size  $\alpha$ .*

In Stata, the commands `test` and `testparm` can be used after `ivregress gmm` to implement Wald tests of linear hypotheses. The commands `nlcom` and `testnl` can be used after `ivregress gmm` to implement Wald tests of nonlinear hypotheses.

## 12.14 Restricted GMM

It is often desirable to impose restrictions on the coefficients. In this section we consider estimation subject to the constraints  $\mathbf{r}(\boldsymbol{\beta}) = \mathbf{0}$ .

The constrained GMM estimator minimizes the GMM criterion subject to the constraint. It is defined as

$$\hat{\boldsymbol{\beta}}_{\text{cgmm}} = \underset{\mathbf{r}(\boldsymbol{\beta})=\mathbf{0}}{\operatorname{argmin}} J(\boldsymbol{\beta}).$$

This is the parameter vector which makes the estimating equations as close to zero as possible with respect to the weighted quadratic distance while imposing the restriction on the parameters.

It is useful to separately consider the cases where  $\mathbf{r}(\boldsymbol{\beta})$  are linear and nonlinear.

First let's consider the linear case, where  $\mathbf{r}(\boldsymbol{\beta}) = \mathbf{R}'\boldsymbol{\beta} - \mathbf{c}$ . Using the methods of Chapter 8 it is straightforward to derive that given any weight matrix  $\mathbf{W}$  the constrained GMM estimator is

$$\hat{\boldsymbol{\beta}}_{\text{cgmm}} = \hat{\boldsymbol{\beta}}_{\text{gmm}} - (\mathbf{X}'\mathbf{Z}\mathbf{W}\mathbf{Z}'\mathbf{X})^{-1} \mathbf{R} \left( \mathbf{R}' (\mathbf{X}'\mathbf{Z}\mathbf{W}\mathbf{Z}'\mathbf{X})^{-1} \mathbf{R} \right)^{-1} \left( \mathbf{R}' \hat{\boldsymbol{\beta}}_{\text{gmm}} - \mathbf{c} \right). \quad (12.15)$$

In particular, when the efficient weight matrix  $\mathbf{W} = \hat{\boldsymbol{\Omega}}^{-1}$  is used the constrained GMM estimator can be written as

$$\hat{\boldsymbol{\beta}}_{\text{cgmm}} = \hat{\boldsymbol{\beta}}_{\text{gmm}} - \hat{\mathbf{V}}_{\boldsymbol{\beta}} \mathbf{R} \left( \mathbf{R}' \hat{\mathbf{V}}_{\boldsymbol{\beta}} \mathbf{R} \right)^{-1} \left( \mathbf{R}' \hat{\boldsymbol{\beta}}_{\text{gmm}} - \mathbf{c} \right) \quad (12.16)$$

which is the same formula (8.28) as efficient minimum distance.

To derive the asymptotic distribution under the assumption that the restriction is true, make the substitution  $\mathbf{c} = \mathbf{R}'\boldsymbol{\beta}$  in (12.15) to find

$$\sqrt{n} \left( \hat{\boldsymbol{\beta}}_{\text{cgmm}} - \boldsymbol{\beta} \right) = \left( \mathbf{I}_k - (\mathbf{X}'\mathbf{Z}\mathbf{W}\mathbf{Z}'\mathbf{X})^{-1} \mathbf{R} \left( \mathbf{R}' (\mathbf{X}'\mathbf{Z}\mathbf{W}\mathbf{Z}'\mathbf{X})^{-1} \mathbf{R} \right)^{-1} \mathbf{R}' \right) \sqrt{n} \left( \hat{\boldsymbol{\beta}}_{\text{gmm}} - \boldsymbol{\beta} \right). \quad (12.17)$$

which is a linear function of  $\sqrt{n} \left( \hat{\boldsymbol{\beta}}_{\text{gmm}} - \boldsymbol{\beta} \right)$ . Since the asymptotic distribution of the latter is known, it is straightforward to derive that of  $\sqrt{n} \left( \hat{\boldsymbol{\beta}}_{\text{cgmm}} - \boldsymbol{\beta} \right)$ . We present the result for the efficient case in Theorem 12.14.1 below.

Second, let's consider the nonlinear case, meaning that  $\mathbf{r}(\boldsymbol{\beta})$  is not an affine function of  $\boldsymbol{\beta}$ . In this case there is (in general) no explicit solution for  $\hat{\boldsymbol{\beta}}_{\text{cgmm}}$ . Instead, the solution needs to be found numerically. Fortunately there are excellent nonlinear constrained optimization solvers which make the task quite feasible. We do not review these here, but can be found in any numerical software system.



For the asymptotic distribution assume again that the restriction  $\mathbf{r}(\boldsymbol{\beta}) = \mathbf{0}$  is true. Then, using the same methods as in the proof of Theorem 8.14.1 we can show that (12.15) approximately holds, in the sense that

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_{\text{cgmm}} - \boldsymbol{\beta}) = \left( \mathbf{I}_k - (\mathbf{X}'\mathbf{Z}\mathbf{W}\mathbf{Z}'\mathbf{X})^{-1} \mathbf{R} \left( \mathbf{R}'(\mathbf{X}'\mathbf{Z}\mathbf{W}\mathbf{Z}'\mathbf{X})^{-1} \mathbf{R} \right)^{-1} \mathbf{R}' \right) \sqrt{n}(\hat{\boldsymbol{\beta}}_{\text{gmm}} - \boldsymbol{\beta}) + o_p(1) \quad (12.18)$$

where  $\mathbf{R} = \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{r}(\boldsymbol{\beta})'$ . Thus the asymptotic distribution of the constrained estimator takes the same form as in the linear case.

**Theorem 12.14.1** *Under Assumptions 11.14.1 and 8.14.1, and  $\boldsymbol{\Omega} > \mathbf{0}$ , for the efficient constrained GMM estimator (12.16)*

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_{\text{cgmm}} - \boldsymbol{\beta}) \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{V}_{\text{cgmm}})$$

as  $n \rightarrow \infty$ , where

$$\mathbf{V}_{\text{cgmm}} = \mathbf{V}_{\boldsymbol{\beta}} - \mathbf{V}_{\boldsymbol{\beta}} \mathbf{R} (\mathbf{R}' \mathbf{V}_{\boldsymbol{\beta}} \mathbf{R})^{-1} \mathbf{R}' \mathbf{V}_{\boldsymbol{\beta}}.$$

The asymptotic covariance matrix is estimated by

$$\begin{aligned} \hat{\mathbf{V}}_{\text{cgmm}} &= \tilde{\mathbf{V}}_{\boldsymbol{\beta}} - \tilde{\mathbf{V}}_{\boldsymbol{\beta}} \hat{\mathbf{R}} \left( \hat{\mathbf{R}}' \tilde{\mathbf{V}}_{\boldsymbol{\beta}} \hat{\mathbf{R}} \right)^{-1} \hat{\mathbf{R}}' \tilde{\mathbf{V}}_{\boldsymbol{\beta}}. \\ \tilde{\mathbf{V}}_{\boldsymbol{\beta}} &= \left( \hat{\mathbf{Q}}' \tilde{\boldsymbol{\Omega}}^{-1} \hat{\mathbf{Q}} \right)^{-1} \\ \tilde{\boldsymbol{\Omega}} &= \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i' \tilde{e}_i^2 \\ \tilde{e}_i &= y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{\text{cgmm}} \\ \hat{\mathbf{R}} &= \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{r}(\hat{\boldsymbol{\beta}}_{\text{cgmm}})'. \end{aligned}$$

## 12.15 Constrained Regression

Take the conventional projection model

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= \mathbf{0}. \end{aligned}$$

We can view this as a very special case of GMM. It is model (12.5) with  $\mathbf{z}_i = \mathbf{x}_i$ . This is just-identified GMM and the estimator is least-squares  $\hat{\boldsymbol{\beta}}_{\text{gmm}} = \hat{\boldsymbol{\beta}}_{\text{ols}}$ .

In Chapter 8 we discussed estimation of the projection model subject to linear constraints  $\mathbf{R}'\boldsymbol{\beta} = \mathbf{c}$ , which includes exclusion restrictions. Since the projection model is a special case of GMM, the constrained projection model is also constrained GMM. From the results of the previous section we find that the efficient constrained GMM estimator is

$$\hat{\boldsymbol{\beta}}_{\text{cgmm}} = \hat{\boldsymbol{\beta}}_{\text{ols}} - \hat{\mathbf{V}}_{\boldsymbol{\beta}} \mathbf{R} \left( \mathbf{R}' \hat{\mathbf{V}}_{\boldsymbol{\beta}} \mathbf{R} \right)^{-1} \left( \mathbf{R}' \hat{\boldsymbol{\beta}}_{\text{ols}} - \mathbf{c} \right) = \hat{\boldsymbol{\beta}}_{\text{emd}},$$

the efficient minimum distance estimator. Thus for linear constraints on the linear projection model, efficient GMM equals efficient minimum distance. Thus one convenient method to implement efficient minimum distance is by using GMM methods.

## 12.16 Distance Test

As in Section 12.13 consider testing the hypothesis  $\mathbb{H}_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$  where  $\boldsymbol{\theta} = \mathbf{r}(\boldsymbol{\beta})$  for a given function  $\mathbf{r}(\boldsymbol{\beta}) : \mathbb{R}^k \rightarrow \Theta \subset \mathbb{R}^q$ . When  $\mathbf{r}(\boldsymbol{\beta})$  is non-linear, a better approach than the Wald statistic is use a criterion-based statistic. This is sometimes called the GMM Distance statistic and sometimes called a LR-like statistic (the LR is for likelihood-ratio). The idea was first put forward by Newey and West (1987).

The idea is to compare the unrestricted and restricted estimators by contrasting the criterion functions. The unrestricted estimator takes the form

$$\hat{\boldsymbol{\beta}}_{\text{gmm}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} J(\boldsymbol{\beta})$$

where

$$\hat{J}(\boldsymbol{\beta}) = n \cdot \bar{\mathbf{g}}_n(\boldsymbol{\beta})' \hat{\boldsymbol{\Omega}}^{-1} \bar{\mathbf{g}}_n(\boldsymbol{\beta})$$

is the unrestricted GMM criterion which depends on an efficient weight matrix estimate  $\hat{\boldsymbol{\Omega}}$ . The minimized value of the criterion is

$$\hat{J} = \hat{J}(\hat{\boldsymbol{\beta}}_{\text{gmm}}).$$

As in Section 12.14, the estimator subject to  $\mathbf{r}(\boldsymbol{\beta}) = \boldsymbol{\theta}_0$  is

$$\hat{\boldsymbol{\beta}}_{\text{cgmm}} = \underset{\mathbf{r}(\boldsymbol{\beta})=\boldsymbol{\theta}_0}{\operatorname{argmin}} \tilde{J}(\boldsymbol{\beta})$$

where

$$\tilde{J}(\boldsymbol{\beta}) = n \cdot \bar{\mathbf{g}}_n(\boldsymbol{\beta})' \tilde{\boldsymbol{\Omega}}^{-1} \bar{\mathbf{g}}_n(\boldsymbol{\beta})$$

which depends on an efficient weight matrix estimate  $\tilde{\boldsymbol{\Omega}}$ . One possibility is to set  $\tilde{\boldsymbol{\Omega}} = \hat{\boldsymbol{\Omega}}$ . The minimized value of the criterion is

$$\tilde{J} = \tilde{J}(\hat{\boldsymbol{\beta}}_{\text{cgmm}}).$$

The GMM distance (or LR-like) statistic is the difference in the criterions

$$D = \tilde{J} - \hat{J}.$$

The distance test shares the useful feature of LR tests in that it is a natural by-product of the computation of alternative models.

The test has the following large sample distribution.

**Theorem 12.16.1** *Under Assumptions 11.14.1 and 8.14.1,  $\boldsymbol{\Omega} > 0$ , and  $\mathbb{H}_0$  holds, then as  $n \rightarrow \infty$ ,*

$$D \xrightarrow{d} \chi_q^2.$$

*For  $c$  satisfying  $\alpha = 1 - G_q(c)$ ,*

$$\Pr(D > c \mid \mathbb{H}_0) \longrightarrow \alpha$$

*so the test “Reject  $\mathbb{H}_0$  if  $D > c$ ” has asymptotic size  $\alpha$ .*

The proof is given in Section 12.24.

Theorem 12.16.1 shows that the distance statistic has a large sample distribution similar to that of Wald and likelihood ratio statistics, and can be interpreted in much the same way. Small values of  $D$  mean that imposing the restriction does not result in a large value of the moment equations. Hence the restrictions appear to be compatible with the data. On the other hand, large values

of  $D$  mean that imposing the restriction results in a much larger value of the moment equations, implying that the restrictions do not appear to be compatible with the data. The finding that the asymptotic distribution is chi-squared means that it is simple to obtain asymptotic critical values and p-values for the test.

We now discuss the choice of weight matrix. As mentioned above, one simple choice is to set  $\tilde{\Omega} = \hat{\Omega}$ . In this case we have the following result.

**Theorem 12.16.2** *If  $\tilde{\Omega} = \hat{\Omega}$  then  $D \geq 0$ . Furthermore, if  $\mathbf{r}$  is linear in  $\beta$ , then  $D$  equals the Wald statistic.*

The statement that  $\tilde{\Omega} = \hat{\Omega}$  implies  $D \geq 0$  follows from the fact that in this case the criterion functions  $\hat{J}(\beta) = \tilde{J}(\beta)$  are identical, so the constrained minimum cannot be smaller than the unconstrained. The statement that linear hypotheses and an efficient weight matrix implies  $D = W$  follows from applying the expression for the constrained GMM estimator (12.16) and using the variance matrix formula (12.10).

This result shows some advantages to using the same weight matrix to estimate both  $\hat{\beta}_{\text{gmm}}$  and  $\hat{\beta}_{\text{cgmm}}$ . In particular, the non-negativity finding motivated Newey and West (1987) to recommend using  $\tilde{\Omega} = \hat{\Omega}$ . However, this is not an important advantage. Alternatively, we can set  $\tilde{\Omega} = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i' \tilde{e}_i^2$  where  $\tilde{e}_i$  are residuals using the constrained estimator. This seems rather natural as in this case  $\hat{J}$  and  $\tilde{J}$  are simple outputs from iterated gmm. In the event that  $D < 0$  the test simply fails to reject  $\mathbb{H}_0$  at any significance level.

As discussed in Section 9.17, for tests of nonlinear hypotheses the Wald statistic can work quite poorly. In particular, the Wald statistic is affected by how the hypothesis  $\mathbf{r}(\beta)$  is formulated. In contrast, the distance statistic  $D$  is not affected by the algebraic formulation of the hypothesis. Current evidence suggests that the  $D$  statistic appears to have good sampling properties, and is a preferred test statistic relative to the Wald statistic for nonlinear hypotheses.

In Stata, the command `estat overid` after `ivregress gmm` can be used to report the value of the GMM criterion  $J$ . By estimating the two nested GMM regressions the values  $\hat{J}$  and  $\tilde{J}$  can be obtained and  $D$  computed.

## 12.17 Continuously-Updated GMM

An alternative to the two-step GMM estimator can be constructed by letting the weight matrix be an explicit function of  $\beta$ . This leads to the criterion function

$$J(\beta) = n \cdot \bar{\mathbf{g}}_n(\beta)' \left( \frac{1}{n} \sum_{i=1}^n \mathbf{g}(\mathbf{w}_i, \beta) \mathbf{g}(\mathbf{w}_i, \beta)' \right)^{-1} \bar{\mathbf{g}}_n(\beta).$$

The  $\hat{\beta}$  which minimizes this function is called the **continuously-updated GMM (CU-GMM) estimator**, and was introduced by L. Hansen, Heaton and Yaron (1996).

A complication is that the continuously-updated criterion  $J(\beta)$  is not quadratic in  $\beta$ . This means that minimization requires numerical methods. It may appear that the CU-GMM estimator is the same as the iterated GMM estimator, but this is not the case at all. They solve distinct first-order conditions, and can be quite different in applications.

Relative to traditional GMM, the CU-GMM estimator has lower bias but thicker distributional tails. While it has received considerable theoretical attention, it is not used commonly in applications.

## 12.18 OverIdentification Test

In Section 11.27 we introduced the Sargan (1958) overidentification test for the 2SLS estimator under the assumption of homoskedasticity. L. Hansen (1982) generalized the test to cover the GMM estimator allowing for general heteroskedasticity.

Recall, overidentified models ( $\ell > k$ ) are special in the sense that there may not be a parameter value  $\beta$  such that the moment condition

$$\mathbb{E}(\mathbf{g}_i(\beta)) = \mathbf{0}$$

holds. Thus the model – the overidentifying restrictions – are testable.

For example, take the linear model  $y_i = \beta_1' \mathbf{x}_{1i} + \beta_2' \mathbf{x}_{2i} + e_i$  with  $\mathbb{E}(\mathbf{x}_{1i}e_i) = \mathbf{0}$  and  $\mathbb{E}(\mathbf{x}_{2i}e_i) = \mathbf{0}$ . It is possible that  $\beta_2 = \mathbf{0}$ , so that the linear equation may be written as  $y_i = \beta_1' \mathbf{x}_{1i} + e_i$ . However, it is possible that  $\beta_2 \neq \mathbf{0}$ , and in this case it would be impossible to find a value of  $\beta_1$  so that both  $\mathbb{E}(\mathbf{x}_{1i}(y_i - \mathbf{x}_{1i}'\beta_1)) = \mathbf{0}$  and  $\mathbb{E}(\mathbf{x}_{2i}(y_i - \mathbf{x}_{1i}'\beta_1)) = \mathbf{0}$  hold simultaneously. In this sense an exclusion restriction can be seen as an overidentifying restriction.

Note that  $\bar{\mathbf{g}}_n \xrightarrow{p} \mathbb{E}(\mathbf{g}_i)$ , and thus  $\bar{\mathbf{g}}_n$  can be used to assess whether or not the hypothesis that  $\mathbb{E}(\mathbf{g}_i) = \mathbf{0}$  is true or not. Assuming that an efficient weight matrix estimate is used, the criterion function at the parameter estimates is

$$\begin{aligned} J &= J(\hat{\beta}_{\text{gmm}}) \\ &= n \bar{\mathbf{g}}_n' \hat{\Omega}^{-1} \bar{\mathbf{g}}_n \end{aligned}$$

is a quadratic form in  $\bar{\mathbf{g}}_n$ , and is thus a natural test statistic for  $\mathbb{H}_0 : \mathbb{E}(\mathbf{g}_i) = \mathbf{0}$ . Note that we assume that the criterion function is constructed with an efficient weight matrix estimate. This is important for the distribution theory.

**Theorem 12.18.1** Under Assumption 11.14.1 and  $\Omega > 0$ , then as  $n \rightarrow \infty$ ,

$$J = J(\hat{\beta}_{\text{gmm}}) \xrightarrow{d} \chi_{\ell-k}^2.$$

For  $c$  satisfying  $\alpha = 1 - G_{\ell-k}(c)$ ,

$$\Pr(J > c \mid \mathbb{H}_0) \longrightarrow \alpha$$

so the test “Reject  $\mathbb{H}_0$  if  $J > c$ ” has asymptotic size  $\alpha$ .

The proof of the theorem is left to Exercise 12.8.

The degrees of freedom of the asymptotic distribution are the number of overidentifying restrictions. If the statistic  $J$  exceeds the chi-square critical value, we can reject the model. Based on this information alone it is unclear what is wrong, but it is typically cause for concern. The GMM overidentification test is a very useful by-product of the GMM methodology, and it is advisable to report the statistic  $J$  whenever GMM is the estimation method. When over-identified models are estimated by GMM, it is customary to report the  $J$  statistic as a general test of model adequacy.

In Stata, the command `estat overid` after `ivregress gmm` can be used to implement the overidentification test. The GMM criterion  $J$  and its asymptotic p-value using the  $\chi_{\ell-k}^2$  distribution are reported.

## 12.19 Subset OverIdentification Tests

In Section 11.28 we introduced subset overidentification tests for the 2SLS estimator under the assumption of homoskedasticity. In this section we describe how to construct analogous tests for the GMM estimator under general heteroskedasticity.

Recall, subset overidentification tests are used when it is desired to focus attention on a subset of instruments whose validity is questioned. Partition  $\mathbf{z}_i = (\mathbf{z}_{ai}, \mathbf{z}_{bi})$  with dimensions  $\ell_a$  and  $\ell_b$ , respectively, where  $\mathbf{z}_{ai}$  contains the instruments which are believed to be uncorrelated with  $e_i$ , and  $\mathbf{z}_{bi}$  contains the instruments which may be correlated with  $e_i$ . It is necessary to select this partition so that  $\ell_a > k$ , so that the instruments  $\mathbf{z}_{ai}$  alone identify the parameters. The instruments  $\mathbf{z}_{bi}$  are potentially valid additional instruments.

Given this partition, the maintained hypothesis is that  $\mathbb{E}(\mathbf{z}_{ai}e_i) = \mathbf{0}$ . The null and alternative hypotheses are

$$\begin{aligned}\mathbb{H}_0 &: \mathbb{E}(\mathbf{z}_{bi}e_i) = \mathbf{0} \\ \mathbb{H}_1 &: \mathbb{E}(\mathbf{z}_{bi}e_i) \neq \mathbf{0}.\end{aligned}$$

The GMM test is constructed as follows. First, estimate the model by efficient GMM with only the smaller set  $\mathbf{z}_{ai}$  of instruments. Let  $\tilde{J}$  denote the resulting GMM criterion. Second, estimate the model by efficient GMM with the full set  $\mathbf{z}_i = (\mathbf{z}_{ai}, \mathbf{z}_{bi})$  of instruments. Let  $\hat{J}$  denote the resulting GMM criterion. The test statistic is the difference in the criterion functions:

$$C = \hat{J} - \tilde{J}.$$

This is similar in form to the GMM distance statistic presented in Section 12.16. The difference is that the distance statistic compares models which differ based on the parameter restrictions, while the  $C$  statistic compares models based on different instrument sets.

Typically, the model with the greater instrument set will produce a larger value for  $J$  so that  $C \geq 0$ . However negative values can algebraically occur. That is okay for this simply leads to a non-rejection of  $\mathbb{H}_0$ .

If the smaller instrument set  $\mathbf{z}_{ai}$  is just-identified so that  $\ell_a = k$  then  $\tilde{J} = 0$  so  $C = \hat{J}$  is simply the standard overidentification test. This is why we have restricted attention to the case  $\ell_a > k$ .

The test has the following large sample distribution.

**Theorem 12.19.1** *Under Assumption 11.14.1,  $\Omega > 0$ , and  $\mathbb{E}(\mathbf{z}_{ai}\mathbf{x}'_i)$  has full rank  $k$ , then as  $n \rightarrow \infty$ ,*

$$C \xrightarrow{d} \chi^2_{\ell_b}.$$

*For  $c$  satisfying  $\alpha = 1 - G_{\ell_b}(c)$ ,*

$$\Pr(C > c \mid \mathbb{H}_0) \longrightarrow \alpha$$

*so the test “Reject  $\mathbb{H}_0$  if  $C > c$ ” has asymptotic size  $\alpha$ .*

The proof of Theorem 12.19.1 is presented in Section 12.24.

In Stata, the command `estat overid zb` after `ivregress gmm` can be used to implement a subset overidentification test, where `zb` is the name(s) of the instrument(s) tested for validity. The statistic  $C$  and its asymptotic p-value using the  $\chi^2_{\ell_b}$  distribution are reported.

## 12.20 Endogeneity Test

In Section 11.25 we introduced tests for endogeneity in the context of 2SLS estimation. Endogeneity tests are simple to implement in the GMM framework as a subset overidentification test. The model is

$$y_i = \mathbf{x}'_{1i}\boldsymbol{\beta}_1 + \mathbf{x}'_{2i}\boldsymbol{\beta}_2 + e_i$$

where the maintained assumption is that the regressors  $\mathbf{x}_{1i}$  and excluded instruments  $\mathbf{z}_{2i}$  are exogenous so that  $\mathbb{E}(\mathbf{x}_{1i}e_i) = \mathbf{0}$  and  $\mathbb{E}(\mathbf{z}_{2i}e_i) = \mathbf{0}$ . The question is whether or not  $\mathbf{x}_{2i}$  is endogenous. Thus the null hypothesis is

$$\mathbb{H}_0 : \mathbb{E}(\mathbf{x}_{2i}e_i) = \mathbf{0}$$

with the alternative

$$\mathbb{H}_1 : \mathbb{E}(\mathbf{x}_{2i}e_i) \neq \mathbf{0}.$$

The GMM test is constructed as follows. First, estimate the model by efficient GMM using  $(\mathbf{x}_{1i}, \mathbf{z}_{2i})$  as instruments for  $(\mathbf{x}_{1i}, \mathbf{x}_{2i})$ . Let  $\tilde{J}$  denote the resulting GMM criterion. Second, estimate the model by efficient GMM using  $(\mathbf{x}_{1i}, \mathbf{x}_{2i}, \mathbf{z}_{2i})$  as instruments for  $(\mathbf{x}_{1i}, \mathbf{x}_{2i})$ . Let  $\hat{J}$  denote the resulting GMM criterion. The test statistic is the difference in the criterion functions:

$$C = \hat{J} - \tilde{J}.$$

The distribution theory for the test is a special case of the theory of overidentification testing.

**Theorem 12.20.1** *Under Assumption 11.14.1,  $\Omega > 0$ , and  $\mathbb{E}(\mathbf{z}_{2i}\mathbf{x}'_{2i})$  has full rank  $k_2$ , then as  $n \rightarrow \infty$ ,*

$$C \xrightarrow{d} \chi^2_{k_2}.$$

*For  $c$  satisfying  $\alpha = 1 - G_{\ell_2}(c)$ ,*

$$\Pr(C > c \mid \mathbb{H}_0) \longrightarrow \alpha$$

*so the test “Reject  $\mathbb{H}_0$  if  $C > c$ ” has asymptotic size  $\alpha$ .*

In Stata, the command `estat endogenous` after `ivregress gmm` can be used to implement the test for endogeneity. The statistic  $C$  and its asymptotic p-value using the  $\chi^2_{k_2}$  distribution are reported.

## 12.21 Subset Endogeneity Test

In Section 11.26 we introduced subset endogeneity tests for 2SLS estimation. GMM tests are simple to implement as subset overidentification tests. The model is

$$y_i = \mathbf{x}'_{1i}\boldsymbol{\beta}_1 + \mathbf{x}'_{2i}\boldsymbol{\beta}_2 + \mathbf{x}'_{3i}\boldsymbol{\beta}_3 + e_i$$

$$\mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0}$$

where the instrument vector is  $\mathbf{z}_i = (\mathbf{x}_{1i}, \mathbf{z}_{2i})$ . The  $k_3 \times 1$  variables  $\mathbf{x}_{3i}$  are treated as endogenous, and the  $k_2 \times 1$  variables  $\mathbf{x}_{2i}$  are treated as potentially endogenous. The hypothesis to test is that  $\mathbf{x}_{2i}$  is exogenous, or

$$\mathbb{H}_0 : \mathbb{E}(\mathbf{x}_{2i}e_i) = \mathbf{0}$$

against

$$\mathbb{H}_1 : \mathbb{E}(\mathbf{x}_{2i}e_i) \neq \mathbf{0}.$$

The test requires that  $\ell_2 \geq (k_2 + k_3)$  so that the model can be estimated under  $\mathbb{H}_1$ .

The GMM test is constructed as follows. First, estimate the model by efficient GMM using  $(\mathbf{x}_{1i}, \mathbf{z}_{2i})$  as instruments for  $(\mathbf{x}_{1i}, \mathbf{x}_{2i}, \mathbf{x}_{3i})$ . Let  $\tilde{J}$  denote the resulting GMM criterion. Second, estimate the model by efficient GMM using  $(\mathbf{x}_{1i}, \mathbf{x}_{2i}, \mathbf{z}_{2i})$  as instruments for  $(\mathbf{x}_{1i}, \mathbf{x}_{2i}, \mathbf{x}_{3i})$ . Let  $\hat{J}$  denote the resulting GMM criterion. The test statistic is the difference in the criterion functions:

$$C = \hat{J} - \tilde{J}.$$

The distribution theory for the test is a special case of the theory of overidentification testing.

**Theorem 12.21.1** *Under Assumption 11.14.1,  $\mathbf{\Omega} > \mathbf{0}$ , and  $\mathbb{E}(\mathbf{z}_{2i}(\mathbf{x}'_{2i}, \mathbf{x}'_{3i}))$  has full rank  $k_2 + k_3$ , then as  $n \rightarrow \infty$ ,*

$$C \xrightarrow{d} \chi^2_{k_2}.$$

*For  $c$  satisfying  $\alpha = 1 - G_{\ell_2}(c)$ ,*

$$\Pr(C > c \mid \mathbb{H}_0) \longrightarrow \alpha$$

*so the test “Reject  $\mathbb{H}_0$  if  $C > c$ ” has asymptotic size  $\alpha$ .*

In Stata, the command `estat endogenous x2` after `ivregress gmm` can be used to implement the test for endogeneity, where `x2` is the name(s) of the variable(s) tested for endogeneity. The statistic  $C$  and its asymptotic p-value using the  $\chi^2_{k_2}$  distribution are reported.

## 12.22 GMM: The General Case

In its most general form, GMM applies whenever an economic or statistical model implies the  $\ell \times 1$  moment condition

$$\mathbb{E}(\mathbf{g}_i(\boldsymbol{\beta})) = \mathbf{0}.$$

Often, this is all that is known. Identification requires  $\ell \geq k = \dim(\boldsymbol{\beta})$ . The GMM estimator minimizes

$$J(\boldsymbol{\beta}) = n \cdot \bar{\mathbf{g}}_n(\boldsymbol{\beta})' \widehat{\mathbf{W}} \bar{\mathbf{g}}_n(\boldsymbol{\beta})$$

for some weight matrix  $\widehat{\mathbf{W}}$ , where

$$\bar{\mathbf{g}}_n(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \mathbf{g}_i(\boldsymbol{\beta}).$$

The efficient GMM estimator can be constructed by setting

$$\widehat{\mathbf{W}} = \left( \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{g}}_i \hat{\mathbf{g}}_i' - \bar{\mathbf{g}}_n \bar{\mathbf{g}}_n' \right)^{-1},$$

with  $\hat{\mathbf{g}}_i = \mathbf{g}(\mathbf{w}_i, \tilde{\boldsymbol{\beta}})$  constructed using a preliminary consistent estimator  $\tilde{\boldsymbol{\beta}}$ , perhaps obtained by first setting  $\widehat{\mathbf{W}} = \mathbf{I}_\ell$ .

As in the case of the linear model, the weight matrix can be iterated until convergence to obtain the iterated GMM estimator.

**Proposition 12.22.1** *Distribution of Nonlinear GMM Estimator**Under general regularity conditions,*

$$\sqrt{n} \left( \hat{\beta}_{\text{gmm}} - \beta \right) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_{\beta}).$$

*where*

$$\mathbf{V}_{\beta} = (\mathbf{Q}' \mathbf{W} \mathbf{Q})^{-1} (\mathbf{Q}' \mathbf{W} \mathbf{\Omega} \mathbf{W} \mathbf{Q}) (\mathbf{Q}' \mathbf{W} \mathbf{Q})^{-1}$$

*with*

$$\mathbf{\Omega} = \mathbb{E}(\mathbf{g}_i \mathbf{g}_i')$$

*and*

$$\mathbf{Q} = \mathbb{E} \left( \frac{\partial}{\partial \beta'} \mathbf{g}_i(\beta) \right).$$

*If the efficient weight matrix is used then*

$$\mathbf{V}_{\beta} = (\mathbf{Q}' \mathbf{\Omega}^{-1} \mathbf{Q})^{-1}.$$

The proof of this result is omitted as it uses more advanced techniques.

The asymptotic covariance matrices can be estimated by sample counterparts of the population matrices. For the case of a general weight matrix,

$$\hat{\mathbf{V}}_{\beta} = \left( \hat{\mathbf{Q}}' \hat{\mathbf{W}} \hat{\mathbf{Q}} \right)^{-1} \left( \hat{\mathbf{Q}}' \hat{\mathbf{W}} \hat{\mathbf{\Omega}} \hat{\mathbf{W}} \hat{\mathbf{Q}} \right) \left( \hat{\mathbf{Q}}' \hat{\mathbf{W}} \hat{\mathbf{Q}} \right)^{-1}$$

where

$$\hat{\mathbf{\Omega}} = \frac{1}{n} \sum_{i=1}^n \left( \mathbf{g}_i(\hat{\beta}) - \bar{\mathbf{g}} \right) \left( \mathbf{g}_i(\hat{\beta}) - \bar{\mathbf{g}} \right)'$$

$$\bar{\mathbf{g}} = n^{-1} \sum_{i=1}^n \mathbf{g}_i(\hat{\beta})$$

and

$$\hat{\mathbf{Q}} = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \beta'} \mathbf{g}_i(\hat{\beta}).$$

For the case of the iterated efficient weight matrix,

$$\hat{\mathbf{V}}_{\beta} = \left( \hat{\mathbf{Q}}' \hat{\mathbf{\Omega}}^{-1} \hat{\mathbf{Q}} \right)^{-1}.$$

All of the methods discussed in this chapter – Wald tests, constrained estimation, Distance tests, overidentification tests, endogeneity tests – apply similarly to the nonlinear GMM estimator (under the same regularity conditions as the latter).

## 12.23 Conditional Moment Equation Models

In many contexts, an economic model implies more than an unconditional moment restriction of the form  $\mathbb{E}(\mathbf{g}(\mathbf{w}_i, \beta)) = \mathbf{0}$ . It implies a conditional moment restriction of the form

$$\mathbb{E}(\mathbf{e}_i(\beta) \mid \mathbf{z}_i) = \mathbf{0}$$

where  $\mathbf{e}_i(\beta)$  is some  $s \times 1$  function of the observation and the parameters. In many cases,  $s = 1$ . The variable  $\mathbf{z}_i$  is often called an **instrument**.



It turns out that this conditional moment restriction is much more powerful, and restrictive, than the unconditional moment equation model discussed throughout this chapter.

For example, the linear model  $y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$  with instruments  $\mathbf{z}_i$  falls into this class under the assumption  $\mathbb{E}(e_i | \mathbf{z}_i) = 0$ . In this case,  $e_i(\boldsymbol{\beta}) = y_i - \mathbf{x}_i' \boldsymbol{\beta}$ .

It is also helpful to realize that conventional regression models also fall into this class, except that in this case  $\mathbf{x}_i = \mathbf{z}_i$ . For example, in linear regression,  $e_i(\boldsymbol{\beta}) = y_i - \mathbf{x}_i' \boldsymbol{\beta}$ , while in a nonlinear regression model  $e_i(\boldsymbol{\beta}) = y_i - g(\mathbf{x}_i, \boldsymbol{\beta})$ . In a joint model of the conditional mean  $\mathbb{E}(y | \mathbf{x}) = \mathbf{x}' \boldsymbol{\beta}$  and variance  $\text{var}(y | \mathbf{x}) = f(\mathbf{x})' \boldsymbol{\gamma}$ , then

$$e_i(\boldsymbol{\beta}, \boldsymbol{\gamma}) = \begin{cases} y_i - \mathbf{x}_i' \boldsymbol{\beta} \\ (y_i - \mathbf{x}_i' \boldsymbol{\beta})^2 - f(\mathbf{x}_i)' \boldsymbol{\gamma} \end{cases}.$$

Here  $s = 2$ .

Given a conditional moment restriction, an unconditional moment restriction can always be constructed. That is for any  $\ell \times 1$  function  $\boldsymbol{\phi}(\mathbf{z}, \boldsymbol{\beta})$ , we can set  $\mathbf{g}_i(\boldsymbol{\beta}) = \boldsymbol{\phi}(\mathbf{z}_i, \boldsymbol{\beta}) e_i(\boldsymbol{\beta})$  which satisfies  $\mathbb{E}(\mathbf{g}_i(\boldsymbol{\beta})) = \mathbf{0}$  and hence defines an unconditional moment equation model. The obvious problem is that the class of functions  $\boldsymbol{\phi}$  is infinite. Which should be selected?

This is equivalent to the problem of selection of the best instruments. If  $z_i \in \mathbb{R}$  is a valid instrument satisfying  $\mathbb{E}(e_i | z_i) = 0$ , then  $z_i, z_i^2, z_i^3, \dots$ , etc., are all valid instruments. Which should be used?

One solution is to construct an infinite list of potent instruments, and then use the first  $k$  instruments. How is  $k$  to be determined? This is an area of theory still under development. A recent study of this problem is Donald and Newey (2001).

Another approach is to construct the **optimal instrument**. The form was uncovered by Chamberlain (1987). Take the case  $s = 1$ . Let

$$\mathbf{R}_i = \mathbb{E} \left( \frac{\partial}{\partial \boldsymbol{\beta}} e_i(\boldsymbol{\beta}) | \mathbf{z}_i \right)$$

and

$$\sigma_i^2 = \mathbb{E} (e_i(\boldsymbol{\beta})^2 | \mathbf{z}_i).$$

Then the “optimal instrument” is

$$\mathbf{A}_i = -\sigma_i^{-2} \mathbf{R}_i$$

so the optimal moment is

$$\mathbf{g}_i(\boldsymbol{\beta}) = \mathbf{A}_i e_i(\boldsymbol{\beta}).$$

Setting  $\mathbf{g}_i(\boldsymbol{\beta})$  to be this choice (which is  $k \times 1$ , so is just-identified) yields the best GMM estimator possible.

In practice,  $\mathbf{A}_i$  is unknown, but its form does help us think about construction of optimal instruments.

In the linear model  $e_i(\boldsymbol{\beta}) = y_i - \mathbf{x}_i' \boldsymbol{\beta}$ , note that

$$\mathbf{R}_i = -\mathbb{E}(\mathbf{x}_i | \mathbf{z}_i)$$

and

$$\sigma_i^2 = \mathbb{E}(e_i^2 | \mathbf{z}_i),$$

so

$$\mathbf{A}_i = \sigma_i^{-2} \mathbb{E}(\mathbf{x}_i | \mathbf{z}_i).$$

In the case of linear regression,  $\mathbf{x}_i = \mathbf{z}_i$ , so  $\mathbf{A}_i = \sigma_i^{-2} \mathbf{z}_i$ . Hence efficient GMM is equivalently to optimal GLS.

In the case of endogenous variables, note that the efficient instrument  $\mathbf{A}_i$  involves the estimation of the conditional mean of  $\mathbf{x}_i$  given  $\mathbf{z}_i$ . In other words, to get the best instrument for  $\mathbf{x}_i$ , we need the best conditional mean model for  $\mathbf{x}_i$  given  $\mathbf{z}_i$ , not just an arbitrary linear projection. The efficient instrument is also inversely proportional to the conditional variance of  $e_i$ . This is the same as the GLS estimator; namely that improved efficiency can be obtained if the observations are weighted inversely to the conditional variance of the errors.

## 12.24 Technical Proofs\*

**Proof of Theorem 12.16.1.** Set

$$\begin{aligned}\tilde{\mathbf{e}} &= \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{cgmm}} \\ \hat{\mathbf{e}} &= \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{gmm}}\end{aligned}$$

By standard covariance matrix analysis  $\hat{\boldsymbol{\Omega}} \xrightarrow{p} \boldsymbol{\Omega}$  and  $\tilde{\boldsymbol{\Omega}} \xrightarrow{p} \boldsymbol{\Omega}$ . Thus we can replace  $\hat{\boldsymbol{\Omega}}$  and  $\tilde{\boldsymbol{\Omega}}$  in the criteria without affecting the asymptotic distribution. With this substitution  $\hat{J}(\boldsymbol{\beta}) = \tilde{J}(\boldsymbol{\beta}) = n \cdot \bar{\mathbf{g}}_n(\boldsymbol{\beta})' \boldsymbol{\Omega}^{-1} \bar{\mathbf{g}}_n(\boldsymbol{\beta})$ . From (12.18) and setting  $\mathbf{W} = \boldsymbol{\Omega}^{-1}$

$$\sqrt{n} \left( \hat{\boldsymbol{\beta}}_{\text{cgmm}} - \boldsymbol{\beta} \right) = \left( \mathbf{I}_k - \mathbf{V}\mathbf{R}(\mathbf{R}'\mathbf{V}\mathbf{R})^{-1}\mathbf{R}' \right) \sqrt{n} \left( \hat{\boldsymbol{\beta}}_{\text{gmm}} - \boldsymbol{\beta} \right) + o_p(1).$$

Thus

$$\begin{aligned}\sqrt{n}\bar{\mathbf{g}}_n(\hat{\boldsymbol{\beta}}_{\text{cgmm}}) &= \frac{1}{\sqrt{n}}\mathbf{Z}'\tilde{\mathbf{e}} \\ &= \frac{1}{\sqrt{n}}\mathbf{Z}'\hat{\mathbf{e}} + \frac{1}{n}\mathbf{Z}'\mathbf{X}\mathbf{V}\mathbf{R}(\mathbf{R}'\mathbf{V}\mathbf{R})^{-1}\mathbf{R}'\sqrt{n} \left( \hat{\boldsymbol{\beta}}_{\text{gmm}} - \boldsymbol{\beta} \right) + o_p(1).\end{aligned}$$

The first-order condition for  $\hat{\boldsymbol{\beta}}_{\text{gmm}}$  is  $\mathbf{X}'\mathbf{Z}\boldsymbol{\Omega}^{-1}\mathbf{Z}'\hat{\mathbf{e}} = 0$  so the two components in this last expression are orthogonal with respect to the weight matrix  $\boldsymbol{\Omega}^{-1}$ . Hence

$$\begin{aligned}\hat{J}(\hat{\boldsymbol{\beta}}_{\text{cgmm}}) &= \left( \frac{1}{\sqrt{n}}\mathbf{Z}'\tilde{\mathbf{e}} \right)' \boldsymbol{\Omega}^{-1} \left( \frac{1}{\sqrt{n}}\mathbf{Z}'\tilde{\mathbf{e}} \right) \\ &= \left( \frac{1}{\sqrt{n}}\mathbf{Z}'\hat{\mathbf{e}} \right)' \boldsymbol{\Omega}^{-1} \left( \frac{1}{\sqrt{n}}\mathbf{Z}'\hat{\mathbf{e}} \right) \\ &\quad + n \left( \hat{\boldsymbol{\beta}}_{\text{gmm}} - \boldsymbol{\beta} \right)' \mathbf{R}(\mathbf{R}'\mathbf{V}\mathbf{R})^{-1} \mathbf{R}' \mathbf{V} \frac{1}{n} \mathbf{X}'\mathbf{Z}\boldsymbol{\Omega}^{-1} \frac{1}{n} \mathbf{Z}'\mathbf{X}\mathbf{V}\mathbf{R}(\mathbf{R}'\mathbf{V}\mathbf{R})^{-1} \mathbf{R}' n \left( \hat{\boldsymbol{\beta}}_{\text{gmm}} - \boldsymbol{\beta} \right) \\ &\quad + o_p(1) \\ &= \hat{J}(\hat{\boldsymbol{\beta}}_{\text{gmm}}) + n \left( \hat{\boldsymbol{\beta}}_{\text{gmm}} - \boldsymbol{\beta} \right)' \mathbf{R}(\mathbf{R}'\mathbf{V}\mathbf{R})^{-1} \mathbf{R}' \left( \hat{\boldsymbol{\beta}}_{\text{gmm}} - \boldsymbol{\beta} \right) + o_p(1).\end{aligned}$$

Thus

$$\begin{aligned}D &= \hat{J}(\hat{\boldsymbol{\beta}}_{\text{cgmm}}) - \hat{J}(\hat{\boldsymbol{\beta}}_{\text{gmm}}) \\ &= n \left( \hat{\boldsymbol{\beta}}_{\text{gmm}} - \boldsymbol{\beta} \right)' \mathbf{R}(\mathbf{R}'\mathbf{V}\mathbf{R})^{-1} \mathbf{R}' \left( \hat{\boldsymbol{\beta}}_{\text{gmm}} - \boldsymbol{\beta} \right) + o_p(1)\end{aligned}$$

which converges in distribution to  $\chi_q^2$  as claimed.  $\blacksquare$

**Proof of Theorem 12.19.1.** Let  $\tilde{\boldsymbol{\beta}}$  denote the GMM estimate obtained with the instrument set

$\mathbf{z}_{ai}$  and let  $\hat{\boldsymbol{\beta}}$  denote the GMM estimates obtained with the instrument set  $\mathbf{z}_i$ . Set

$$\begin{aligned}\tilde{\mathbf{e}} &= \mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}} \\ \hat{\mathbf{e}} &= \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} \\ \tilde{\boldsymbol{\Omega}} &= n^{-1} \sum_{i=1}^n \mathbf{z}_{ai} \mathbf{z}_{ai}' \tilde{e}_i^2 \\ \hat{\boldsymbol{\Omega}} &= n^{-1} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i' \hat{e}_i^2\end{aligned}$$

Let  $\mathbf{R}$  be the  $\ell \times \ell_a$  selector matrix so that  $\mathbf{z}_{ai} = \mathbf{R}' \mathbf{z}_i$ . Note that

$$\tilde{\boldsymbol{\Omega}} = \mathbf{R}' n^{-1} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i' \tilde{e}_i^2 \mathbf{R}.$$

By standard covariance matrix analysis,  $\hat{\boldsymbol{\Omega}} \xrightarrow{p} \boldsymbol{\Omega}$  and  $\tilde{\boldsymbol{\Omega}} \xrightarrow{p} \mathbf{R}' \boldsymbol{\Omega} \mathbf{R}$ . Also,  $\frac{1}{n} \mathbf{Z}' \mathbf{X} \xrightarrow{p} \mathbf{Q}$ , say. By the CLT,  $n^{-1/2} \mathbf{Z}' \mathbf{e} \xrightarrow{d} \mathbf{Z}$  where  $\mathbf{Z} \sim \mathbf{N}(\mathbf{0}, \boldsymbol{\Omega})$ . Then

$$\begin{aligned}n^{-1/2} \mathbf{Z}' \hat{\mathbf{e}} &= \left( \mathbf{I}_\ell - \left( \frac{1}{n} \mathbf{Z}' \mathbf{X} \right) \left( \frac{1}{n} \mathbf{X}' \mathbf{Z} \tilde{\boldsymbol{\Omega}}^{-1} \frac{1}{n} \mathbf{Z}' \mathbf{X} \right)^{-1} \left( \frac{1}{n} \mathbf{X}' \mathbf{Z} \right) \hat{\boldsymbol{\Omega}}^{-1} \right) n^{-1/2} \mathbf{Z}' \mathbf{e} \\ &\xrightarrow{d} \left( \mathbf{I}_\ell - \mathbf{Q} (\mathbf{Q}' \boldsymbol{\Omega}^{-1} \mathbf{Q})^{-1} \mathbf{Q}' \boldsymbol{\Omega}^{-1} \right) \mathbf{Z}\end{aligned}$$

and

$$\begin{aligned}n^{-1/2} \mathbf{Z}'_a \tilde{\mathbf{e}} &= \mathbf{R}' \left( \mathbf{I}_\ell - \left( \frac{1}{n} \mathbf{Z}' \mathbf{X} \right) \left( \frac{1}{n} \mathbf{X}' \mathbf{Z} \tilde{\boldsymbol{\Omega}}^{-1} \frac{1}{n} \mathbf{Z}' \mathbf{X} \right)^{-1} \left( \frac{1}{n} \mathbf{X}' \mathbf{Z} \right) \tilde{\boldsymbol{\Omega}}^{-1} \mathbf{R}' \right) n^{-1/2} \mathbf{Z}' \mathbf{e} \\ &\xrightarrow{d} \mathbf{R}' \left( \mathbf{I}_\ell - \mathbf{Q} (\mathbf{Q}' \mathbf{R} (\mathbf{R}' \boldsymbol{\Omega} \mathbf{R})^{-1} \mathbf{R}' \mathbf{Q})^{-1} \mathbf{Q}' \mathbf{R} (\mathbf{R}' \boldsymbol{\Omega} \mathbf{R})^{-1} \mathbf{R}' \right) \mathbf{Z}\end{aligned}$$

jointly. Thus

$$\hat{\mathbf{J}} \xrightarrow{d} \mathbf{Z}' \left( \boldsymbol{\Omega}^{-1} - \boldsymbol{\Omega}^{-1} \mathbf{Q} (\mathbf{Q}' \boldsymbol{\Omega}^{-1} \mathbf{Q})^{-1} \mathbf{Q}' \boldsymbol{\Omega}^{-1} \right) \mathbf{Z}$$

and

$$\tilde{\mathbf{J}} \xrightarrow{d} \mathbf{Z}' \left( \mathbf{R} (\mathbf{R}' \boldsymbol{\Omega} \mathbf{R})^{-1} \mathbf{R}' - \mathbf{R} (\mathbf{R}' \boldsymbol{\Omega} \mathbf{R})^{-1} \mathbf{R}' \mathbf{Q} (\mathbf{Q}' \mathbf{R} (\mathbf{R}' \boldsymbol{\Omega} \mathbf{R})^{-1} \mathbf{R}' \mathbf{Q})^{-1} \mathbf{Q}' \mathbf{R} (\mathbf{R}' \boldsymbol{\Omega} \mathbf{R})^{-1} \mathbf{R}' \right) \mathbf{Z}.$$

By linear rotations of  $\mathbf{Z}$  and  $\mathbf{R}$  we can set  $\boldsymbol{\Omega} = \mathbf{I}_\ell$  to simplify the notation. It follows that

$$\mathbf{C} \xrightarrow{d} \mathbf{Z}' \mathbf{A} \mathbf{Z}$$

where

$$\mathbf{A} = \left( \mathbf{I}_\ell - \mathbf{P}_Q - \mathbf{P}_R + \mathbf{P}_R \mathbf{Q} (\mathbf{Q}' \mathbf{P}_R \mathbf{Q})^{-1} \mathbf{Q}' \mathbf{P}_R \right),$$

$\mathbf{P}_R = \mathbf{R} (\mathbf{R}' \mathbf{R})^{-1} \mathbf{R}'$ ,  $\mathbf{P}_Q = \mathbf{Q} (\mathbf{Q}' \mathbf{Q})^{-1} \mathbf{Q}'$ , and  $\mathbf{Z} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_\ell)$ . This is a quadratic form in a standard normal vector, and the matrix  $\mathbf{A}$  is idempotent (this is straightforward to check). It is thus distributed as  $\chi_d^2$  with degrees of freedom  $d$  equal to the rank of  $\mathbf{A}$ . This is

$$\begin{aligned}\text{rank}(\mathbf{A}) &= \text{tr} \left( \mathbf{I}_\ell - \mathbf{P}_Q - \mathbf{P}_R + \mathbf{P}_R \mathbf{Q} (\mathbf{Q}' \mathbf{P}_R \mathbf{Q})^{-1} \mathbf{Q}' \mathbf{P}_R \right) \\ &= \ell - k - \ell_a + k \\ &= \ell_b.\end{aligned}$$

Thus the asymptotic distribution of  $\mathbf{C}$  is  $\chi_{\ell_b}^2$  as claimed.  $\blacksquare$

## Exercises

**Exercise 12.1** Take the model

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= \mathbf{0} \\ e_i^2 &= \mathbf{z}_i' \boldsymbol{\gamma} + \eta_i \\ \mathbb{E}(\mathbf{z}_i \eta_i) &= \mathbf{0}. \end{aligned}$$

Find the method of moments estimators  $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\gamma}})$  for  $(\boldsymbol{\beta}, \boldsymbol{\gamma})$ .

**Exercise 12.2** Take the single equation

$$\begin{aligned} \mathbf{y} &= \mathbf{X} \boldsymbol{\beta} + \mathbf{e} \\ \mathbb{E}(\mathbf{e} \mid \mathbf{Z}) &= \mathbf{0} \end{aligned}$$

Assume  $\mathbb{E}(e_i^2 \mid \mathbf{z}_i) = \sigma^2$ . Show that if  $\hat{\boldsymbol{\beta}}_{\text{gmm}}$  is the GMM estimated by GMM with weight matrix  $\mathbf{W}_n = (\mathbf{Z}' \mathbf{Z})^{-1}$ , then

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sigma^2 (\mathbf{Q}' \mathbf{M}^{-1} \mathbf{Q})^{-1})$$

where  $\mathbf{Q} = \mathbb{E}(\mathbf{z}_i \mathbf{x}_i')$  and  $\mathbf{M} = \mathbb{E}(\mathbf{z}_i \mathbf{z}_i')$ .

**Exercise 12.3** Take the model  $y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$  with  $\mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0}$ . Let  $\tilde{e}_i = y_i - \mathbf{x}_i' \tilde{\boldsymbol{\beta}}$  where  $\tilde{\boldsymbol{\beta}}$  is consistent for  $\boldsymbol{\beta}$  (e.g. a GMM estimator with arbitrary weight matrix). Define an estimate of the optimal GMM weight matrix

$$\widehat{\mathbf{W}} = \left( \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i' \tilde{e}_i^2 \right)^{-1}.$$

Show that  $\widehat{\mathbf{W}} \xrightarrow{p} \boldsymbol{\Omega}^{-1}$  where  $\boldsymbol{\Omega} = \mathbb{E}(\mathbf{z}_i \mathbf{z}_i' e_i^2)$ .

**Exercise 12.4** In the linear model estimated by GMM with general weight matrix  $\mathbf{W}$ , the asymptotic variance of  $\hat{\boldsymbol{\beta}}_{\text{GMM}}$  is

$$\mathbf{V} = (\mathbf{Q}' \mathbf{W} \mathbf{Q})^{-1} \mathbf{Q}' \mathbf{W} \boldsymbol{\Omega} \mathbf{W} \mathbf{Q} (\mathbf{Q}' \mathbf{W} \mathbf{Q})^{-1}$$

- Let  $\mathbf{V}_0$  be this matrix when  $\mathbf{W} = \boldsymbol{\Omega}^{-1}$ . Show that  $\mathbf{V}_0 = (\mathbf{Q}' \boldsymbol{\Omega}^{-1} \mathbf{Q})^{-1}$ .
- We want to show that for any  $\mathbf{W}$ ,  $\mathbf{V} - \mathbf{V}_0$  is positive semi-definite (for then  $\mathbf{V}_0$  is the smaller possible covariance matrix and  $\mathbf{W} = \boldsymbol{\Omega}^{-1}$  is the efficient weight matrix). To do this, start by finding matrices  $\mathbf{A}$  and  $\mathbf{B}$  such that  $\mathbf{V} = \mathbf{A}' \boldsymbol{\Omega} \mathbf{A}$  and  $\mathbf{V}_0 = \mathbf{B}' \boldsymbol{\Omega} \mathbf{B}$ .
- Show that  $\mathbf{B}' \boldsymbol{\Omega} \mathbf{A} = \mathbf{B}' \boldsymbol{\Omega} \mathbf{B}$  and therefore that  $\mathbf{B}' \boldsymbol{\Omega} (\mathbf{A} - \mathbf{B}) = \mathbf{0}$ .
- Use the expressions  $\mathbf{V} = \mathbf{A}' \boldsymbol{\Omega} \mathbf{A}$ ,  $\mathbf{A} = \mathbf{B} + (\mathbf{A} - \mathbf{B})$ , and  $\mathbf{B}' \boldsymbol{\Omega} (\mathbf{A} - \mathbf{B}) = \mathbf{0}$  to show that  $\mathbf{V} \geq \mathbf{V}_0$ .

**Exercise 12.5** The equation of interest is

$$\begin{aligned} y_i &= \mathbf{m}(\mathbf{x}_i, \boldsymbol{\beta}) + e_i \\ \mathbb{E}(\mathbf{z}_i e_i) &= \mathbf{0}. \end{aligned}$$

The observed data is  $(y_i, \mathbf{z}_i, \mathbf{x}_i)$ .  $\mathbf{z}_i$  is  $\ell \times 1$  and  $\boldsymbol{\beta}$  is  $k \times 1$ ,  $\ell \geq k$ . Show how to construct an efficient GMM estimator for  $\boldsymbol{\beta}$ .

**Exercise 12.6** As a continuation of Exercise 11.7, derive the efficient GMM estimator using the instrument  $\mathbf{z}_i = (x_i \quad x_i^2)'$ . Does this differ from 2SLS and/or OLS?

**Exercise 12.7** In the linear model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$  with  $\mathbb{E}(\mathbf{x}_i e_i) = \mathbf{0}$ , a Generalized Method of Moments (GMM) criterion function for  $\boldsymbol{\beta}$  is defined as

$$J(\boldsymbol{\beta}) = \frac{1}{n} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' \mathbf{X} \hat{\boldsymbol{\Omega}}^{-1} \mathbf{X}' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \quad (12.19)$$

where  $\hat{\boldsymbol{\Omega}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \hat{e}_i^2$ ,  $\hat{e}_i = y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}$  are the OLS residuals, and  $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$  is LS. The GMM estimator of  $\boldsymbol{\beta}$ , subject to the restriction  $\mathbf{r}(\boldsymbol{\beta}) = \mathbf{0}$ , is defined as

$$\tilde{\boldsymbol{\beta}} = \underset{\mathbf{r}(\boldsymbol{\beta})=\mathbf{0}}{\operatorname{argmin}} J_n(\boldsymbol{\beta}).$$

The GMM test statistic (the distance statistic) of the hypothesis  $\mathbf{r}(\boldsymbol{\beta}) = \mathbf{0}$  is

$$D = J(\tilde{\boldsymbol{\beta}}) = \min_{\mathbf{r}(\boldsymbol{\beta})=\mathbf{0}} J(\boldsymbol{\beta}). \quad (12.20)$$

(a) Show that you can rewrite  $J(\boldsymbol{\beta})$  in (12.19) as

$$J(\boldsymbol{\beta}) = n (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})' \hat{\mathbf{V}}_{\boldsymbol{\beta}}^{-1} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})$$

thus  $\tilde{\boldsymbol{\beta}}$  is the same as the minimum distance estimator.

(b) Show that under linear hypotheses the distance statistic  $D$  in (12.20) equals the Wald statistic.

**Exercise 12.8** Take the linear model

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(\mathbf{z}_i e_i) &= \mathbf{0}. \end{aligned}$$

and consider the GMM estimator  $\hat{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta}$ . Let

$$J = n \bar{\mathbf{g}}_n(\hat{\boldsymbol{\beta}})' \hat{\boldsymbol{\Omega}}^{-1} \bar{\mathbf{g}}_n(\hat{\boldsymbol{\beta}})$$

denote the test of overidentifying restrictions. Show that  $J \xrightarrow{d} \chi_{\ell-k}^2$  as  $n \rightarrow \infty$  by demonstrating each of the following:

(a) Since  $\boldsymbol{\Omega} > 0$ , we can write  $\boldsymbol{\Omega}^{-1} = \mathbf{C}\mathbf{C}'$  and  $\boldsymbol{\Omega} = \mathbf{C}'^{-1}\mathbf{C}^{-1}$

(b)  $J = n \left( \mathbf{C}' \bar{\mathbf{g}}_n(\hat{\boldsymbol{\beta}}) \right)' \left( \mathbf{C}' \hat{\boldsymbol{\Omega}} \mathbf{C} \right)^{-1} \mathbf{C}' \bar{\mathbf{g}}_n(\hat{\boldsymbol{\beta}})$

(c)  $\mathbf{C}' \bar{\mathbf{g}}_n(\hat{\boldsymbol{\beta}}) = \mathbf{D}_n \mathbf{C}' \bar{\mathbf{g}}_n(\boldsymbol{\beta})$  where

$$\mathbf{D}_n = \mathbf{I}_{\ell} - \mathbf{C}' \left( \frac{1}{n} \mathbf{Z}' \mathbf{X} \right) \left( \left( \frac{1}{n} \mathbf{X}' \mathbf{Z} \right) \hat{\boldsymbol{\Omega}}^{-1} \left( \frac{1}{n} \mathbf{Z}' \mathbf{X} \right) \right)^{-1} \left( \frac{1}{n} \mathbf{X}' \mathbf{Z} \right) \hat{\boldsymbol{\Omega}}^{-1} \mathbf{C}'^{-1}$$

$$\bar{\mathbf{g}}_n(\boldsymbol{\beta}) = \frac{1}{n} \mathbf{Z}' \mathbf{e}.$$

(d)  $\mathbf{D}_n \xrightarrow{p} \mathbf{I}_{\ell} - \mathbf{R}(\mathbf{R}'\mathbf{R})^{-1} \mathbf{R}'$  where  $\mathbf{R} = \mathbf{C}' \mathbb{E}(\mathbf{z}_i \mathbf{x}_i')$

(e)  $n^{1/2} \mathbf{C}' \bar{\mathbf{g}}_n(\boldsymbol{\beta}) \xrightarrow{d} \mathbf{u} \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_{\ell})$

$$(f) \quad J \xrightarrow{d} \mathbf{u}' \left( \mathbf{I}_\ell - \mathbf{R}(\mathbf{R}'\mathbf{R})^{-1} \mathbf{R}' \right) \mathbf{u}$$

$$(g) \quad \mathbf{u}' \left( \mathbf{I}_\ell - \mathbf{R}(\mathbf{R}'\mathbf{R})^{-1} \mathbf{R}' \right) \mathbf{u} \sim \chi_{\ell-k}^2.$$

Hint:  $\mathbf{I}_\ell - \mathbf{R}(\mathbf{R}'\mathbf{R})^{-1} \mathbf{R}'$  is a projection matrix.

**Exercise 12.9** Take the model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0}$$

$y_i$  scalar,  $\mathbf{x}_i$  a  $k$  vector and  $\mathbf{z}_i$  an  $\ell$  vector,  $\ell \geq k$ . Assume iid observations. Consider the statistic

$$J_n(\boldsymbol{\beta}) = n \overline{\mathbf{m}}_n(\boldsymbol{\beta})' \mathbf{W} \overline{\mathbf{m}}_n(\boldsymbol{\beta}) \\ \overline{\mathbf{m}}_n(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i (y_i - \mathbf{x}_i' \boldsymbol{\beta})$$

for some weight matrix  $\mathbf{W} > 0$ .

(a) Take the hypothesis

$$\mathbb{H}_0 : \boldsymbol{\beta} = \boldsymbol{\beta}_0$$

Derive the asymptotic distribution of  $J_n(\boldsymbol{\beta}_0)$  under  $\mathbb{H}_0$  as  $n \rightarrow \infty$ .

(b) What choice for  $\mathbf{W}$  yields a known asymptotic distribution in part (a)? (Be specific about degrees of freedom.)

(c) Write down an appropriate estimator  $\widehat{\mathbf{W}}$  for  $\mathbf{W}$  which takes advantage of  $\mathbb{H}_0$ . (You do not need to demonstrate consistency or unbiasedness.)

(d) Describe an asymptotic test of  $\mathbb{H}_0$  against  $\mathbb{H}_1 : \boldsymbol{\beta} \neq \boldsymbol{\beta}_0$  based on this statistic.

(e) Use the result in part (d) to construct a confidence region for  $\boldsymbol{\beta}$ . What can you say about the form of this region? For example, does the confidence region take the form of an ellipse, similar to conventional confidence regions?

**Exercise 12.10** Consider the model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(\mathbf{z}_i e_i) = \mathbf{0} \tag{12.21}$$

$$\mathbf{R}' \boldsymbol{\beta} = \mathbf{0} \tag{12.22}$$

with  $y_i$  scalar,  $\mathbf{x}_i$  a  $k$  vector and  $\mathbf{z}_i$  an  $\ell$  vector with  $\ell > k$ . The matrix  $\mathbf{R}$  is  $k \times \ell$  with  $1 \leq \ell < k$ . You have a random sample  $(y_i, \mathbf{x}_i, \mathbf{z}_i : i = 1, \dots, n)$ .

For simplicity, assume the “efficient” weight matrix  $\mathbf{W} = (\mathbb{E}(\mathbf{z}_i \mathbf{z}_i' e_i^2))^{-1}$  is known.

(a) Write out the GMM estimator  $\widehat{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta}$  given the moment conditions (12.21) but ignoring constraint (12.22).

(b) Write out the GMM estimator  $\widetilde{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta}$  given the moment conditions (12.21) and constraint (12.22).

(c) Find the asymptotic distribution of  $\sqrt{n}(\widetilde{\boldsymbol{\beta}} - \boldsymbol{\beta})$  as  $n \rightarrow \infty$  under the assumption that (12.21) and (12.22) are correct.

**Exercise 12.11** The observed data is  $\{y_i, x_i, z_i\} \in \mathbb{R} \times \mathbb{R}^k \times \mathbb{R}^\ell$ ,  $k > 1$  and  $\ell > k > 1$ ,  $i = 1, \dots, n$ . The model is

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(\mathbf{z}_i e_i) &= 0 \end{aligned} \tag{12.23}$$

(a) Given a weight matrix  $\mathbf{W} > 0$ , write down the GMM estimator  $\hat{\boldsymbol{\beta}}$  for  $\boldsymbol{\beta}$ .

(b) Suppose the model is misspecified in that

$$\begin{aligned} e_i &= \delta n^{-1/2} + u_i \\ \mathbb{E}(u_i \mid \mathbf{z}_i) &= 0 \end{aligned} \tag{12.24}$$

with  $\boldsymbol{\mu}_z = \mathbb{E}(\mathbf{z}_i) \neq 0$  and  $\delta \neq 0$ . Show that (12.24) implies (12.23) is false.

(c) Express  $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$  as a function of  $\mathbf{W}$ ,  $n$ ,  $\delta$ , and the variables  $(\mathbf{x}_i, \mathbf{z}_i, u_i)$ .

(d) Find the asymptotic distribution of  $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$  under Assumption (12.24).

**Exercise 12.12** The model is

$$\begin{aligned} y_i &= z_i \beta + x_i \gamma + e_i \\ \mathbb{E}(e_i \mid x_i) &= 0 \end{aligned}$$

Thus  $z_i$  is potentially endogenous and  $x_i$  is exogenous. Assume that  $z_i$  and  $x_i$  are scalar. Someone suggests estimating  $(\beta, \gamma)$  by GMM, using the pair  $(x_i, x_i^2)$  as the instruments. Is this feasible? Under what conditions, if any, (in addition to those described above) is this a valid estimator?

**Exercise 12.13** The observations are iid,  $(y_i, \mathbf{x}_i, \mathbf{q}_i : i = 1, \dots, n)$ , where  $\mathbf{x}_i$  is  $k \times 1$  and  $\mathbf{q}_i$  is  $m \times 1$ . The model is

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= 0 \\ \mathbb{E}(\mathbf{q}_i e_i) &= 0 \end{aligned}$$

Find the efficient GMM estimator for  $\boldsymbol{\beta}$ .

**Exercise 12.14** You want to estimate  $\mu = \mathbb{E}(y_i)$  under the assumption that  $\mathbb{E}(x_i) = 0$ , where  $y_i$  and  $x_i$  are scalar and observed from a random sample. Find an efficient GMM estimator for  $\mu$ .

**Exercise 12.15** Consider the model

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(\mathbf{z}_i e_i) &= 0 \\ \mathbf{R}' \boldsymbol{\beta} &= 0 \end{aligned}$$

The dimensions are  $\mathbf{x} \in \mathbb{R}^k$ ,  $\mathbf{z} \in \mathbb{R}^\ell$ ,  $\ell > k$ . The matrix  $\mathbf{R}$  is  $k \times q$ ,  $1 \leq q < k$ . Derive an efficient GMM estimator for  $\boldsymbol{\beta}$  for this model.

**Exercise 12.16** Take the linear equation  $y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$ , and consider the following estimators of  $\boldsymbol{\beta}$ .

1.  $\hat{\boldsymbol{\beta}}$ : 2SLS using the instruments  $\mathbf{z}_{1i}$
2.  $\tilde{\boldsymbol{\beta}}$ : 2SLS using the instruments  $\mathbf{z}_{1i}$

3.  $\bar{\beta}$ : GMM using the instruments  $\mathbf{z}_i = (\mathbf{z}_{1i}, \mathbf{z}_{2i})$  and the weight matrix

$$\mathbf{W} = \begin{pmatrix} (\mathbf{Z}'_1 \mathbf{Z}_1)^{-1} \lambda & 0 \\ 0 & (\mathbf{Z}'_2 \mathbf{Z}_2)^{-1} (1 - \lambda) \end{pmatrix}$$

for  $\lambda \in (0, 1)$ . Find an expression for  $\bar{\beta}$  which shows that it is a specific weighted average of  $\hat{\beta}$  and  $\hat{\beta}$ .

**Exercise 12.17** Consider the just-identified model

$$\begin{aligned} y_i &= \mathbf{x}'_{1i} \beta_1 + \mathbf{x}'_{2i} \beta_2 + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= 0 \end{aligned}$$

where  $\mathbf{x}_i = (\mathbf{x}'_{1i} \ \mathbf{x}'_{2i})'$  and  $\mathbf{z}_i$  are  $k \times 1$ . We want to test  $\mathbb{H}_0 : \beta_1 = 0$ . Three econometricians are called to advise on how to test  $\mathbb{H}_0$ .

- Econometrician 1 proposes testing  $\mathbb{H}_0$  by a Wald statistic.
- Econometrician 2 suggests testing  $\mathbb{H}_0$  by the GMM Distance Statistic.
- Econometrician 3 suggests testing  $\mathbb{H}_0$  using the test of overidentifying restrictions.

You are asked to settle this dispute. Explain the advantages and/or disadvantages of the different procedures, in this specific context.

**Exercise 12.18** Take the model

$$\begin{aligned} y_i &= \mathbf{x}'_i \beta + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= \mathbf{0} \\ \beta &= \mathbf{Q} \theta \end{aligned}$$

where  $\beta$  is  $k \times 1$ ,  $\mathbf{Q}$  is  $k \times m$  with  $m < k$ , and  $\mathbf{Q}$  is known. Assume that the observations  $(y_i, \mathbf{x}_i)$  are i.i.d. across  $i = 1, \dots, n$ .

Under these assumptions, what is the efficient estimator of  $\theta$ ?

**Exercise 12.19** Take the model

$$\begin{aligned} y_i &= \theta + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= 0 \end{aligned}$$

with  $(y_i, \mathbf{x}_i)$  a random sample.  $y_i$  is real-valued and  $\mathbf{x}_i$  is  $k \times 1$ ,  $k > 1$ .

- Find the efficient GMM estimator of  $\theta$ .
- Is this model over-identified or just-identified?
- Find the GMM test statistic for over-identification.

**Exercise 12.20** Continuation of Exercise 11.23, based on the empirical work reported in Acemoglu, Johnson and Robinson (2001)

- Re-estimate the model estimated part (j) by efficient GMM. I suggest that you use the 2SLS estimates as the first-step to get the weight matrix, and then calculate the GMM estimator from this weight matrix without further iteration. Report the estimates and standard errors.



- (b) Calculate and report the  $J$  statistic for overidentification.
- (c) Compare the GMM and 2SLS estimates. Discuss your findings

**Exercise 12.21** Continuation of Exercise 11.24, which involved estimation of a wage equation by 2SLS.

- (a) Re-estimate the model in part (a) by efficient GMM. Do the results change meaningfully?
- (b) Re-estimate the model in part (d) by efficient GMM. Do the results change meaningfully?
- (c) Report the  $J$  statistic for overidentification.

# Chapter 13

## The Bootstrap

### 13.1 Definition of the Bootstrap

Let  $F$  denote the distribution function for the population of observations  $(y_i, \mathbf{x}_i)$ . Let

$$T = T((y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n), F)$$

be a statistic of interest, for example an estimator  $\hat{\theta}$  or a t-statistic  $(\hat{\theta} - \theta) / s(\hat{\theta})$ . Note that we write  $T$  as possibly a function of  $F$ . For example, the t-statistic is a function of the parameter  $\theta = \theta(F)$  which itself is a function of  $F$ .

The exact CDF of  $T$  when the data are sampled from the distribution  $F$  is

$$G_n(u, F) = \Pr(T \leq u \mid F)$$

In general,  $G_n(u, F)$  depends on  $F$  and  $n$ , meaning that  $G$  changes as  $F$  or  $n$  changes.

Ideally, inference would be based on  $G_n(u, F)$ . This is generally impossible since  $F$  is unknown.

Asymptotic inference is based on approximating  $G_n(u, F)$  with  $G(u, F) = \lim_{n \rightarrow \infty} G_n(u, F)$ . When  $G(u, F) = G(u)$  does not depend on  $F$ , we say that  $T$  is asymptotically pivotal and use the distribution function  $G(u)$  for inferential purposes.

In a seminal contribution, Efron (1979) proposed the bootstrap, which makes a different approximation. The unknown  $F$  is replaced by a consistent estimate  $\hat{F}$  (one choice is discussed in the next section). Plugged into  $G_n(u, F)$  we obtain

$$G_n^*(u) = G_n(u, \hat{F}). \quad (13.1)$$

We call  $G_n^*$  the bootstrap distribution. Bootstrap inference is based on  $G_n^*(u)$ .

Let  $(y_i^*, \mathbf{x}_i^*)$  denote random variables from the distribution  $\hat{F}$ . A random sample  $\{(y_i^*, \mathbf{x}_i^*) : i = 1, \dots, n\}$  from this distribution is called the **bootstrap data**. The statistic  $T^* = T((y_1^*, \mathbf{x}_1^*), \dots, (y_n^*, \mathbf{x}_n^*), \hat{F})$  constructed on this sample is a random variable with distribution  $G_n^*$ . That is,  $\Pr(T^* \leq u) = G_n^*(u)$ . We call  $T^*$  the **bootstrap statistic**. The distribution of  $T^*$  is identical to that of  $T$  when the true CDF is  $\hat{F}$  rather than  $F$ .

The bootstrap distribution is itself random, as it depends on the sample through the estimator  $\hat{F}$ .

In the next sections we describe computation of the bootstrap distribution.

### 13.2 The Empirical Distribution Function

Recall that  $F(y, \mathbf{x}) = \Pr(y_i \leq y, \mathbf{x}_i \leq \mathbf{x}) = \mathbb{E}(1(y_i \leq y) 1(\mathbf{x}_i \leq \mathbf{x}))$ , where  $1(\cdot)$  is the indicator function. This is a population moment. The method of moments estimator is the corresponding

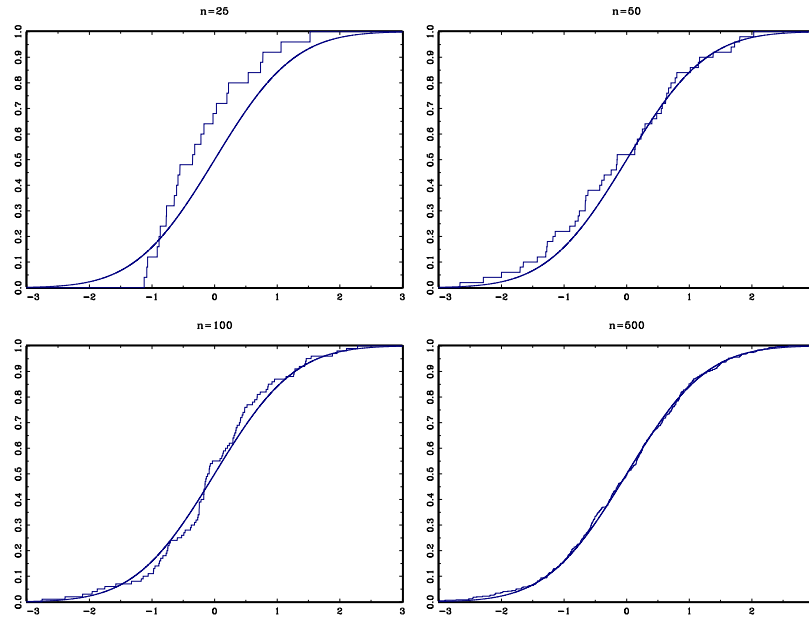


Figure 13.1: Empirical Distribution Functions

sample moment:

$$\hat{F}(y, \mathbf{x}) = \frac{1}{n} \sum_{i=1}^n 1(y_i \leq y) 1(\mathbf{x}_i \leq \mathbf{x}). \quad (13.2)$$

$\hat{F}(y, \mathbf{x})$  is called the empirical distribution function (EDF) and is a nonparametric estimate of  $F$ . Note that while  $F$  may be either discrete or continuous,  $\hat{F}$  is by construction a step function.

The EDF is a consistent estimator of the CDF. To see this, note that for any  $(y, \mathbf{x})$ ,  $1(y_i \leq y) 1(\mathbf{x}_i \leq \mathbf{x})$  is an iid random variable with expectation  $F(y, \mathbf{x})$ . Thus by the WLLN (Theorem 6.4.2),  $\hat{F}(y, \mathbf{x}) \xrightarrow{p} F(y, \mathbf{x})$ . Furthermore, by the CLT (Theorem 6.8.1),

$$\sqrt{n} \left( \hat{F}(y, \mathbf{x}) - F(y, \mathbf{x}) \right) \xrightarrow{d} N(0, F(y, \mathbf{x}) (1 - F(y, \mathbf{x}))).$$

To see the effect of sample size on the EDF, in Figure 13.1, I have plotted the EDF and true CDF for three random samples of size  $n = 25, 50, 100$ , and  $500$ . The random draws are from the  $N(0, 1)$  distribution. For  $n = 25$ , the EDF is only a crude approximation to the CDF, but the approximation appears to improve for the large  $n$ . In general, as the sample size gets larger, the EDF step function gets uniformly close to the true CDF.

The EDF is a valid discrete probability distribution which puts probability mass  $1/n$  at each pair  $(y_i, \mathbf{x}_i)$ ,  $i = 1, \dots, n$ . Notationally, it is helpful to think of a random pair  $(y_i^*, \mathbf{x}_i^*)$  with the distribution  $\hat{F}$ . That is,

$$\Pr(y_i^* \leq y, \mathbf{x}_i^* \leq \mathbf{x}) = \hat{F}(y, \mathbf{x}).$$

We can easily calculate the moments of functions of  $(y_i^*, \mathbf{x}_i^*)$ :

$$\begin{aligned} \mathbb{E}(h(y_i^*, \mathbf{x}_i^*)) &= \int h(y, \mathbf{x}) d\hat{F}(y, \mathbf{x}) \\ &= \sum_{i=1}^n h(y_i, \mathbf{x}_i) \Pr(y_i^* = y_i, \mathbf{x}_i^* = \mathbf{x}_i) \\ &= \frac{1}{n} \sum_{i=1}^n h(y_i, \mathbf{x}_i), \end{aligned}$$

the empirical sample average.

### 13.3 Nonparametric Bootstrap

The **nonparametric bootstrap** is obtained when the bootstrap distribution (13.1) is defined using the EDF (13.2) as the estimate  $\hat{F}$  of  $F$ .

Since the EDF  $\hat{F}$  is a multinomial (with  $n$  support points), in principle the distribution  $G_n^*$  could be calculated by direct methods. However, as there are  $\binom{2n-1}{n}$  possible samples  $\{(y_1^*, \mathbf{x}_1^*), \dots, (y_n^*, \mathbf{x}_n^*)\}$ , such a calculation is computationally infeasible. The popular alternative is to use simulation to approximate the distribution. The algorithm is identical to our discussion of Monte Carlo simulation, with the following points of clarification:

- The sample size  $n$  used for the simulation is the same as the sample size.
- The random vectors  $(y_i^*, \mathbf{x}_i^*)$  are drawn randomly from the empirical distribution. This is equivalent to sampling a pair  $(y_i, \mathbf{x}_i)$  randomly from the sample.

The bootstrap statistic  $T^* = T((y_1^*, \mathbf{x}_1^*), \dots, (y_n^*, \mathbf{x}_n^*), \hat{F})$  is calculated for each bootstrap sample. This is repeated  $B$  times.  $B$  is known as the number of bootstrap replications. A theory for the determination of the number of bootstrap replications  $B$  has been developed by Andrews and Buchinsky (2000). It is desirable for  $B$  to be large, so long as the computational costs are reasonable.  $B = 1000$  typically suffices.

When the statistic  $T$  is a function of  $F$ , it is typically through dependence on a parameter. For example, the t-ratio  $(\hat{\theta} - \theta) / s(\hat{\theta})$  depends on  $\theta$ . As the bootstrap statistic replaces  $F$  with  $\hat{F}$ , it similarly replaces  $\theta$  with  $\theta^* = \theta(\hat{F})$ , the value of  $\theta$  implied by  $\hat{F}$ . Typically  $\theta^* = \hat{\theta}$ , the parameter estimate. (When in doubt use  $\hat{\theta}$ .)

Sampling from the EDF is particularly easy. Since  $\hat{F}$  is a discrete probability distribution putting probability mass  $1/n$  at each sample point, sampling from the EDF is equivalent to random sampling a pair  $(y_i, \mathbf{x}_i)$  from the observed data **with replacement**. In consequence, a bootstrap sample  $\{(y_1^*, \mathbf{x}_1^*), \dots, (y_n^*, \mathbf{x}_n^*)\}$  will necessarily have some ties and multiple values, which is generally not a problem.

### 13.4 Bootstrap Estimation of Bias and Variance

The bias of  $\hat{\theta}$  is  $\tau = \mathbb{E}(\hat{\theta} - \theta)$ . The bootstrap counterparts are  $\hat{\theta}^* = \hat{\theta}((y_1^*, \mathbf{x}_1^*), \dots, (y_n^*, \mathbf{x}_n^*))$  and  $\tau^* = \mathbb{E}(\hat{\theta}^* - \theta^*)$ . The latter can be estimated by the simulation described in the previous section. This estimator is

$$\begin{aligned} \hat{\tau}^* &= \frac{1}{B} \sum_{b=1}^B (\hat{\theta}_b^* - \hat{\theta}) \\ &= \overline{\hat{\theta}^*} - \hat{\theta}. \end{aligned}$$

If  $\hat{\theta}$  is biased, it might be desirable to construct a biased-corrected estimator for  $\theta$  (one with reduced bias). Ideally, this would be

$$\tilde{\theta} = \hat{\theta} - \tau,$$

but  $\tau$  is unknown. The (estimated) bootstrap biased-corrected estimator is

$$\begin{aligned} \tilde{\theta}^* &= \hat{\theta} - \hat{\tau}^* \\ &= \hat{\theta} - (\overline{\hat{\theta}^*} - \hat{\theta}) \\ &= 2\hat{\theta} - \overline{\hat{\theta}^*}. \end{aligned}$$

Note, in particular, that the biased-corrected estimator is not  $\widehat{\theta}^*$ . Intuitively, the bootstrap makes the following experiment. Suppose that  $\widehat{\theta}$  is the truth. Then what is the average value of  $\widehat{\theta}$  calculated from such samples? The answer is  $\widehat{\theta}^*$ . If this is lower than  $\widehat{\theta}$ , this suggests that the estimator is downward-biased, so a biased-corrected estimator of  $\theta$  should be larger than  $\widehat{\theta}$ , and the best guess is the difference between  $\widehat{\theta}$  and  $\widehat{\theta}^*$ . Similarly if  $\widehat{\theta}^*$  is higher than  $\widehat{\theta}$ , then the estimator is upward-biased and the biased-corrected estimator should be lower than  $\widehat{\theta}$ .

Recall that variance of  $\widehat{\theta}$  is

$$V_n = \mathbb{E} \left( (\widehat{\theta} - \mathbb{E}(\widehat{\theta}))^2 \right).$$

The bootstrap analog is the variance of  $\widehat{\theta}^*$  which is

$$V^* = \mathbb{E} \left( (\widehat{\theta}^* - \mathbb{E}(\widehat{\theta}^*))^2 \right).$$

The simulation estimate is

$$\widehat{V}_n^* = \frac{1}{B} \sum_{b=1}^B \left( \widehat{\theta}_b^* - \widehat{\theta}^* \right)^2.$$

A bootstrap standard error for  $\widehat{\theta}$  is the square root of the bootstrap estimate of variance,  $s^*(\widehat{\theta}) = \sqrt{\widehat{V}_n^*}$ . These are frequently reported in applied economics instead of asymptotic standard errors.

## 13.5 Percentile Intervals

Consider an estimator  $\widehat{\theta}$  for  $\theta$  and suppose we wish to construct a confidence interval for  $\theta$ . Let  $G_n(u, F)$  denote the distribution of  $\widehat{\theta}$  and let  $q(\alpha) = q(\alpha, F)$  denote its quantile function. This is the function which solves

$$G_n(q(\alpha), F) = \alpha.$$

Let  $q^*(\alpha) = q(\alpha, \widehat{F})$  denote the quantile function of the bootstrap distribution. Note that this function will change depending on the underlying statistic  $T$  whose distribution is  $G_n$ .

In  $100(1 - \alpha)\%$  of samples,  $\widehat{\theta}$  lies in the region  $[q(\alpha/2), q(1 - \alpha/2)]$ . This motivates a confidence interval proposed by Efron:

$$\widehat{C}_1 = [q^*(\alpha/2), \quad q^*(1 - \alpha/2)].$$

This is often called the **percentile confidence interval**.

Computationally, the quantile  $q^*(\alpha)$  is estimated by  $\widehat{q}^*(\alpha)$ , the  $\alpha^{th}$  sample quantile of the simulated statistics  $\{T_1^*, \dots, T_B^*\}$ , as discussed in the section on Monte Carlo simulation. The  $1 - \alpha$  Efron percentile interval is then  $[\widehat{q}^*(\alpha/2), \quad \widehat{q}^*(1 - \alpha/2)]$ .

The interval  $\widehat{C}_1$  is a popular bootstrap confidence interval often used in empirical practice. This is because it is easy to compute, simple to motivate, was popularized by Efron early in the history of the bootstrap, and also has the feature that it is translation invariant. That is, if we define  $\phi = f(\theta)$  as the parameter of interest for a monotonically increasing function  $f$ , then percentile method applied to this problem will produce the confidence interval  $[f(q^*(\alpha/2)), \quad f(q^*(1 - \alpha/2))]$ , which is a naturally good property.

However, as we show now,  $\widehat{C}_1$  can work poorly unless the sampling distribution of  $\widehat{\theta}$  is symmetric about  $\theta$ .

It will be useful if we introduce an alternative definition of  $\widehat{C}_1$ . Let  $q(\alpha)$  and  $q^*(\alpha)$  be the quantile functions of  $\widehat{\theta} - \theta$  and  $\widehat{\theta}^* - \widehat{\theta}$  (These are the original quantiles, with  $\theta$  and  $\widehat{\theta}$  subtracted.) Then  $\widehat{C}_1$  can alternatively be written as

$$\widehat{C}_1 = [\widehat{\theta} + q^*(\alpha/2), \quad \widehat{\theta} + q^*(1 - \alpha/2)].$$

This is a bootstrap estimate of the “ideal” confidence interval

$$\widehat{C}_1^0 = [\widehat{\theta} + q(\alpha/2), \widehat{\theta} + q(1 - \alpha/2)].$$

The latter has coverage probability

$$\begin{aligned} \Pr(\theta \in \widehat{C}_1^0) &= \Pr(\widehat{\theta} + q(\alpha/2) \leq \theta \leq \widehat{\theta} + q(1 - \alpha/2)) \\ &= \Pr(-q(1 - \alpha/2) \leq \widehat{\theta} - \theta \leq -q(\alpha/2)) \\ &= G_n(-q(\alpha/2), F) - G_n(-q(1 - \alpha/2), F) \end{aligned}$$

which generally is not  $1 - \alpha$ ! There is one important exception. If  $\widehat{\theta} - \theta$  has a symmetric distribution about 0, then  $G_n(-u, F) = 1 - G_n(u, F)$ , so

$$\begin{aligned} \Pr(\theta \in \widehat{C}_1^0) &= G_n(-q(\alpha/2), F) - G_n(-q(1 - \alpha/2), F) \\ &= (1 - G_n(q(\alpha/2), F)) - (1 - G_n(q(1 - \alpha/2), F)) \\ &= \left(1 - \frac{\alpha}{2}\right) - \left(1 - \left(1 - \frac{\alpha}{2}\right)\right) \\ &= 1 - \alpha \end{aligned}$$

and this idealized confidence interval is accurate. Therefore,  $\widehat{C}_1^0$  and  $\widehat{C}_1$  are designed for the case that  $\widehat{\theta}$  has a symmetric distribution about  $\theta$ .

When  $\widehat{\theta}$  does not have a symmetric distribution,  $\widehat{C}_1$  may perform quite poorly.

However, by the translation invariance argument presented above, it also follows that if there exists some monotonically increasing transformation  $f(\cdot)$  such that  $f(\widehat{\theta})$  is symmetrically distributed about  $f(\theta)$ , then the idealized percentile bootstrap method will be accurate.

Based on these arguments, many argue that the percentile interval should not be used unless the sampling distribution is close to unbiased and symmetric.

The problems with the percentile method can be circumvented, at least in principle, by an alternative method. Again, let  $q(\alpha)$  and  $q^*(\alpha)$  be the quantile functions of  $\widehat{\theta} - \theta$  and  $\widehat{\theta}^* - \widehat{\theta}$ . Then

$$\begin{aligned} 1 - \alpha &= \Pr(q(\alpha/2) \leq \widehat{\theta} - \theta \leq q_n(1 - \alpha/2)) \\ &= \Pr(\widehat{\theta} - q(1 - \alpha/2) \leq \theta \leq \widehat{\theta} - q(\alpha/2)), \end{aligned}$$

so an exact  $1 - \alpha$  confidence interval for  $\theta$  is

$$\widehat{C}_2^0 = [\widehat{\theta} - q(1 - \alpha/2), \widehat{\theta} - q(\alpha/2)].$$

This motivates a bootstrap analog

$$\widehat{C}_2 = [\widehat{\theta} - q^*(1 - \alpha/2), \widehat{\theta} - q^*(\alpha/2)].$$

Notice that generally this is very different from the Efron interval  $\widehat{C}_1$ ! They coincide in the special case that  $G_n^*(u)$  is symmetric about  $\widehat{\theta}$ , but otherwise they differ.

Computationally, this interval can be estimated from a bootstrap simulation by sorting the bootstrap statistics  $T^* = \widehat{\theta}^* - \widehat{\theta}$ . These are sorted to yield the quantile estimates  $\widehat{q}^*(.025)$  and  $\widehat{q}^*(.975)$ . The 95% confidence interval is then  $[\widehat{\theta} - \widehat{q}^*(.975), \widehat{\theta} - \widehat{q}^*(.025)]$ .

This confidence interval is discussed in most theoretical treatments of the bootstrap, but is not widely used in practice.

### 13.6 Percentile-t Equal-Tailed Interval

Suppose we want to test  $\mathbb{H}_0 : \theta = \theta_0$  against  $\mathbb{H}_1 : \theta < \theta_0$  at size  $\alpha$ . We would set  $T(\theta) = (\hat{\theta} - \theta) / s(\hat{\theta})$  and reject  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  if  $T(\theta_0) < c$ , where  $c$  would be selected so that

$$\Pr(T(\theta_0) < c) = \alpha.$$

Thus  $c = q(\alpha)$ . Since this is unknown, a bootstrap test replaces  $q(\alpha)$  with the bootstrap estimate  $q^*(\alpha)$ , and the test rejects if  $T(\theta_0) < q^*(\alpha)$ .

Similarly, if the alternative is  $\mathbb{H}_1 : \theta > \theta_0$ , the bootstrap test rejects if  $T(\theta_0) > q^*(1 - \alpha)$ .

Computationally, these critical values can be estimated from a bootstrap simulation by sorting the bootstrap t-statistics  $T^* = (\hat{\theta}^* - \hat{\theta}) / s(\hat{\theta}^*)$ . Note, and this is important, that the bootstrap test statistic is centered at the estimate  $\hat{\theta}$ , and the standard error  $s(\hat{\theta}^*)$  is calculated on the bootstrap sample. These t-statistics are sorted to find the estimated quantiles  $\hat{q}^*(\alpha)$  and/or  $\hat{q}^*(1 - \alpha)$ .

Let  $T(\theta) = (\hat{\theta} - \theta) / s(\hat{\theta})$ . Then taking the intersection of two one-sided intervals,

$$\begin{aligned} 1 - \alpha &= \Pr(q(\alpha/2) \leq T(\theta_0) \leq q(1 - \alpha/2)) \\ &= \Pr\left(q(\alpha/2) \leq (\hat{\theta} - \theta_0) / s(\hat{\theta}) \leq q(1 - \alpha/2)\right) \\ &= \Pr\left(\hat{\theta} - s(\hat{\theta})q(1 - \alpha/2) \leq \theta_0 \leq \hat{\theta} - s(\hat{\theta})q(\alpha/2)\right). \end{aligned}$$

An exact  $(1 - \alpha)\%$  confidence interval for  $\theta$  is

$$\hat{C}_3^0 = [\hat{\theta} - s(\hat{\theta})q(1 - \alpha/2), \quad \hat{\theta} - s(\hat{\theta})q(\alpha/2)].$$

This motivates a bootstrap analog

$$\hat{C}_3 = [\hat{\theta} - s(\hat{\theta})q^*(1 - \alpha/2), \quad \hat{\theta} - s(\hat{\theta})q^*(\alpha/2)].$$

This is often called a **percentile-t confidence interval**. It is **equal-tailed** or **central** since the probability that  $\theta$  is below the left endpoint approximately equals the probability that  $\theta$  is above the right endpoint, each  $\alpha/2$ .

Computationally, this is based on the critical values from the one-sided hypothesis tests, discussed above.

### 13.7 Symmetric Percentile-t Intervals

Suppose we want to test  $\mathbb{H}_0 : \theta = \theta_0$  against  $\mathbb{H}_1 : \theta \neq \theta_0$  at size  $\alpha$ . We would set  $T(\theta) = (\hat{\theta} - \theta) / s(\hat{\theta})$  and reject  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  if  $|T(\theta_0)| > c$ , where  $c$  would be selected so that

$$\Pr(|T(\theta_0)| > c) = \alpha.$$

Note that

$$\begin{aligned} \Pr(|T(\theta_0)| < c) &= \Pr(-c < T(\theta_0) < c) \\ &= G_n(c) - G_n(-c) \\ &\equiv \overline{G}_n(c), \end{aligned}$$

which is a symmetric distribution function. The ideal critical value  $c = q(1 - \alpha)$  solves the equation

$$\overline{G}_n(q(1 - \alpha)) = 1 - \alpha.$$

Equivalently,  $q(1 - \alpha)$  is the  $1 - \alpha$  quantile of the distribution of  $|T(\theta_0)|$ .

The bootstrap estimate is  $q^*(1 - \alpha)$ , the  $1 - \alpha$  quantile of the distribution of  $|T^*|$ , or the number which solves the equation

$$\overline{G}_n^*(q^*(1 - \alpha)) = G_n^*(q^*(1 - \alpha)) - G_n^*(-q^*(1 - \alpha)) = 1 - \alpha.$$

Computationally,  $q^*(1 - \alpha)$  is estimated from a bootstrap simulation by sorting the bootstrap t-statistics  $|T^*| = |\hat{\theta}^* - \hat{\theta}|/s(\hat{\theta}^*)$ , and taking the  $1 - \alpha$  quantile. The bootstrap test rejects if  $|T(\theta_0)| > q^*(1 - \alpha)$ .

Let

$$\hat{C}_4 = [\hat{\theta} - s(\hat{\theta})q^*(1 - \alpha), \quad \hat{\theta} + s(\hat{\theta})q^*(1 - \alpha)],$$

where  $q_n^*(1 - \alpha)$  is the bootstrap critical value for a two-sided hypothesis test.  $\hat{C}_4$  is called the symmetric percentile-t interval. It is designed to work well since

$$\begin{aligned} \Pr(\theta \in \hat{C}_4) &= \Pr(\hat{\theta} - s(\hat{\theta})q^*(1 - \alpha) \leq \theta \leq \hat{\theta} + s(\hat{\theta})q^*(1 - \alpha)) \\ &= \Pr(|T(\theta)| < q^*(1 - \alpha)) \\ &\simeq \Pr(|T(\theta)| < q(1 - \alpha)) \\ &= 1 - \alpha. \end{aligned}$$

If  $\theta$  is a vector, then to test  $\mathbb{H}_0 : \theta = \theta_0$  against  $\mathbb{H}_1 : \theta \neq \theta_0$  at size  $\alpha$ , we would use a Wald statistic

$$W(\theta) = n(\hat{\theta} - \theta)' \hat{\mathbf{V}}_{\theta}^{-1} (\hat{\theta} - \theta)$$

or a similar asymptotically chi-square statistic. The ideal test rejects if  $W \geq q(1 - \alpha)$ , where  $q(1 - \alpha)$  is the  $1 - \alpha$  quantile of the distribution of  $W$ . The bootstrap test rejects if  $W \geq q^*(1 - \alpha)$ , where  $q^*(1 - \alpha)$  is the  $1 - \alpha$  quantile of the distribution of

$$W^* = n(\hat{\theta}^* - \hat{\theta})' \hat{\mathbf{V}}_{\theta}^{*-1} (\hat{\theta}^* - \hat{\theta}).$$

Computationally, the critical value  $q^*(1 - \alpha)$  is found as the quantile from simulated values of  $W^*$ . Note in the simulation that the Wald statistic is a quadratic form in  $(\hat{\theta}^* - \hat{\theta})$ , not  $(\hat{\theta}^* - \theta_0)$ . (The latter is a common mistake made by practitioners.)

## 13.8 Asymptotic Expansions

Let  $T \in \mathbb{R}$  be a statistic such that

$$T \xrightarrow{d} N(0, \sigma^2). \quad (13.3)$$

In some cases, such as when  $T$  is a t-ratio, then  $\sigma^2 = 1$ . In other cases  $\sigma^2$  is unknown. Equivalently, writing  $T \sim G_n(u, F)$  then for each  $u$  and  $F$

$$\lim_{n \rightarrow \infty} G_n(u, F) = \Phi\left(\frac{u}{\sigma}\right),$$

or

$$G_n(u, F) = \Phi\left(\frac{u}{\sigma}\right) + o(1). \quad (13.4)$$

While (13.4) says that  $G_n$  converges to  $\Phi\left(\frac{u}{\sigma}\right)$  as  $n \rightarrow \infty$ , it says nothing, however, about the rate of convergence, or the size of the divergence for any particular sample size  $n$ . A better asymptotic approximation may be obtained through an **asymptotic expansion**.



Notationally, it is useful to recall the stochastic order notation of Section 6.13. Also, it is convenient to define even and odd functions. We say that a function  $g(u)$  is **even** if  $g(-u) = g(u)$ , and a function  $h(u)$  is **odd** if  $h(-u) = -h(u)$ . The derivative of an even function is odd, and vice-versa.

**Theorem 13.8.1** *Under regularity conditions and (13.3),*

$$G_n(u, F) = \Phi\left(\frac{u}{\sigma}\right) + \frac{1}{n^{1/2}}g_1(u, F) + \frac{1}{n}g_2(u, F) + O(n^{-3/2})$$

*uniformly over  $u$ , where  $g_1$  is an even function of  $u$ , and  $g_2$  is an odd function of  $u$ . Moreover,  $g_1$  and  $g_2$  are differentiable functions of  $u$  and continuous in  $F$  relative to the supremum norm on the space of distribution functions.*

The expansion in Theorem 13.8.1 is often called an **Edgeworth expansion**.

We can interpret Theorem 13.8.1 as follows. First,  $G_n(u, F)$  converges to the normal limit at rate  $n^{1/2}$ . To a second order of approximation,

$$G_n(u, F) \approx \Phi\left(\frac{u}{\sigma}\right) + n^{-1/2}g_1(u, F).$$

Since the derivative of  $g_1$  is odd, the density function is skewed. To a third order of approximation,

$$G_n(u, F) \approx \Phi\left(\frac{u}{\sigma}\right) + n^{-1/2}g_1(u, F) + n^{-1}g_2(u, F)$$

which adds a symmetric non-normal component to the approximate density (for example, adding leptokurtosis).

As a side note, when  $T = \sqrt{n}(\bar{X}_n - \mu)/\sigma$ , a standardized sample mean, then

$$\begin{aligned} g_1(u) &= -\frac{1}{6}\kappa_3(u^2 - 1)\phi(u) \\ g_2(u) &= -\left(\frac{1}{24}\kappa_4(u^3 - 3u) + \frac{1}{72}\kappa_3^2(u^5 - 10u^3 + 15u)\right)\phi(u) \end{aligned}$$

where  $\phi(u)$  is the standard normal pdf, and

$$\begin{aligned} \kappa_3 &= \mathbb{E}\left((X - \mu)^3\right)/\sigma^3 \\ \kappa_4 &= \mathbb{E}\left((X - \mu)^4\right)/\sigma^4 - 3 \end{aligned}$$

the standardized skewness and excess kurtosis of the distribution of  $X$ . Note that when  $\kappa_3 = 0$  and  $\kappa_4 = 0$ , then  $g_1 = 0$  and  $g_2 = 0$ , so the second-order Edgeworth expansion corresponds to the normal distribution.

### Francis Edgeworth

Francis Ysidro Edgeworth (1845-1926) of Ireland, founding editor of the *Economic Journal*, was a profound economic and statistical theorist, developing the theories of indifference curves and asymptotic expansions. He also could be viewed as the first econometrician due to his early use of mathematical statistics in the study of economic data.

### 13.9 One-Sided Tests

Using the expansion of Theorem 13.8.1, we can assess the accuracy of one-sided hypothesis tests and confidence regions based on an asymptotically normal t-ratio  $T$ . An asymptotic test is based on  $\Phi(u)$ .

To the second order, the exact distribution is

$$\Pr(T < u) = G_n(u, F) = \Phi(u) + \frac{1}{n^{1/2}}g_1(u, F) + O(n^{-1})$$

since  $\sigma = 1$ . The difference is

$$\begin{aligned}\Phi(u) - G_n(u, F) &= \frac{1}{n^{1/2}}g_1(u, F) + O(n^{-1}) \\ &= O(n^{-1/2}),\end{aligned}$$

so the order of the error is  $O(n^{-1/2})$ .

A bootstrap test is based on  $G_n^*(u)$ , which from Theorem 13.8.1 has the expansion

$$G_n^*(u) = G_n(u, \hat{F}) = \Phi(u) + \frac{1}{n^{1/2}}g_1(u, \hat{F}) + O(n^{-1}).$$

Because  $\Phi(u)$  appears in both expansions, the difference between the bootstrap distribution and the true distribution is

$$G_n^*(u) - G_n(u, F) = \frac{1}{n^{1/2}} \left( g_1(u, \hat{F}) - g_1(u, F) \right) + O(n^{-1}).$$

Since  $\hat{F}$  converges to  $F$  at rate  $\sqrt{n}$ , and  $g_1$  is continuous with respect to  $F$ , the difference  $(g_1(u, \hat{F}) - g_1(u, F))$  converges to 0 at rate  $\sqrt{n}$ . Heuristically,

$$\begin{aligned}g_1(u, \hat{F}) - g_1(u, F) &\approx \frac{\partial}{\partial F}g_1(u, F) (\hat{F} - F) \\ &= O(n^{-1/2}),\end{aligned}$$

The “derivative”  $\frac{\partial}{\partial F}g_1(u, F)$  is only heuristic, as  $F$  is a function. We conclude that

$$G_n^*(u) - G_n(u, F) = O(n^{-1}),$$

or

$$\Pr(T^* \leq u) = \Pr(T \leq u) + O(n^{-1}),$$

which is an improved rate of convergence over the asymptotic test (which converged at rate  $O(n^{-1/2})$ ). This rate can be used to show that one-tailed bootstrap inference based on the t-ratio achieves a so-called **asymptotic refinement** – the Type I error of the test converges at a faster rate than an analogous asymptotic test.

### 13.10 Symmetric Two-Sided Tests

If a random variable  $y$  has distribution function  $H(u) = \Pr(y \leq u)$ , then the random variable  $|y|$  has distribution function

$$\overline{H}(u) = H(u) - H(-u)$$

since

$$\begin{aligned}\Pr(|y| \leq u) &= \Pr(-u \leq y \leq u) \\ &= \Pr(y \leq u) - \Pr(y \leq -u) \\ &= H(u) - H(-u).\end{aligned}$$

For example, if  $Z \sim N(0, 1)$ , then  $|Z|$  has distribution function

$$\overline{\Phi}(u) = \Phi(u) - \Phi(-u) = 2\Phi(u) - 1.$$

Similarly, if  $T$  has exact distribution  $G_n(u, F)$ , then  $|T|$  has the distribution function

$$\overline{G}_n(u, F) = G_n(u, F) - G_n(-u, F).$$

A two-sided hypothesis test rejects  $\mathbb{H}_0$  for large values of  $|T|$ . Since  $T \xrightarrow{d} Z$ , then  $|T| \xrightarrow{d} |Z| \sim \overline{\Phi}$ . Thus asymptotic critical values are taken from the  $\overline{\Phi}$  distribution, and exact critical values are taken from the  $\overline{G}_n(u, F)$  distribution. From Theorem 13.8.1, we can calculate that

$$\begin{aligned} \overline{G}_n(u, F) &= G_n(u, F) - G_n(-u, F) \\ &= \left( \Phi(u) + \frac{1}{n^{1/2}}g_1(u, F) + \frac{1}{n}g_2(u, F) \right) \\ &\quad - \left( \Phi(-u) + \frac{1}{n^{1/2}}g_1(-u, F) + \frac{1}{n}g_2(-u, F) \right) + O(n^{-3/2}) \\ &= \overline{\Phi}(u) + \frac{2}{n}g_2(u, F) + O(n^{-3/2}), \end{aligned} \tag{13.5}$$

where the simplifications are because  $g_1$  is even and  $g_2$  is odd. Hence the difference between the asymptotic distribution and the exact distribution is

$$\overline{\Phi}(u) - \overline{G}_n(u, F_0) = \frac{2}{n}g_2(u, F_0) + O(n^{-3/2}) = O(n^{-1}).$$

The order of the error is  $O(n^{-1})$ .

Interestingly, the asymptotic two-sided test has a better coverage rate than the asymptotic one-sided test. This is because the first term in the asymptotic expansion,  $g_1$ , is an even function, meaning that the errors in the two directions exactly cancel out.

Applying (13.5) to the bootstrap distribution, we find

$$\overline{G}_n^*(u) = \overline{G}_n(u, \hat{F}) = \overline{\Phi}(u) + \frac{2}{n}g_2(u, \hat{F}) + O(n^{-3/2}).$$

Thus the difference between the bootstrap and exact distributions is

$$\begin{aligned} \overline{G}_n^*(u) - \overline{G}_n(u, F) &= \frac{2}{n} \left( g_2(u, \hat{F}) - g_2(u, F) \right) + O(n^{-3/2}) \\ &= O(n^{-3/2}), \end{aligned}$$

the last equality because  $\hat{F}$  converges to  $F$  at rate  $\sqrt{n}$ , and  $g_2$  is continuous in  $F$ . Another way of writing this is

$$\Pr(|T^*| < u) = \Pr(|T| < u) + O(n^{-3/2})$$

so the error from using the bootstrap distribution (relative to the true unknown distribution) is  $O(n^{-3/2})$ . This is in contrast to the use of the asymptotic distribution, whose error is  $O(n^{-1})$ . Thus a two-sided bootstrap test also achieves an asymptotic refinement, similar to a one-sided test.

A reader might get confused between the two simultaneous effects. Two-sided tests have better rates of convergence than the one-sided tests, and bootstrap tests have better rates of convergence than asymptotic tests.

The analysis shows that there may be a trade-off between one-sided and two-sided tests. Two-sided tests will have more accurate size (Reported Type I error), but one-sided tests might have more power against alternatives of interest. Confidence intervals based on the bootstrap can be asymmetric if based on one-sided tests (equal-tailed intervals) and can therefore be more informative and have smaller length than symmetric intervals. Therefore, the choice between symmetric and equal-tailed confidence intervals is unclear, and needs to be determined on a case-by-case basis.

### 13.11 Percentile Confidence Intervals

To evaluate the coverage rate of the percentile interval, set  $T = \sqrt{n}(\hat{\theta} - \theta)$ . We know that  $T \xrightarrow{d} N(0, V)$ , which is not pivotal, as it depends on the unknown  $V$ . Theorem 13.8.1 shows that a first-order approximation

$$G_n(u, F) = \Phi\left(\frac{u}{\sigma}\right) + O(n^{-1/2}),$$

where  $\sigma = \sqrt{V}$ , and for the bootstrap

$$G_n^*(u) = G_n(u, \hat{F}) = \Phi\left(\frac{u}{\hat{\sigma}}\right) + O(n^{-1/2}),$$

where  $\hat{\sigma} = V(\hat{F})$  is the bootstrap estimate of  $\sigma$ . The difference is

$$\begin{aligned} G_n^*(u) - G_n(u, F) &= \Phi\left(\frac{u}{\hat{\sigma}}\right) - \Phi\left(\frac{u}{\sigma}\right) + O(n^{-1/2}) \\ &= -\phi\left(\frac{u}{\sigma}\right) \frac{u}{\sigma} (\hat{\sigma} - \sigma) + O(n^{-1/2}) \\ &= O(n^{-1/2}) \end{aligned}$$

Hence the order of the error is  $O(n^{-1/2})$ .

The good news is that the percentile-type methods (if appropriately used) can yield  $\sqrt{n}$ -convergent asymptotic inference. Yet these methods do not require the calculation of standard errors! This means that in contexts where standard errors are not available or are difficult to calculate, the percentile bootstrap methods provide an attractive inference method.

The bad news is that the rate of convergence is disappointing. It is no better than the rate obtained from an asymptotic one-sided confidence region. Therefore if standard errors are available, it is unclear if there are any benefits from using the percentile bootstrap over simple asymptotic methods.

Based on these arguments, the theoretical literature (e.g. Hall, 1992, Horowitz, 2001) tends to advocate the use of the percentile-t bootstrap methods rather than percentile methods.

### 13.12 Bootstrap Methods for Regression Models

The bootstrap methods we have discussed have set  $G_n^*(u) = G_n(u, \hat{F})$ , where  $\hat{F}$  is the EDF. Any other consistent estimate of  $F$  may be used to define a feasible bootstrap estimator. The advantage of the EDF is that it is fully nonparametric, it imposes no conditions, and works in nearly any context. But since it is fully nonparametric, it may be inefficient in contexts where more is known about  $F$ . We discuss bootstrap methods appropriate for the linear regression model

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \mathbb{E}(e_i \mid \mathbf{x}_i) &= 0. \end{aligned}$$

The non-parametric bootstrap resamples the observations  $(y_i^*, \mathbf{x}_i^*)$  from the EDF, which implies

$$\begin{aligned} y_i^* &= \mathbf{x}_i^{*'} \hat{\boldsymbol{\beta}} + e_i^* \\ \mathbb{E}(\mathbf{x}_i^* e_i^*) &= \mathbf{0} \end{aligned}$$

but generally

$$\mathbb{E}(e_i^* \mid \mathbf{x}_i^*) \neq 0.$$

The bootstrap distribution does not impose the regression assumption, and is thus an inefficient estimator of the true distribution (when in fact the regression assumption is true.)

One approach to this problem is to impose the very strong assumption that the error  $\varepsilon_i$  is independent of the regressor  $\mathbf{x}_i$ . The advantage is that in this case it is straightforward to construct bootstrap distributions. The disadvantage is that the bootstrap distribution may be a poor approximation when the error is not independent of the regressors.

To impose independence, it is sufficient to sample the  $\mathbf{x}_i^*$  and  $e_i^*$  independently, and then create  $y_i^* = \mathbf{x}_i^{*\prime} \hat{\boldsymbol{\beta}} + e_i^*$ . There are different ways to impose independence. A non-parametric method is to sample the bootstrap errors  $e_i^*$  randomly from the OLS residuals  $\{\hat{e}_1, \dots, \hat{e}_n\}$ . A parametric method is to generate the bootstrap errors  $e_i^*$  from a parametric distribution, such as the normal  $e_i^* \sim N(0, \hat{\sigma}^2)$ .

For the regressors  $\mathbf{x}_i^*$ , a nonparametric method is to sample the  $\mathbf{x}_i^*$  randomly from the EDF or sample values  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ . A parametric method is to sample  $\mathbf{x}_i^*$  from an estimated parametric distribution. A third approach sets  $\mathbf{x}_i^* = \mathbf{x}_i$ . This is equivalent to treating the regressors as **fixed in repeated samples**. If this is done, then all inferential statements are made conditionally on the observed values of the regressors, which is a valid statistical approach. It does not really matter, however, whether or not the  $\mathbf{x}_i$  are really “fixed” or random.

The methods discussed above are unattractive for most applications in econometrics because they impose the stringent assumption that  $\mathbf{x}_i$  and  $e_i$  are independent. Typically what is desirable is to impose only the regression condition  $\mathbb{E}(e_i | \mathbf{x}_i) = 0$ . Unfortunately this is a harder problem.

One proposal which imposes the regression condition without independence is the **Wild Bootstrap**. The idea is to construct a conditional distribution for  $e_i^*$  so that

$$\begin{aligned}\mathbb{E}(e_i^* | \mathbf{x}_i) &= 0 \\ \mathbb{E}(e_i^{*2} | \mathbf{x}_i) &= \hat{e}_i^2 \\ \mathbb{E}(e_i^{*3} | \mathbf{x}_i) &= \hat{e}_i^3.\end{aligned}$$

A conditional distribution with these features will preserve the main important features of the data. This can be achieved using a two-point distribution of the form

$$\begin{aligned}\Pr\left(e_i^* = \left(\frac{1 + \sqrt{5}}{2}\right) \hat{e}_i\right) &= \frac{\sqrt{5} - 1}{2\sqrt{5}} \\ \Pr\left(e_i^* = \left(\frac{1 - \sqrt{5}}{2}\right) \hat{e}_i\right) &= \frac{\sqrt{5} + 1}{2\sqrt{5}}\end{aligned}$$

For each  $\mathbf{x}_i$ , you sample  $e_i^*$  using this two-point distribution.

### 13.13 Bootstrap GMM Inference

Consider an unconditional moment model

$$\mathbb{E}(\mathbf{g}_i(\boldsymbol{\beta})) = 0$$

and let  $\hat{\boldsymbol{\beta}}$  be the 2SLS or GMM estimator of  $\boldsymbol{\beta}$ . Using the EDF of  $\mathbf{w}_i = (y_i, \mathbf{z}_i, \mathbf{x}_i)$ , we can apply bootstrap methods to compute estimates of the bias and variance of  $\hat{\boldsymbol{\beta}}$  and construct confidence intervals for  $\boldsymbol{\beta}$ , identically as in the regression model. However, caution should be applied when interpreting such results.

A straightforward application of the nonparametric bootstrap works in the sense of consistently achieving the first-order asymptotic distribution. This has been shown by Hahn (1996). However, it fails to achieve an asymptotic refinement when the model is over-identified, jeopardizing the theoretical justification for percentile-t methods. Furthermore, the bootstrap applied  $J$  test will yield the wrong answer.

The problem is that in the sample,  $\hat{\beta}$  is the “true” value and yet  $\bar{g}_n(\hat{\beta}) \neq 0$ . Thus according to random variables  $(y_i^*, z_i^*, x_i^*)$  drawn from the EDF  $F_n$ ,

$$\mathbb{E} \left( g_i(\hat{\beta}) \right) = \bar{g}_n(\hat{\beta}) \neq \mathbf{0}.$$

This means that  $(y_i^*, z_i^*, x_i^*)$  do not satisfy the same moment conditions as the population distribution.

A correction suggested by Hall and Horowitz (1996) can solve the problem. Given the bootstrap sample  $(y^*, Z^*, X^*)$ , define the bootstrap GMM criterion

$$J^*(\beta) = n \cdot \left( \bar{g}_n^*(\beta) - \bar{g}_n(\hat{\beta}) \right)' \widehat{W}^* \left( \bar{g}_n^*(\beta) - \bar{g}_n(\hat{\beta}) \right)$$

where  $\bar{g}_n(\hat{\beta})$  is from the in-sample data, not from the bootstrap data.

Let  $\hat{\beta}^*$  minimize  $J^*(\beta)$ , and define all statistics and tests accordingly. In the linear model, this implies that the bootstrap estimator is

$$\hat{\beta}^* = (X^{*'} Z^* W^* Z^{*'} X^*)^{-1} \left( X^{*'} Z^* \widehat{W}^* (Z^{*'} y^* - Z^{*'} \hat{e}) \right).$$

where  $\hat{e} = y - X\hat{\beta}$  are the in-sample residuals. The bootstrap J statistic is  $J^*(\hat{\beta}^*)$ .

## Exercises

**Exercise 13.1** Let  $\widehat{F}(\mathbf{x})$  denote the EDF of a random sample. Show that

$$\sqrt{n} \left( \widehat{F}(\mathbf{x}) - F(\mathbf{x}) \right) \xrightarrow{d} N(0, F(\mathbf{x})(1 - F(\mathbf{x}))).$$

**Exercise 13.2** Take a random sample  $\{y_1, \dots, y_n\}$  with  $\mu = \mathbb{E}(y_i)$  and  $\sigma^2 = \text{var}(y_i)$  and set  $\bar{y}_n = n^{-1} \sum_{i=1}^n y_i$ . Find the population moments  $\mathbb{E}(\bar{y}_n)$  and  $\text{var}(\bar{y}_n)$ . Let  $\{y_1^*, \dots, y_n^*\}$  be a random sample from the empirical distribution function and set  $\bar{y}_n^* = n^{-1} \sum_{i=1}^n y_i^*$ . Find the bootstrap moments  $\mathbb{E}(\bar{y}_n^*)$  and  $\text{var}(\bar{y}_n^*)$ .

**Exercise 13.3** Consider the following bootstrap procedure for a regression of  $y_i$  on  $\mathbf{x}_i$ . Let  $\widehat{\beta}$  denote the OLS estimator from the regression of  $\mathbf{y}$  on  $\mathbf{X}$ , and  $\widehat{\mathbf{e}} = \mathbf{y} - \mathbf{X}\widehat{\beta}$  the OLS residuals.

- Draw a random vector  $(\mathbf{x}^*, e^*)$  from the pair  $\{(\mathbf{x}_i, \widehat{e}_i) : i = 1, \dots, n\}$ . That is, draw a random integer  $i'$  from  $[1, 2, \dots, n]$ , and set  $\mathbf{x}^* = \mathbf{x}_{i'}$  and  $e^* = \widehat{e}_{i'}$ . Set  $y^* = \mathbf{x}^{*'}\widehat{\beta} + e^*$ . Draw (with replacement)  $n$  such vectors, creating a random bootstrap data set  $(\mathbf{y}^*, \mathbf{X}^*)$ .
- Regress  $\mathbf{y}^*$  on  $\mathbf{X}^*$ , yielding OLS estimates  $\widehat{\beta}^*$  and any other statistic of interest.

Show that this bootstrap procedure is (numerically) identical to the non-parametric bootstrap.

**Exercise 13.4** Consider the following bootstrap procedure. Using the non-parametric bootstrap, generate bootstrap samples, calculate the estimate  $\widehat{\theta}^*$  on these samples and then calculate

$$T^* = (\widehat{\theta}^* - \widehat{\theta})/s(\widehat{\theta}),$$

where  $s(\widehat{\theta})$  is the standard error in the original data. Let  $q^*(.05)$  and  $q^*(.95)$  denote the 5% and 95% quantiles of  $T^*$ , and define the bootstrap confidence interval

$$\widehat{C} = \left[ \widehat{\theta} - s(\widehat{\theta})q^*(.95), \quad \widehat{\theta} - s(\widehat{\theta})q^*(.05) \right].$$

Show that  $\widehat{C}$  exactly equals the Alternative percentile interval (not the percentile-t interval).

**Exercise 13.5** You want to test  $\mathbb{H}_0 : \theta = 0$  against  $\mathbb{H}_1 : \theta > 0$ . The test for  $\mathbb{H}_0$  is to reject if  $T_n = \widehat{\theta}/s(\widehat{\theta}) > c$  where  $c$  is picked so that Type I error is  $\alpha$ . You do this as follows. Using the non-parametric bootstrap, you generate bootstrap samples, calculate the estimates  $\widehat{\theta}^*$  on these samples and then calculate

$$T^* = \widehat{\theta}^*/s(\widehat{\theta}^*).$$

Let  $q^*(.95)$  denote the 95% quantile of  $T^*$ . You replace  $c$  with  $q^*(.95)$ , and thus reject  $\mathbb{H}_0$  if  $T_n = \widehat{\theta}/s(\widehat{\theta}) > q^*(.95)$ . What is wrong with this procedure?

**Exercise 13.6** Suppose that in an application,  $\widehat{\theta} = 1.2$  and  $s(\widehat{\theta}) = .2$ . Using the non-parametric bootstrap, 1000 samples are generated from the bootstrap distribution, and  $\widehat{\theta}^*$  is calculated on each sample. The  $\widehat{\theta}^*$  are sorted, and the 2.5% and 97.5% quantiles of the  $\widehat{\theta}^*$  are .75 and 1.3, respectively.

- Report the 95% Efron Percentile interval for  $\theta$ .
- Report the 95% Alternative Percentile interval for  $\theta$ .
- With the given information, can you report the 95% Percentile-t interval for  $\theta$ ?

**Exercise 13.7** Consider the model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(e_i | \mathbf{x}_i) = 0$$

with  $y_i$  scalar and  $\mathbf{x}_i$  a  $k$  vector. You have a random sample  $(y_i, \mathbf{x}_i : i = 1, \dots, n)$ . You are interested in estimating the regression function  $m(\mathbf{x}) = E(y_i | \mathbf{x}_i = \mathbf{x})$  at a fixed vector  $\mathbf{x}$  and constructing a 95% confidence interval.

- Write the standard estimator and asymptotic confidence interval for  $m(\mathbf{x})$ .
- Describe the percentile bootstrap confidence interval for  $m(\mathbf{x})$ .
- Describe the percentile-t bootstrap confidence interval for  $m(\mathbf{x})$ .

**Exercise 13.8** The observed data is  $\{y_i, x_i\} \in \mathbb{R} \times \mathbb{R}^k$ ,  $k > 1$ ,  $i = 1, \dots, n$ . Take the model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(x_i e_i) = 0$$

$$\mu_3 = \mathbb{E}(e_i^3)$$

- Write down an estimator for  $\mu_3$ .
- Explain how to use the Efron percentile method to construct a 90% confidence interval for  $\mu_3$  in this specific model.

**Exercise 13.9** Take the model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(x_i e_i) = 0$$

$$\mathbb{E}(e_i^2) = \sigma^2$$

Describe the bootstrap percentile confidence interval for  $\sigma^2$ .

**Exercise 13.10** The model is

$$y_i = \mathbf{x}_{1i}' \boldsymbol{\beta}_1 + \mathbf{x}_{2i}' \boldsymbol{\beta}_2 + e_i$$

$$\mathbb{E}(\mathbf{x}_i e_i) = 0$$

with  $x_{2i}$  scalar. Describe how to test  $\mathbb{H}_0 : \beta_2 = 0$  against  $\mathbb{H}_1 : \beta_2 \neq 0$  using the nonparametric bootstrap.

**Exercise 13.11** The model is

$$y_i = \mathbf{x}_{1i}' \boldsymbol{\beta}_1 + x_{2i} \beta_2 + e_i$$

$$\mathbb{E}(\mathbf{x}_i e_i) = 0$$

with both  $\mathbf{x}_{1i}$  and  $\mathbf{x}_{2i}$   $k \times 1$ . Describe how to test  $\mathbb{H}_0 : \beta_1 = \beta_2$  against  $\mathbb{H}_1 : \beta_1 \neq \beta_2$  using the nonparametric bootstrap.

**Exercise 13.12** Suppose a PhD student has a sample  $(y_i, x_i, z_i : i = 1, \dots, n)$  and estimates by OLS the equation

$$y_i = z_i \hat{\alpha} + x_i' \hat{\boldsymbol{\beta}} + \hat{e}_i$$

where  $\alpha$  is the coefficient of interest and she is interested in testing  $\mathbb{H}_0 : \alpha = 0$  against  $\mathbb{H}_1 : \alpha \neq 0$ . She obtains  $\hat{\alpha} = 2.0$  with standard error  $s(\hat{\alpha}) = 1.0$  so the value of the t-ratio for  $\mathbb{H}_0$  is  $T = \hat{\alpha}/s(\hat{\alpha}) = 2.0$ . To assess significance, the student decides to use the bootstrap. She uses the following algorithm



1. Samples  $(y_i^*, x_i^*, z_i^*)$  randomly from the observations. (Random sampling with replacement). Creates a random sample with  $n$  observations.
2. On this pseudo-sample, estimates the equation

$$y_i^* = z_i^* \hat{\alpha}^* + x_i^{*'} \hat{\beta}^* + e_i^*$$

by OLS and computes standard errors, including  $s(\hat{\alpha}^*)$ . The t-ratio for  $\mathbb{H}_0$ ,  $T^* = \hat{\alpha}^*/s(\hat{\alpha}^*)$  is computed and stored.

3. This is repeated  $B = 9999$  times.
4. The 95% empirical quantile  $\hat{q}_{.95}^*$  of the bootstrap absolute t-ratios  $|T^*|$  is computed. It is  $\hat{q}_{.95}^* = 3.5$ .
5. The student notes that while  $|T| = 2 > 1.96$  (and thus an asymptotic 5% size test rejects  $\mathbb{H}_0$ ),  $|T| = 2 < \hat{q}_{.95}^* = 3.5$  and thus the bootstrap test does not reject  $\mathbb{H}_0$ . As the bootstrap is more reliable, the student concludes that  $\mathbb{H}_0$  cannot be rejected in favor of  $\mathbb{H}_1$ .

Question: Do you agree with the student's method and reasoning? Do you see an error in her method?

**Exercise 13.13** Take the model

$$\begin{aligned} y_i &= x_{1i}\beta_1 + x_{2i}\beta_2 + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= 0 \end{aligned}$$

The parameter of interest is  $\theta = \beta_1\beta_2$ . Show how to construct a confidence interval for  $\theta$  using the following three methods.

1. Asymptotic Theory
2. Percentile Bootstrap
3. Equal-Tailed Percentile-t Bootstrap.

Your answer should be specific to this problem, not general.

**Exercise 13.14** Let  $y_i$  be iid,  $\mu = \mathbb{E}(y_i) > 0$ , and  $\theta = \mu^{-1}$ . Let  $\hat{\mu} = \bar{Y}_n$  be the sample mean and  $\hat{\theta} = \hat{\mu}^{-1}$ .

- (a) Is  $\hat{\theta}$  unbiased for  $\theta$ ?
- (b) If  $\hat{\theta}$  is biased, can you determine the direction of the bias  $\mathbb{E}(\hat{\theta} - \theta)$  (up or down)?
- (c) Could the nonparametric bootstrap be used to estimate the bias? If so, explain how.

**Exercise 13.15** Take the model

$$\begin{aligned} y_i &= x_{1i}\beta_1 + x_{2i}\beta_2 + e_i \\ \mathbb{E}(\mathbf{x}_i e_i) &= 0 \\ \theta &= \frac{\beta_1}{\beta_2} \end{aligned}$$

Assume that the observations  $(y_i, x_{1i}, x_{2i})$  are i.i.d. across  $i = 1, \dots, n$ . Describe how you would construct the percentile-t bootstrap confidence interval for  $\theta$ .

**Exercise 13.16** The model is iid data,  $i = 1, \dots, n$ ,

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(e_i \mid \mathbf{x}_i) = 0$$

Does the presence of conditional heteroskedasticity invalidate the application of the non-parametric bootstrap? Explain.

**Exercise 13.17** The RESET specification test for nonlinearity in a random sample is the following. The null hypothesis is a linear regression

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(e_i \mid \mathbf{x}_i) = 0$$

The parameter  $\boldsymbol{\beta}$  is estimated by OLS yielding predicted values  $\hat{y}_i$ . Then a second-stage least-squares regression is estimated including both  $\mathbf{x}_i$  and  $\hat{y}_i$

$$y_i = \mathbf{x}_i' \tilde{\boldsymbol{\beta}} + (\hat{y}_i)^2 \tilde{\gamma} + \tilde{e}_i$$

The RESET test statistic  $R$  is the squared t-ratio on  $\tilde{\gamma}$ .

A colleague suggests obtaining the critical value for the test using the bootstrap. He proposes the following bootstrap implementation.

- Draw  $n$  observations  $(y_i^*, \mathbf{x}_i^*)$  randomly from the observed sample pairs  $(y_i, \mathbf{x}_i)$  to create a bootstrap sample.
- Compute the statistic  $R^*$  on this bootstrap sample as described above.
- Repeat this 999 times. Sort the bootstrap statistics  $R^*$ , take number 950 (the 95% percentile) and use this as the critical value.
- Reject the null hypothesis if  $R$  exceeds this critical value, otherwise do not reject.

Is this procedure a correct implementation of the bootstrap in this context? If not, propose a modified bootstrap.

**Exercise 13.18** The model is

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(\mathbf{x}_i e_i) \neq 0,$$

so the regressor  $\mathbf{x}_i$  is endogenous. We know that in this case, the OLS estimator is biased for the parameter  $\boldsymbol{\beta}$ . We also know that the non-parametric bootstrap is (generally) a good method to estimate bias, and thereby make bias-adjusted. Explain whether or not the non-parametric bootstrap can be used to estimate the bias of OLS in the above context.

**Exercise 13.19** The datafile `hprice1.txt` contains data on house prices (sales), with variables listed in the file `hprice1.pdf`. Estimate a linear regression of price on the number of bedrooms, lot size, size of house, and the colonial dummy. Calculate 95% confidence intervals for the regression coefficients using both the asymptotic normal approximation and the percentile-t bootstrap.

# Chapter 14

## Univariate Time Series

A time series  $y_t$  is a process observed in sequence over time,  $t = 1, \dots, T$ . To indicate the dependence on time, we adopt new notation, and use the subscript  $t$  to denote the individual observation, and  $T$  to denote the number of observations.

Because of the sequential nature of time series, we expect that  $y_t$  and  $y_{t-1}$  are *not* independent, so classical assumptions are not valid.

We can separate time series into two categories: univariate ( $y_t \in \mathbb{R}$  is scalar); and multivariate ( $y_t \in \mathbb{R}^m$  is vector-valued). The primary model for univariate time series is autoregressions (ARs). The primary model for multivariate time series is vector autoregressions (VARs).

### 14.1 Stationarity and Ergodicity

**Definition 14.1.1**  $\{y_t\}$  is **covariance (weakly) stationary** if

$$\mathbb{E}(y_t) = \mu$$

is independent of  $t$ , and

$$\text{cov}(y_t, y_{t-k}) = \gamma(k)$$

is independent of  $t$  for all  $k$ .  $\gamma(k)$  is called the **autocovariance function**.

$$\rho(k) = \gamma(k)/\gamma(0) = \text{corr}(y_t, y_{t-k})$$

is the **autocorrelation function**.

**Definition 14.1.2**  $\{y_t\}$  is **strictly stationary** if the joint distribution of  $(y_t, \dots, y_{t-k})$  is independent of  $t$  for all  $k$ .

**Definition 14.1.3** A stationary time series is **ergodic** if  $\gamma(k) \rightarrow 0$  as  $k \rightarrow \infty$ .

The following two theorems are essential to the analysis of stationary time series. The proofs are rather difficult, however.

**Theorem 14.1.1** *If  $y_t$  is strictly stationary and ergodic and  $x_t = f(y_t, y_{t-1}, \dots)$  is a random variable, then  $x_t$  is strictly stationary and ergodic.*

**Theorem 14.1.2** (Ergodic Theorem). *If  $y_t$  is strictly stationary and ergodic and  $\mathbb{E}|y_t| < \infty$ , then as  $T \rightarrow \infty$ ,*

$$\frac{1}{T} \sum_{t=1}^T y_t \xrightarrow{p} \mathbb{E}(y_t).$$

This allows us to consistently estimate parameters using time-series moments:  
The sample mean:

$$\hat{\mu} = \frac{1}{T} \sum_{t=1}^T y_t$$

The sample autocovariance

$$\hat{\gamma}(k) = \frac{1}{T} \sum_{t=1}^T (y_t - \hat{\mu})(y_{t-k} - \hat{\mu}).$$

The sample autocorrelation

$$\hat{\rho}(k) = \frac{\hat{\gamma}(k)}{\hat{\gamma}(0)}.$$

**Theorem 14.1.3** *If  $y_t$  is strictly stationary and ergodic and  $\mathbb{E}(y_t^2) < \infty$ , then as  $T \rightarrow \infty$ ,*

1.  $\hat{\mu} \xrightarrow{p} \mathbb{E}(y_t)$ ;
2.  $\hat{\gamma}(k) \xrightarrow{p} \gamma(k)$ ;
3.  $\hat{\rho}(k) \xrightarrow{p} \rho(k)$ .

---

**Proof of Theorem 14.1.3.** Part (1) is a direct consequence of the Ergodic theorem. For Part (2), note that

$$\begin{aligned} \hat{\gamma}(k) &= \frac{1}{T} \sum_{t=1}^T (y_t - \hat{\mu})(y_{t-k} - \hat{\mu}) \\ &= \frac{1}{T} \sum_{t=1}^T y_t y_{t-k} - \frac{1}{T} \sum_{t=1}^T y_t \hat{\mu} - \frac{1}{T} \sum_{t=1}^T y_{t-k} \hat{\mu} + \hat{\mu}^2. \end{aligned}$$

By Theorem 14.1.1 above, the sequence  $y_t y_{t-k}$  is strictly stationary and ergodic, and it has a finite mean by the assumption that  $\mathbb{E}(y_t^2) < \infty$ . Thus an application of the Ergodic Theorem yields

$$\frac{1}{T} \sum_{t=1}^T y_t y_{t-k} \xrightarrow{p} \mathbb{E}(y_t y_{t-k}).$$

Thus

$$\hat{\gamma}(k) \xrightarrow{p} \mathbb{E}(y_t y_{t-k}) - \mu^2 - \mu^2 + \mu^2 = \mathbb{E}(y_t y_{t-k}) - \mu^2 = \gamma(k).$$

Part (3) follows by the continuous mapping theorem:  $\hat{\rho}(k) = \hat{\gamma}(k)/\hat{\gamma}(0) \xrightarrow{p} \gamma(k)/\gamma(0) = \rho(k)$ .

---

## 14.2 Autoregressions

In time-series, the series  $\{..., y_1, y_2, ..., y_T, ...\}$  are jointly random. We consider the conditional expectation

$$\mathbb{E}(y_t \mid \mathcal{F}_{t-1})$$

where  $\mathcal{F}_{t-1} = \{y_{t-1}, y_{t-2}, \dots\}$  is the past history of the series.

An autoregressive (AR) model specifies that only a finite number of past lags matter:

$$\mathbb{E}(y_t \mid \mathcal{F}_{t-1}) = \mathbb{E}(y_t \mid y_{t-1}, \dots, y_{t-k}).$$

A linear AR model (the most common type used in practice) specifies linearity:

$$\mathbb{E}(y_t \mid \mathcal{F}_{t-1}) = \alpha_0 + \alpha_1 y_{t-1} + \alpha_2 y_{t-2} + \dots + \alpha_k y_{t-k}.$$

Letting

$$e_t = y_t - \mathbb{E}(y_t \mid \mathcal{F}_{t-1}),$$

then we have the autoregressive model

$$y_t = \alpha_0 + \alpha_1 y_{t-1} + \alpha_2 y_{t-2} + \dots + \alpha_k y_{t-k} + e_t$$

$$\mathbb{E}(e_t \mid \mathcal{F}_{t-1}) = 0.$$

The last property defines a special time-series process.

**Definition 14.2.1**  $e_t$  is a *martingale difference sequence (MDS)* if  $\mathbb{E}(e_t \mid \mathcal{F}_{t-1}) = 0$ .

Regression errors are naturally a MDS. Some time-series processes may be a MDS as a consequence of optimizing behavior. For example, some versions of the life-cycle hypothesis imply that either changes in consumption, or consumption growth rates, should be a MDS. Most asset pricing models imply that asset returns should be the sum of a constant plus a MDS.

The MDS property for the regression error plays the same role in a time-series regression as does the conditional mean-zero property for the regression error in a cross-section regression. In fact, it is even more important in the time-series context, as it is difficult to derive distribution theories without this property.

A useful property of a MDS is that  $e_t$  is uncorrelated with any function of the lagged information  $\mathcal{F}_{t-1}$ . Thus for  $k > 0$ ,  $\mathbb{E}(y_{t-k} e_t) = 0$ .

### 14.3 Stationarity of AR(1) Process

A mean-zero AR(1) is

$$y_t = \alpha y_{t-1} + e_t.$$

Assume that  $e_t$  is iid,  $\mathbb{E}(e_t) = 0$  and  $\mathbb{E}(e_t^2) = \sigma^2 < \infty$ .

By back-substitution, we find

$$\begin{aligned} y_t &= e_t + \alpha e_{t-1} + \alpha^2 e_{t-2} + \dots \\ &= \sum_{k=0}^{\infty} \alpha^k e_{t-k}. \end{aligned}$$

Loosely speaking, this series converges if the sequence  $\alpha^k e_{t-k}$  gets small as  $k \rightarrow \infty$ . This occurs when  $|\alpha| < 1$ .

**Theorem 14.3.1** *If and only if  $|\alpha| < 1$  then  $y_t$  is strictly stationary and ergodic.*

We can compute the moments of  $y_t$  using the infinite sum:

$$\begin{aligned} \mathbb{E}(y_t) &= \sum_{k=0}^{\infty} \alpha^k \mathbb{E}(e_{t-k}) = 0 \\ \text{var}(y_t) &= \sum_{k=0}^{\infty} \alpha^{2k} \text{var}(e_{t-k}) = \frac{\sigma^2}{1 - \alpha^2}. \end{aligned}$$

If the equation for  $y_t$  has an intercept, the above results are unchanged, except that the mean of  $y_t$  can be computed from the relationship

$$\mathbb{E}(y_t) = \alpha_0 + \alpha_1 \mathbb{E}(y_{t-1}),$$

and solving for  $\mathbb{E}(y_t) = \mathbb{E}(y_{t-1})$  we find  $\mathbb{E}(y_t) = \alpha_0 / (1 - \alpha_1)$ .

### 14.4 Lag Operator

An algebraic construct which is useful for the analysis of autoregressive models is the lag operator.

**Definition 14.4.1** *The **lag operator**  $L$  satisfies  $Ly_t = y_{t-1}$ .*

Defining  $L^2 = LL$ , we see that  $L^2 y_t = Ly_{t-1} = y_{t-2}$ . In general,  $L^k y_t = y_{t-k}$ .

The AR(1) model can be written in the format

$$y_t - \alpha y_{t-1} = e_t$$

or

$$(1 - \alpha L) y_t = e_t.$$

The operator  $\alpha(L) = (1 - \alpha L)$  is a polynomial in the operator  $L$ . We say that the *root* of the polynomial is  $1/\alpha$ , since  $\rho(z) = 0$  when  $z = 1/\alpha$ . We call  $\alpha(L)$  the autoregressive polynomial of  $y_t$ .

From Theorem 14.3.1, an AR(1) is stationary iff  $|\alpha| < 1$ . Note that an equivalent way to say this is that an AR(1) is stationary iff the root of the autoregressive polynomial is larger than one (in absolute value).

## 14.5 Stationarity of AR(k)

The AR(k) model is

$$y_t = \alpha_1 y_{t-1} + \alpha_2 y_{t-2} + \cdots + \alpha_k y_{t-k} + e_t.$$

Using the lag operator,

$$y_t - \alpha_1 L y_t - \alpha_2 L^2 y_t - \cdots - \alpha_k L^k y_t = e_t,$$

or

$$\alpha(L) y_t = e_t$$

where

$$\rho(L) = 1 - \alpha_1 L - \alpha_2 L^2 - \cdots - \alpha_k L^k.$$

We call  $\alpha(L)$  the autoregressive polynomial of  $y_t$ .

The *Fundamental Theorem of Algebra* says that any polynomial can be factored as

$$\alpha(z) = (1 - \lambda_1^{-1} z) (1 - \lambda_2^{-1} z) \cdots (1 - \lambda_k^{-1} z)$$

where the  $\lambda_1, \dots, \lambda_k$  are the complex *roots* of  $\alpha(z)$ , which satisfy  $\alpha(\lambda_j) = 0$ .

We know that an AR(1) is stationary iff the absolute value of the root of its autoregressive polynomial is larger than one. For an AR(k), the requirement is that all roots are larger than one. Let  $|\lambda|$  denote the modulus of a complex number  $\lambda$ .

**Theorem 14.5.1** *The AR(k) is strictly stationary and ergodic if and only if  $|\lambda_j| > 1$  for all  $j$ .*

One way of stating this is that “All roots lie outside the unit circle.”

If one of the roots equals 1, we say that  $\alpha(L)$ , and hence  $y_t$ , “has a unit root”. This is a special case of non-stationarity, and is of great interest in applied time series.

## 14.6 Estimation

Let

$$\begin{aligned} \mathbf{x}_t &= \begin{pmatrix} 1 & y_{t-1} & y_{t-2} & \cdots & y_{t-k} \end{pmatrix}' \\ \boldsymbol{\beta} &= \begin{pmatrix} \alpha_0 & \alpha_1 & \alpha_2 & \cdots & \alpha_k \end{pmatrix}'. \end{aligned}$$

Then the model can be written as

$$y_t = \mathbf{x}_t' \boldsymbol{\beta} + e_t.$$

The OLS estimator is

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}.$$

To study  $\hat{\boldsymbol{\beta}}$ , it is helpful to define the process  $u_t = \mathbf{x}_t e_t$ . Note that  $u_t$  is a MDS, since

$$\mathbb{E}(u_t | \mathcal{F}_{t-1}) = \mathbb{E}(\mathbf{x}_t e_t | \mathcal{F}_{t-1}) = \mathbf{x}_t \mathbb{E}(e_t | \mathcal{F}_{t-1}) = 0.$$

By Theorem 14.1.1, it is also strictly stationary and ergodic. Thus

$$\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t e_t = \frac{1}{T} \sum_{t=1}^T u_t \xrightarrow{p} \mathbb{E}(u_t) = 0. \quad (14.1)$$

The vector  $\mathbf{x}_t$  is strictly stationary and ergodic, and by Theorem 14.1.1, so is  $\mathbf{x}_t \mathbf{x}'_t$ . Thus by the Ergodic Theorem,

$$\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \mathbf{x}'_t \xrightarrow{p} \mathbb{E}(\mathbf{x}_t \mathbf{x}'_t) = \mathbf{Q}.$$

Combined with (14.1) and the continuous mapping theorem, we see that

$$\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} = \left( \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \mathbf{x}'_t \right)^{-1} \left( \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t e_t \right) \xrightarrow{p} \mathbf{Q}^{-1} \mathbf{0} = \mathbf{0}.$$

We have shown the following:

**Theorem 14.6.1** *If the AR(k) process  $y_t$  is strictly stationary and ergodic and  $\mathbb{E}(y_t^2) < \infty$ , then  $\hat{\boldsymbol{\beta}} \xrightarrow{p} \boldsymbol{\beta}$  as  $T \rightarrow \infty$ .*

## 14.7 Asymptotic Distribution

**Theorem 14.7.1** *MDS CLT. If  $\mathbf{u}_t$  is a strictly stationary and ergodic MDS and  $\mathbb{E}(\mathbf{u}_t \mathbf{u}'_t) = \boldsymbol{\Omega} < \infty$ , then as  $T \rightarrow \infty$ ,*

$$\frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{u}_t \xrightarrow{d} \mathbf{N}(\mathbf{0}, \boldsymbol{\Omega}).$$

Since  $\mathbf{x}_t e_t$  is a MDS, we can apply Theorem 14.7.1 to see that

$$\frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{x}_t e_t \xrightarrow{d} \mathbf{N}(\mathbf{0}, \boldsymbol{\Omega}),$$

where

$$\boldsymbol{\Omega} = \mathbb{E}(\mathbf{x}_t \mathbf{x}'_t e_t^2).$$

**Theorem 14.7.2** *If the AR(k) process  $y_t$  is strictly stationary and ergodic and  $\mathbb{E}(y_t^4) < \infty$ , then as  $T \rightarrow \infty$ ,*

$$\sqrt{T}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{Q}^{-1} \boldsymbol{\Omega} \mathbf{Q}^{-1}).$$

This is identical in form to the asymptotic distribution of OLS in cross-section regression. The implication is that asymptotic inference is the same. In particular, the asymptotic covariance matrix is estimated just as in the cross-section case.



## 14.8 Bootstrap for Autoregressions

In the non-parametric bootstrap, we constructed the bootstrap sample by randomly resampling from the data values  $\{y_t, \mathbf{x}_t\}$ . This creates an iid bootstrap sample. Clearly, this cannot work in a time-series application, as this imposes inappropriate independence.

Briefly, there are two popular methods to implement bootstrap resampling for time-series data.

### Method 1: Model-Based (Parametric) Bootstrap.

1. Estimate  $\hat{\beta}$  and residuals  $\hat{e}_t$ .
2. Fix an initial condition  $(y_{-k+1}, y_{-k+2}, \dots, y_0)$ .
3. Simulate iid draws  $e_i^*$  from the empirical distribution of the residuals  $\{\hat{e}_1, \dots, \hat{e}_T\}$ .
4. Create the bootstrap series  $y_t^*$  by the recursive formula

$$y_t^* = \hat{\alpha}_0 + \hat{\alpha}_1 y_{t-1}^* + \hat{\alpha}_2 y_{t-2}^* + \dots + \hat{\alpha}_k y_{t-k}^* + e_t^*.$$

This construction imposes homoskedasticity on the errors  $e_i^*$ , which may be different than the properties of the actual  $e_i$ . It also presumes that the AR(k) structure is the truth.

### Method 2: Block Resampling

1. Divide the sample into  $T/m$  blocks of length  $m$ .
2. Resample complete blocks. For each simulated sample, draw  $T/m$  blocks.
3. Paste the blocks together to create the bootstrap time-series  $y_t^*$ .
4. This allows for arbitrary stationary serial correlation, heteroskedasticity, and for model-misspecification.
5. The results may be sensitive to the block length, and the way that the data are partitioned into blocks.
6. May not work well in small samples.

## 14.9 Trend Stationarity

$$y_t = \mu_0 + \mu_1 t + S_t \tag{14.2}$$

$$S_t = \rho_1 S_{t-1} + \rho_2 S_{t-2} + \dots + \rho_k S_{t-k} + e_t, \tag{14.3}$$

or

$$y_t = \alpha_0 + \alpha_1 t + \rho_1 y_{t-1} + \rho_2 y_{t-2} + \dots + \rho_k y_{t-k} + e_t. \tag{14.4}$$

There are two essentially equivalent ways to estimate the autoregressive parameters  $(\alpha_1, \dots, \alpha_k)$ .

- You can estimate (14.4) by OLS.
- You can estimate (14.2)-(14.3) sequentially by OLS. That is, first estimate (14.2), get the residual  $\hat{S}_t$ , and then perform regression (14.3) replacing  $S_t$  with  $\hat{S}_t$ . This procedure is sometimes called *Detrending*.

The reason why these two procedures are (essentially) the same is the Frisch-Waugh-Lovell theorem.

### Seasonal Effects

There are three popular methods to deal with seasonal data.

- Include dummy variables for each season. This presumes that “seasonality” does not change over the sample.
- Use “seasonally adjusted” data. The seasonal factor is typically estimated by a two-sided weighted average of the data for that season in neighboring years. Thus the seasonally adjusted data is a “filtered” series. This is a flexible approach which can extract a wide range of seasonal factors. The seasonal adjustment, however, also alters the time-series correlations of the data.
- First apply a seasonal differencing operator. If  $s$  is the number of seasons (typically  $s = 4$  or  $s = 12$ ),

$$\Delta_s y_t = y_t - y_{t-s},$$

or the season-to-season change. The series  $\Delta_s y_t$  is clearly free of seasonality. But the long-run trend is also eliminated, and perhaps this was of relevance.

## 14.10 Testing for Omitted Serial Correlation

For simplicity, let the null hypothesis be an AR(1):

$$y_t = \alpha_0 + \alpha_1 y_{t-1} + u_t. \quad (14.5)$$

We are interested in the question if the error  $u_t$  is serially correlated. We model this as an AR(1):

$$u_t = \theta u_{t-1} + e_t \quad (14.6)$$

with  $e_t$  a MDS. The hypothesis of no omitted serial correlation is

$$\mathbb{H}_0 : \theta = 0$$

$$\mathbb{H}_1 : \theta \neq 0.$$

We want to test  $\mathbb{H}_0$  against  $\mathbb{H}_1$ .

To combine (14.5) and (14.6), we take (14.5) and lag the equation once:

$$y_{t-1} = \alpha_0 + \alpha_1 y_{t-2} + u_{t-1}.$$

We then multiply this by  $\theta$  and subtract from (14.5), to find

$$y_t - \theta y_{t-1} = \alpha_0 - \theta \alpha_0 + \alpha_1 y_{t-1} - \theta \alpha_1 y_{t-1} + u_t - \theta u_{t-1},$$

or

$$y_t = \alpha_0(1 - \theta) + (\alpha_1 + \theta) y_{t-1} - \theta \alpha_1 y_{t-2} + e_t = AR(2).$$

Thus under  $\mathbb{H}_0$ ,  $y_t$  is an AR(1), and under  $\mathbb{H}_1$  it is an AR(2).  $\mathbb{H}_0$  may be expressed as the restriction that the coefficient on  $y_{t-2}$  is zero.

An appropriate test of  $\mathbb{H}_0$  against  $\mathbb{H}_1$  is therefore a Wald test that the coefficient on  $y_{t-2}$  is zero. (A simple exclusion test).

In general, if the null hypothesis is that  $y_t$  is an AR(k), and the alternative is that the error is an AR(m), this is the same as saying that under the alternative  $y_t$  is an AR(k+m), and this is equivalent to the restriction that the coefficients on  $y_{t-k-1}, \dots, y_{t-k-m}$  are jointly zero. An appropriate test is the Wald test of this restriction.

## 14.11 Model Selection

What is the appropriate choice of  $k$  in practice? This is a problem of model selection.

A good choice is to minimize the AIC information criterion

$$AIC(k) = \log \hat{\sigma}^2(k) + \frac{2k}{T},$$

where  $\hat{\sigma}^2(k)$  is the estimated residual variance from an AR(k)

One ambiguity in defining the AIC criterion is that the sample available for estimation changes as  $k$  changes. (If you increase  $k$ , you need more initial conditions.) This can induce strange behavior in the AIC. The appropriate remedy is to fix a upper value  $\bar{k}$ , and then reserve the first  $\bar{k}$  as initial conditions, and then estimate the models AR(1), AR(2), ..., AR( $\bar{k}$ ) on this (unified) sample.

## 14.12 Autoregressive Unit Roots

The AR(k) model is

$$\begin{aligned}\alpha(L)y_t &= \alpha_0 + e_t \\ \alpha(L) &= 1 - \alpha_1 L - \dots - \alpha_k L^k.\end{aligned}$$

As we discussed before,  $y_t$  has a unit root when  $\alpha(1) = 0$ , or

$$\alpha_1 + \alpha_2 + \dots + \alpha_k = 1.$$

In this case,  $y_t$  is non-stationary. The ergodic theorem and MDS CLT do not apply, and test statistics are asymptotically non-normal.

A helpful way to write the equation is the so-called Dickey-Fuller reparameterization:

$$\Delta y_t = \rho_0 y_{t-1} + \rho_1 \Delta y_{t-1} + \dots + \rho_{k-1} \Delta y_{t-(k-1)} + e_t. \quad (14.7)$$

These models are equivalent linear transformations of one another. The DF parameterization is convenient because the parameter  $\rho_0$  summarizes the information about the unit root, since  $\alpha(1) = -\rho_0$ . To see this, observe that the lag polynomial for the  $y_t$  computed from (14.7) is

$$(1 - L) - \rho_0 L - \rho_1 (L - L^2) - \dots - \rho_{k-1} (L^{k-1} - L^k)$$

But this must equal  $\rho(L)$ , as the models are equivalent. Thus

$$\alpha(1) = (1 - 1) - \rho_0 - (1 - 1) - \dots - (1 - 1) = -\rho_0.$$

Hence, the hypothesis of a unit root in  $y_t$  can be stated as

$$\mathbb{H}_0 : \rho_0 = 0.$$

Note that the model is stationary if  $\rho_0 < 0$ . So the natural alternative is

$$\mathbb{H}_1 : \rho_0 < 0.$$

Under  $\mathbb{H}_0$ , the model for  $y_t$  is

$$\Delta y_t = \mu + \rho_1 \Delta y_{t-1} + \dots + \rho_{k-1} \Delta y_{t-(k-1)} + e_t,$$

which is an AR(k-1) in the first-difference  $\Delta y_t$ . Thus if  $y_t$  has a (single) unit root, then  $\Delta y_t$  is a stationary AR process. Because of this property, we say that if  $y_t$  is non-stationary but  $\Delta^d y_t$  is stationary, then  $y_t$  is “integrated of order  $d$ ”, or  $I(d)$ . Thus a time series with unit root is  $I(1)$ .

Since  $\alpha_0$  is the parameter of a linear regression, the natural test statistic is the t-statistic for  $\mathbb{H}_0$  from OLS estimation of (14.7). Indeed, this is the most popular unit root test, and is called the Augmented Dickey-Fuller (ADF) test for a unit root.

It would seem natural to assess the significance of the ADF statistic using the normal table. However, under  $\mathbb{H}_0$ ,  $y_t$  is non-stationary, so conventional normal asymptotics are invalid. An alternative asymptotic framework has been developed to deal with non-stationary data. We do not have the time to develop this theory in detail, but simply assert the main results.

**Theorem 14.12.1 Dickey-Fuller Theorem.**

If  $\rho_0 = 0$  then as  $T \rightarrow \infty$ ,

$$T\hat{\rho}_0 \xrightarrow{d} (1 - \rho_1 - \rho_2 - \cdots - \rho_{k-1}) DF_\alpha$$

$$ADF = \frac{\hat{\rho}_0}{s(\hat{\rho}_0)} \rightarrow DF_t.$$

The limit distributions  $DF_\alpha$  and  $DF_t$  are non-normal. They are skewed to the left, and have negative means.

The first result states that  $\hat{\rho}_0$  converges to its true value (of zero) at rate  $T$ , rather than the conventional rate of  $T^{1/2}$ . This is called a “super-consistent” rate of convergence.

The second result states that the t-statistic for  $\hat{\rho}_0$  converges to a limit distribution which is non-normal, but does not depend on the parameters  $\rho$ . This distribution has been extensively tabulated, and may be used for testing the hypothesis  $\mathbb{H}_0$ . Note: The standard error  $s(\hat{\rho}_0)$  is the conventional (“homoskedastic”) standard error. But the theorem does not require an assumption of homoskedasticity. Thus the Dickey-Fuller test is robust to heteroskedasticity.

Since the alternative hypothesis is one-sided, the ADF test rejects  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  when  $ADF < c$ , where  $c$  is the critical value from the ADF table. If the test rejects  $\mathbb{H}_0$ , this means that the evidence points to  $y_t$  being stationary. If the test does not reject  $\mathbb{H}_0$ , a common conclusion is that the data suggests that  $y_t$  is non-stationary. This is not really a correct conclusion, however. All we can say is that there is insufficient evidence to conclude whether the data are stationary or not.

We have described the test for the setting of with an intercept. Another popular setting includes as well a linear time trend. This model is

$$\Delta y_t = \mu_1 + \mu_2 t + \rho_0 y_{t-1} + \rho_1 \Delta y_{t-1} + \cdots + \rho_{k-1} \Delta y_{t-(k-1)} + e_t. \quad (14.8)$$

This is natural when the alternative hypothesis is that the series is stationary about a linear time trend. If the series has a linear trend (e.g. GDP, Stock Prices), then the series itself is non-stationary, but it may be stationary around the linear time trend. In this context, it is a silly waste of time to fit an AR model to the level of the series without a time trend, as the AR model cannot conceivably describe this data. The natural solution is to include a time trend in the fitted OLS equation. When conducting the ADF test, this means that it is computed as the t-ratio for  $\rho_0$  from OLS estimation of (14.8).

If a time trend is included, the test procedure is the same, but different critical values are required. The ADF test has a different distribution when the time trend has been included, and a different table should be consulted.

Most texts include as well the critical values for the extreme polar case where the intercept has been omitted from the model. These are included for completeness (from a pedagogical perspective) but have no relevance for empirical practice where intercepts are always included.

## Chapter 15

# Multivariate Time Series

A multivariate time series  $\mathbf{y}_t$  is a vector process  $m \times 1$ . Let  $\mathcal{F}_{t-1} = (\mathbf{y}_{t-1}, \mathbf{y}_{t-2}, \dots)$  be all lagged information at time  $t$ . The typical goal is to find the conditional expectation  $\mathbb{E}(\mathbf{y}_t | \mathcal{F}_{t-1})$ . Note that since  $\mathbf{y}_t$  is a vector, this conditional expectation is also a vector.

### 15.1 Vector Autoregressions (VARs)

A VAR model specifies that the conditional mean is a function of only a finite number of lags:

$$\mathbb{E}(\mathbf{y}_t | \mathcal{F}_{t-1}) = \mathbb{E}(\mathbf{y}_t | \mathbf{y}_{t-1}, \dots, \mathbf{y}_{t-k}).$$

A linear VAR specifies that this conditional mean is linear in the arguments:

$$\mathbb{E}(\mathbf{y}_t | \mathbf{y}_{t-1}, \dots, \mathbf{y}_{t-k}) = \mathbf{a}_0 + \mathbf{A}_1 \mathbf{y}_{t-1} + \mathbf{A}_2 \mathbf{y}_{t-2} + \dots + \mathbf{A}_k \mathbf{y}_{t-k}.$$

Observe that  $\mathbf{a}_0$  is  $m \times 1$ , and each of  $\mathbf{A}_1$  through  $\mathbf{A}_k$  are  $m \times m$  matrices.

Defining the  $m \times 1$  regression error

$$\mathbf{e}_t = \mathbf{y}_t - \mathbb{E}(\mathbf{y}_t | \mathcal{F}_{t-1}),$$

we have the VAR model

$$\begin{aligned} \mathbf{y}_t &= \mathbf{a}_0 + \mathbf{A}_1 \mathbf{y}_{t-1} + \mathbf{A}_2 \mathbf{y}_{t-2} + \dots + \mathbf{A}_k \mathbf{y}_{t-k} + \mathbf{e}_t \\ \mathbb{E}(\mathbf{e}_t | \mathcal{F}_{t-1}) &= \mathbf{0}. \end{aligned}$$

Alternatively, defining the  $mk + 1$  vector

$$\mathbf{x}_t = \begin{pmatrix} 1 \\ \mathbf{y}_{t-1} \\ \mathbf{y}_{t-2} \\ \vdots \\ \mathbf{y}_{t-k} \end{pmatrix}$$

and the  $m \times (mk + 1)$  matrix

$$\mathbf{A} = \begin{pmatrix} \mathbf{a}_0 & \mathbf{A}_1 & \mathbf{A}_2 & \dots & \mathbf{A}_k \end{pmatrix},$$

then

$$\mathbf{y}_t = \mathbf{A} \mathbf{x}_t + \mathbf{e}_t.$$

The VAR model is a system of  $m$  equations. One way to write this is to let  $a'_j$  be the  $j$ th row of  $\mathbf{A}$ . Then the VAR system can be written as the equations

$$Y_{jt} = a'_j \mathbf{x}_t + e_{jt}.$$

Unrestricted VARs were introduced to econometrics by Sims (1980).

## 15.2 Estimation

Consider the moment conditions

$$\mathbb{E}(\mathbf{x}_t e_{jt}) = \mathbf{0},$$

$j = 1, \dots, m$ . These are implied by the VAR model, either as a regression, or as a linear projection.

The GMM estimator corresponding to these moment conditions is equation-by-equation OLS

$$\hat{\mathbf{a}}_j = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}_j.$$

An alternative way to compute this is as follows. Note that

$$\hat{\mathbf{a}}'_j = \mathbf{y}'_j \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1}.$$

And if we stack these to create the estimate  $\hat{\mathbf{A}}$ , we find

$$\begin{aligned} \hat{\mathbf{A}} &= \begin{pmatrix} \mathbf{y}'_1 \\ \mathbf{y}'_2 \\ \vdots \\ \mathbf{y}'_{m+1} \end{pmatrix} \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \\ &= \mathbf{Y}'\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1}, \end{aligned}$$

where

$$\mathbf{Y} = \begin{pmatrix} \mathbf{y}_1 & \mathbf{y}_2 & \cdots & \mathbf{y}_m \end{pmatrix}$$

the  $T \times m$  matrix of the stacked  $\mathbf{y}'_t$ .

This (system) estimator is known as the SUR (Seemingly Unrelated Regressions) estimator, and was originally derived by Zellner (1962)

## 15.3 Restricted VARs

The unrestricted VAR is a system of  $m$  equations, each with the same set of regressors. A restricted VAR imposes restrictions on the system. For example, some regressors may be excluded from some of the equations. Restrictions may be imposed on individual equations, or across equations. The GMM framework gives a convenient method to impose such restrictions on estimation.

## 15.4 Single Equation from a VAR

Often, we are only interested in a single equation out of a VAR system. This takes the form

$$y_{jt} = \mathbf{a}'_j \mathbf{x}_t + e_t,$$

and  $\mathbf{x}_t$  consists of lagged values of  $y_{jt}$  and the other  $y'_{lt}$ s. In this case, it is convenient to re-define the variables. Let  $y_t = y_{jt}$ , and  $\mathbf{z}_t$  be the other variables. Let  $e_t = e_{jt}$  and  $\beta = \mathbf{a}_j$ . Then the single equation takes the form

$$y_t = \mathbf{x}'_t \beta + e_t, \tag{15.1}$$

and

$$\mathbf{x}_t = \begin{bmatrix} 1 & \mathbf{y}_{t-1} & \cdots & \mathbf{y}_{t-k} & \mathbf{z}'_{t-1} & \cdots & \mathbf{z}'_{t-k} \end{bmatrix}'.$$

This is just a conventional regression with time series data.

## 15.5 Testing for Omitted Serial Correlation

Consider the problem of testing for omitted serial correlation in equation (15.1). Suppose that  $e_t$  is an AR(1). Then

$$\begin{aligned} y_t &= \mathbf{x}'_t \boldsymbol{\beta} + e_t \\ e_t &= \theta e_{t-1} + u_t \\ \mathbb{E}(u_t \mid \mathcal{F}_{t-1}) &= 0. \end{aligned} \tag{15.2}$$

Then the null and alternative are

$$\mathbb{H}_0 : \theta = 0 \quad \mathbb{H}_1 : \theta \neq 0.$$

Take the equation  $y_t = \mathbf{x}'_t \boldsymbol{\beta} + e_t$ , and subtract off the equation once lagged multiplied by  $\theta$ , to get

$$\begin{aligned} y_t - \theta y_{t-1} &= (\mathbf{x}'_t \boldsymbol{\beta} + e_t) - \theta (\mathbf{x}'_{t-1} \boldsymbol{\beta} + e_{t-1}) \\ &= \mathbf{x}'_t \boldsymbol{\beta} - \theta \mathbf{x}'_{t-1} \boldsymbol{\beta} + e_t - \theta e_{t-1}, \end{aligned}$$

or

$$y_t = \theta y_{t-1} + \mathbf{x}'_t \boldsymbol{\beta} + \mathbf{x}'_{t-1} \boldsymbol{\gamma} + u_t, \tag{15.3}$$

which is a valid regression model.

So testing  $\mathbb{H}_0$  versus  $\mathbb{H}_1$  is equivalent to testing for the significance of adding  $(y_{t-1}, \mathbf{x}_{t-1})$  to the regression. This can be done by a Wald test. We see that an appropriate, general, and simple way to test for omitted serial correlation is to test the significance of extra lagged values of the dependent variable and regressors.

You may have heard of the Durbin-Watson test for omitted serial correlation, which once was very popular, and is still routinely reported by conventional regression packages. The DW test is appropriate only when regression  $y_t = \mathbf{x}'_t \boldsymbol{\beta} + e_t$  is not dynamic (has no lagged values on the RHS), and  $e_t$  is iid  $N(0, \sigma^2)$ . Otherwise it is invalid.

Another interesting fact is that (15.2) is a special case of (15.3), under the restriction  $\boldsymbol{\gamma} = -\boldsymbol{\beta}\theta$ . This restriction, which is called a common factor restriction, may be tested if desired. If valid, the model (15.2) may be estimated by iterated GLS. (A simple version of this estimator is called Cochrane-Orcutt.) Since the common factor restriction appears arbitrary, and is typically rejected empirically, direct estimation of (15.2) is uncommon in recent applications.

## 15.6 Selection of Lag Length in an VAR

If you want a data-dependent rule to pick the lag length  $k$  in a VAR, you may either use a testing-based approach (using, for example, the Wald statistic), or an information criterion approach. The formula for the AIC and BIC are

$$\begin{aligned} AIC(k) &= \log \det \left( \widehat{\boldsymbol{\Omega}}(k) \right) + 2 \frac{p}{T} \\ BIC(k) &= \log \det \left( \widehat{\boldsymbol{\Omega}}(k) \right) + \frac{p \log(T)}{T} \\ \widehat{\boldsymbol{\Omega}}(k) &= \frac{1}{T} \sum_{t=1}^T \widehat{\mathbf{e}}_t(k) \widehat{\mathbf{e}}_t(k)' \\ p &= m(km + 1) \end{aligned}$$

where  $p$  is the number of parameters in the model, and  $\widehat{\mathbf{e}}_t(k)$  is the OLS residual vector from the model with  $k$  lags. The log determinant is the criterion from the multivariate normal likelihood.

## 15.7 Granger Causality

Partition the data vector into  $(\mathbf{y}_t, \mathbf{z}_t)$ . Define the two information sets

$$\begin{aligned}\mathcal{F}_{1t} &= (\mathbf{y}_t, \mathbf{y}_{t-1}, \mathbf{y}_{t-2}, \dots) \\ \mathcal{F}_{2t} &= (\mathbf{y}_t, \mathbf{z}_t, \mathbf{y}_{t-1}, \mathbf{z}_{t-1}, \mathbf{y}_{t-2}, \mathbf{z}_{t-2}, \dots)\end{aligned}$$

The information set  $\mathcal{F}_{1t}$  is generated only by the history of  $\mathbf{y}_t$ , and the information set  $\mathcal{F}_{2t}$  is generated by both  $\mathbf{y}_t$  and  $\mathbf{z}_t$ . The latter has more information.

We say that  $\mathbf{z}_t$  does not **Granger-cause**  $\mathbf{y}_t$  if

$$\mathbb{E}(\mathbf{y}_t \mid \mathcal{F}_{1,t-1}) = \mathbb{E}(\mathbf{y}_t \mid \mathcal{F}_{2,t-1}).$$

That is, conditional on information in lagged  $\mathbf{y}_t$ , lagged  $\mathbf{z}_t$  does not help to forecast  $\mathbf{y}_t$ . If this condition does not hold, then we say that  $\mathbf{z}_t$  Granger-causes  $\mathbf{y}_t$ .

The reason why we call this “Granger Causality” rather than “causality” is because this is not a physical or structure definition of causality. If  $\mathbf{z}_t$  is some sort of forecast of the future, such as a futures price, then  $\mathbf{z}_t$  may help to forecast  $\mathbf{y}_t$  even though it does not “cause”  $\mathbf{y}_t$ . This definition of causality was developed by Granger (1969) and Sims (1972).

In a linear VAR, the equation for  $\mathbf{y}_t$  is

$$\mathbf{y}_t = \alpha + \rho_1 \mathbf{y}_{t-1} + \dots + \rho_k \mathbf{y}_{t-k} + \mathbf{z}'_{t-1} \boldsymbol{\gamma}_1 + \dots + \mathbf{z}'_{t-k} \boldsymbol{\gamma}_k + e_t.$$

In this equation,  $\mathbf{z}_t$  does not Granger-cause  $\mathbf{y}_t$  if and only if

$$\mathbb{H}_0 : \boldsymbol{\gamma}_1 = \boldsymbol{\gamma}_2 = \dots = \boldsymbol{\gamma}_k = 0.$$

This may be tested using an exclusion (Wald) test.

This idea can be applied to blocks of variables. That is,  $\mathbf{y}_t$  and/or  $\mathbf{z}_t$  can be vectors. The hypothesis can be tested by using the appropriate multivariate Wald test.

If it is found that  $\mathbf{z}_t$  does not Granger-cause  $\mathbf{y}_t$ , then we deduce that our time-series model of  $\mathbb{E}(\mathbf{y}_t \mid \mathcal{F}_{t-1})$  does not require the use of  $\mathbf{z}_t$ . Note, however, that  $\mathbf{z}_t$  may still be useful to explain other features of  $\mathbf{y}_t$ , such as the conditional variance.

### Clive W. J. Granger

Clive Granger (1934-2009) of England was one of the leading figures in time-series econometrics, and co-winner in 2003 of the Nobel Memorial Prize in Economic Sciences (along with Robert Engle). In addition to formalizing the definition of causality known as Granger causality, he invented the concept of cointegration, introduced spectral methods into econometrics, and formalized methods for the combination of forecasts.

## 15.8 Cointegration

The idea of cointegration is due to Granger (1981), and was articulated in detail by Engle and Granger (1987).



**Definition 15.8.1** The  $m \times 1$  series  $\mathbf{y}_t$  is *cointegrated* if  $\mathbf{y}_t$  is  $I(1)$  yet there exists  $\boldsymbol{\beta}$ ,  $m \times r$ , of rank  $r$ , such that  $\mathbf{z}_t = \boldsymbol{\beta}'\mathbf{y}_t$  is  $I(0)$ . The  $r$  vectors in  $\boldsymbol{\beta}$  are called the *cointegrating vectors*.

If the series  $\mathbf{y}_t$  is not cointegrated, then  $r = 0$ . If  $r = m$ , then  $\mathbf{y}_t$  is  $I(0)$ . For  $0 < r < m$ ,  $\mathbf{y}_t$  is  $I(1)$  and cointegrated.

In some cases, it may be believed that  $\boldsymbol{\beta}$  is known a priori. Often,  $\boldsymbol{\beta} = (1 \ -1)'$ . For example, if  $\mathbf{y}_t$  is a pair of interest rates, then  $\boldsymbol{\beta} = (1 \ -1)'$  specifies that the spread (the difference in returns) is stationary. If  $\mathbf{y} = (\log(C) \ \log(I))'$ , then  $\boldsymbol{\beta} = (1 \ -1)'$  specifies that  $\log(C/I)$  is stationary.

In other cases,  $\boldsymbol{\beta}$  may not be known.

If  $\mathbf{y}_t$  is cointegrated with a single cointegrating vector ( $r = 1$ ), then it turns out that  $\boldsymbol{\beta}$  can be consistently estimated by an OLS regression of one component of  $\mathbf{y}_t$  on the others. Thus  $\mathbf{y}_t = (Y_{1t}, Y_{2t})$  and  $\boldsymbol{\beta} = (\beta_1 \ \beta_2)$  and normalize  $\beta_1 = 1$ . Then  $\hat{\beta}_2 = (\mathbf{y}_2'\mathbf{y}_2)^{-1}\mathbf{y}_2'\mathbf{y}_1 \xrightarrow{p} \beta_2$ . Furthermore this estimator is super-consistent:  $T(\hat{\beta}_2 - \beta_2) = O_p(1)$ , as first shown by Stock (1987). While OLS is not, in general, a good method to estimate  $\boldsymbol{\beta}$ , it is useful in the construction of alternative estimators and tests.

We are often interested in testing the hypothesis of no cointegration:

$$\mathbb{H}_0 : r = 0$$

$$\mathbb{H}_1 : r > 0.$$

Suppose that  $\boldsymbol{\beta}$  is known, so  $\mathbf{z}_t = \boldsymbol{\beta}'\mathbf{y}_t$  is known. Then under  $\mathbb{H}_0$   $\mathbf{z}_t$  is  $I(1)$ , yet under  $\mathbb{H}_1$   $\mathbf{z}_t$  is  $I(0)$ . Thus  $\mathbb{H}_0$  can be tested using a univariate ADF test on  $\mathbf{z}_t$ .

When  $\boldsymbol{\beta}$  is unknown, Engle and Granger (1987) suggested using an ADF test on the estimated residual  $\hat{z}_t = \hat{\boldsymbol{\beta}}'\mathbf{y}_t$ , from OLS of  $y_{1t}$  on  $y_{2t}$ . Their justification was Stock's result that  $\hat{\boldsymbol{\beta}}$  is super-consistent under  $\mathbb{H}_1$ . Under  $\mathbb{H}_0$ , however,  $\hat{\boldsymbol{\beta}}$  is not consistent, so the ADF critical values are not appropriate. The asymptotic distribution was worked out by Phillips and Ouliaris (1990).

When the data have time trends, it may be necessary to include a time trend in the estimated cointegrating regression. Whether or not the time trend is included, the asymptotic distribution of the test is affected by the presence of the time trend. The asymptotic distribution was worked out in B. Hansen (1992).

## 15.9 Cointegrated VARs

We can write a VAR as

$$\begin{aligned} \mathbf{A}(L)\mathbf{y}_t &= \mathbf{e}_t \\ \mathbf{A}(L) &= \mathbf{I} - \mathbf{A}_1L - \mathbf{A}_2L^2 - \dots - \mathbf{A}_kL^k \end{aligned}$$

or alternatively as

$$\Delta\mathbf{y}_t = \boldsymbol{\Pi}\mathbf{y}_{t-1} + \mathbf{D}(L)\Delta\mathbf{y}_{t-1} + \mathbf{e}_t$$

where

$$\begin{aligned} \boldsymbol{\Pi} &= -\mathbf{A}(1) \\ &= -\mathbf{I} + \mathbf{A}_1 + \mathbf{A}_2 + \dots + \mathbf{A}_k. \end{aligned}$$

**Theorem 15.9.1 Granger Representation Theorem**

$\mathbf{y}_t$  is cointegrated with  $m \times r$   $\boldsymbol{\beta}$  if and only if  $\text{rank}(\boldsymbol{\Pi}) = r$  and  $\boldsymbol{\Pi} = \boldsymbol{\alpha}\boldsymbol{\beta}'$  where  $\boldsymbol{\alpha}$  is  $m \times r$ ,  $\text{rank}(\boldsymbol{\alpha}) = r$ .

Thus cointegration imposes a restriction upon the parameters of a VAR. The restricted model can be written as

$$\begin{aligned}\Delta \mathbf{y}_t &= \boldsymbol{\alpha}\boldsymbol{\beta}'\mathbf{y}_{t-1} + \mathbf{D}(\text{L})\Delta \mathbf{y}_{t-1} + \mathbf{e}_t \\ \Delta \mathbf{y}_t &= \boldsymbol{\alpha}\mathbf{z}_{t-1} + \mathbf{D}(\text{L})\Delta \mathbf{y}_{t-1} + \mathbf{e}_t.\end{aligned}$$

If  $\boldsymbol{\beta}$  is known, this can be estimated by OLS of  $\Delta \mathbf{y}_t$  on  $\mathbf{z}_{t-1}$  and the lags of  $\Delta \mathbf{y}_t$ .

If  $\boldsymbol{\beta}$  is unknown, then estimation is done by “reduced rank regression”, which is least-squares subject to the stated restriction. Equivalently, this is the MLE of the restricted parameters under the assumption that  $\mathbf{e}_t$  is iid  $N(\mathbf{0}, \boldsymbol{\Omega})$ .

One difficulty is that  $\boldsymbol{\beta}$  is not identified without normalization. When  $r = 1$ , we typically just normalize one element to equal unity. When  $r > 1$ , this does not work, and different authors have adopted different identification schemes.

In the context of a cointegrated VAR estimated by reduced rank regression, it is simple to test for cointegration by testing the rank of  $\boldsymbol{\Pi}$ . These tests are constructed as likelihood ratio (LR) tests. As they were discovered by Johansen (1988, 1991, 1995), they are typically called the “Johansen Max and Trace” tests. Their asymptotic distributions are non-standard, and are similar to the Dickey-Fuller distributions.

# Chapter 16

## Panel Data

A panel is a set of observations on individuals, collected over time. An observation is the pair  $\{y_{it}, \mathbf{x}_{it}\}$ , where the  $i$  subscript denotes the individual, and the  $t$  subscript denotes time. A panel may be **balanced**:

$$\{y_{it}, \mathbf{x}_{it}\} : t = 1, \dots, T; \quad i = 1, \dots, n,$$

or **unbalanced**:

$$\{y_{it}, \mathbf{x}_{it}\} : \text{For } i = 1, \dots, n, \quad t = \underline{t}_i, \dots, \bar{t}_i.$$

### 16.1 Individual-Effects Model

The standard panel data specification is that there is an individual-specific effect which enters linearly in the regression

$$y_{it} = \mathbf{x}_{it}'\boldsymbol{\beta} + u_i + e_{it}.$$

The typical maintained assumptions are that the individuals  $i$  are mutually independent, that  $u_i$  and  $e_{it}$  are independent, that  $e_{it}$  is iid across individuals and time, and that  $e_{it}$  is uncorrelated with  $\mathbf{x}_{it}$ .

OLS of  $y_{it}$  on  $\mathbf{x}_{it}$  is called **pooled estimation**. It is consistent if

$$\mathbb{E}(\mathbf{x}_{it}u_i) = 0 \tag{16.1}$$

If this condition fails, then OLS is inconsistent. (16.1) fails if the individual-specific unobserved effect  $u_i$  is correlated with the observed explanatory variables  $\mathbf{x}_{it}$ . This is often believed to be plausible if  $u_i$  is an omitted variable.

If (16.1) is true, however, OLS can be improved upon via a GLS technique. In either event, OLS appears a poor estimation choice.

Condition (16.1) is called the **random effects hypothesis**. It is a strong assumption, and most applied researchers try to avoid its use.

### 16.2 Fixed Effects

This is the most common technique for estimation of non-dynamic linear panel regressions.

The motivation is to allow  $u_i$  to be arbitrary, and have arbitrary correlated with  $\mathbf{x}_i$ . The goal is to eliminate  $u_i$  from the estimator, and thus achieve invariance.

There are several derivations of the estimator.

First, let

$$d_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{else} \end{cases},$$

and

$$\mathbf{d}_i = \begin{pmatrix} d_{i1} \\ \vdots \\ d_{in} \end{pmatrix},$$

an  $n \times 1$  dummy vector with a “1” in the  $i^{th}$  place. Let

$$\mathbf{u} = \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix}.$$

Then note that

$$u_i = \mathbf{d}_i' \mathbf{u},$$

and

$$y_{it} = \mathbf{x}_{it}' \boldsymbol{\beta} + \mathbf{d}_i' \mathbf{u} + e_{it}. \quad (16.2)$$

Observe that

$$\mathbb{E}(e_{it} \mid \mathbf{x}_{it}, \mathbf{d}_i) = 0,$$

so (16.2) is a valid regression, with  $\mathbf{d}_i$  as a regressor along with  $\mathbf{x}_i$ .

OLS on (16.2) yields estimator  $(\hat{\boldsymbol{\beta}}, \hat{\mathbf{u}})$ . Conventional inference applies.

Observe that

- This is generally consistent.
- If  $\mathbf{x}_{it}$  contains an intercept, it will be collinear with  $\mathbf{d}_i$ , so the intercept is typically omitted from  $\mathbf{x}_{it}$ .
- Any regressor in  $\mathbf{x}_{it}$  which is constant over time for all individuals (e.g., their gender) will be collinear with  $\mathbf{d}_i$ , so will have to be omitted.
- There are  $n + k$  regression parameters, which is quite large as typically  $n$  is very large.

Computationally, you do not want to actually implement conventional OLS estimation, as the parameter space is too large. OLS estimation of  $\boldsymbol{\beta}$  proceeds by the FWL theorem. Stacking the observations together:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{D}\mathbf{u} + \mathbf{e},$$

then by the FWL theorem,

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\mathbf{X}'(\mathbf{I} - \mathbf{P}_D)\mathbf{X})^{-1} (\mathbf{X}'(\mathbf{I} - \mathbf{P}_D)\mathbf{y}) \\ &= (\mathbf{X}^*{}'\mathbf{X}^*)^{-1} (\mathbf{X}^*{}'\mathbf{y}^*), \end{aligned}$$

where

$$\begin{aligned} \mathbf{y}^* &= \mathbf{y} - \mathbf{D}(\mathbf{D}'\mathbf{D})^{-1}\mathbf{D}'\mathbf{y} \\ \mathbf{X}^* &= \mathbf{X} - \mathbf{D}(\mathbf{D}'\mathbf{D})^{-1}\mathbf{D}'\mathbf{X}. \end{aligned}$$

Since the regression of  $y_{it}$  on  $\mathbf{d}_i$  is a regression onto individual-specific dummies, the predicted value from these regressions is the individual specific mean  $\bar{y}_i$ , and the residual is the demean value

$$y_{it}^* = y_{it} - \bar{y}_i.$$

The fixed effects estimator  $\hat{\boldsymbol{\beta}}$  is OLS of  $y_{it}^*$  on  $\mathbf{x}_{it}^*$ , the dependent variable and regressors in deviation-from-mean form.

Another derivation of the estimator is to take the equation

$$y_{it} = \mathbf{x}'_{it}\boldsymbol{\beta} + u_i + e_{it},$$

and then take individual-specific means by taking the average for the  $i^{th}$  individual:

$$\frac{1}{T_i} \sum_{t=\underline{t}_i}^{\bar{t}_i} y_{it} = \frac{1}{T_i} \sum_{t=\underline{t}_i}^{\bar{t}_i} \mathbf{x}'_{it}\boldsymbol{\beta} + u_i + \frac{1}{T_i} \sum_{t=\underline{t}_i}^{\bar{t}_i} e_{it}$$

or

$$\bar{y}_i = \bar{\mathbf{x}}'_i\boldsymbol{\beta} + u_i + \bar{e}_i.$$

Subtracting, we find

$$y_{it}^* = \mathbf{x}'_{it}\boldsymbol{\beta} + e_{it}^*,$$

which is free of the individual-effect  $u_i$ .

### 16.3 Dynamic Panel Regression

A dynamic panel regression has a lagged dependent variable

$$y_{it} = \alpha y_{it-1} + \mathbf{x}'_{it}\boldsymbol{\beta} + u_i + e_{it}. \quad (16.3)$$

This is a model suitable for studying dynamic behavior of individual agents.

Unfortunately, the fixed effects estimator is inconsistent, at least if  $T$  is held finite as  $n \rightarrow \infty$ . This is because the sample mean of  $y_{it-1}$  is correlated with that of  $e_{it}$ .

The standard approach to estimate a dynamic panel is to combine first-differencing with IV or GMM. Taking first-differences of (16.3) eliminates the individual-specific effect:

$$\Delta y_{it} = \alpha \Delta y_{it-1} + \Delta \mathbf{x}'_{it}\boldsymbol{\beta} + \Delta e_{it}. \quad (16.4)$$

However, if  $e_{it}$  is iid, then it will be correlated with  $\Delta y_{it-1}$ :

$$\mathbb{E}(\Delta y_{it-1} \Delta e_{it}) = \mathbb{E}((y_{it-1} - y_{it-2})(e_{it} - e_{it-1})) = -\mathbb{E}(y_{it-1} e_{it-1}) = -\sigma_e^2.$$

So OLS on (16.4) will be inconsistent.

But if there are valid instruments, then IV or GMM can be used to estimate the equation. Typically, we use lags of the dependent variable, two periods back, as  $y_{t-2}$  is uncorrelated with  $\Delta e_{it}$ . Thus values of  $y_{it-k}$ ,  $k \geq 2$ , are valid instruments.

Hence a valid estimator of  $\alpha$  and  $\boldsymbol{\beta}$  is to estimate (16.4) by IV using  $y_{t-2}$  as an instrument for  $\Delta y_{t-1}$  (which is just identified). Alternatively, GMM using  $y_{t-2}$  and  $y_{t-3}$  as instruments (which is overidentified, but loses a time-series observation).

A more sophisticated GMM estimator recognizes that for time-periods later in the sample, there are more instruments available, so the instrument list should be different for each equation. This is conveniently organized by the GMM principle, as this enables the moments from the different time-periods to be stacked together to create a list of all the moment conditions. A simple application of GMM yields the parameter estimates and standard errors.

## Exercises

**Exercise 16.1** Consider the model

$$y_{it} = \mathbf{x}'_{it}\boldsymbol{\beta} + u_i + e_{it}$$
$$\mathbb{E}(\mathbf{z}_{it}e_{it}) = 0$$

for  $i = 1, \dots, n$  and  $t = 1, \dots, T$ . The individual effect  $u_i$  is treated as fixed. Assume  $\mathbf{x}_{it}$  and  $\mathbf{z}_{it}$  are  $k \times 1$  vectors.

Write out an appropriate estimator for  $\boldsymbol{\beta}$ .

# Chapter 17

## NonParametric Regression

### 17.1 Introduction

When components of  $\mathbf{x}$  are continuously distributed then the conditional expectation function

$$\mathbb{E}(y_i \mid \mathbf{x}_i = \mathbf{x}) = m(\mathbf{x})$$

can take any nonlinear shape. Unless an economic model restricts the form of  $m(\mathbf{x})$  to a parametric function, the CEF is inherently **nonparametric**, meaning that the function  $m(\mathbf{x})$  is an element of an infinite dimensional class. In this situation, how can we estimate  $m(\mathbf{x})$ ? What is a suitable method, if we acknowledge that  $m(\mathbf{x})$  is nonparametric?

There are two main classes of nonparametric regression estimators: kernel estimators, and series estimators. In this chapter we introduce kernel methods.

To get started, suppose that there is a single real-valued regressor  $x_i$ . We consider the case of vector-valued regressors later.

### 17.2 Binned Estimator

For clarity, fix the point  $x$  and consider estimation of the single point  $m(x)$ . This is the mean of  $y_i$  for random pairs  $(y_i, x_i)$  such that  $x_i = x$ . If the distribution of  $x_i$  were discrete then we could estimate  $m(x)$  by taking the average of the sub-sample of observations  $y_i$  for which  $x_i = x$ . But when  $x_i$  is continuous then the probability is zero that  $x_i$  exactly equals any specific  $x$ . So there is no sub-sample of observations with  $x_i = x$  and we cannot simply take the average of the corresponding  $y_i$  values. However, if the CEF  $m(x)$  is continuous, then it should be possible to get a good approximation by taking the average of the observations for which  $x_i$  is *close* to  $x$ , perhaps for the observations for which  $|x_i - x| \leq h$  for some small  $h > 0$ . We call  $h$  a **bandwidth**. This estimator can be written as

$$\hat{m}(x) = \frac{\sum_{i=1}^n 1(|x_i - x| \leq h) y_i}{\sum_{i=1}^n 1(|x_i - x| \leq h)} \quad (17.1)$$

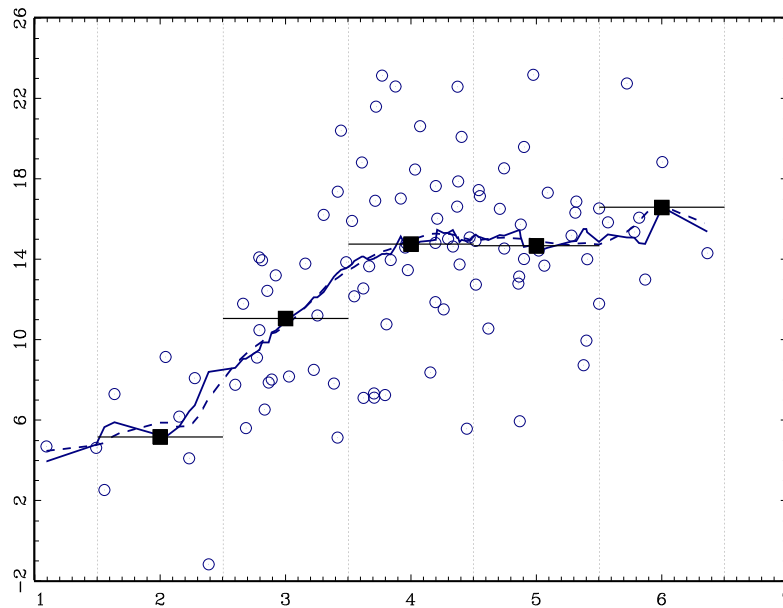
where  $1(\cdot)$  is the indicator function. Alternatively, (17.1) can be written as

$$\hat{m}(x) = \sum_{i=1}^n w_i(x) y_i \quad (17.2)$$

where

$$w_i(x) = \frac{1(|x_i - x| \leq h)}{\sum_{j=1}^n 1(|x_j - x| \leq h)}.$$

Notice that  $\sum_{i=1}^n w_i(x) = 1$ , so (17.2) is a weighted average of the  $y_i$ .

Figure 17.1: Scatter of  $(y_i, x_i)$  and Nadaraya-Watson regression

It is possible that for some values of  $x$  there are no values of  $x_i$  such that  $|x_i - x| \leq h$ , which implies that  $\sum_{i=1}^n 1(|x_i - x| \leq h) = 0$ . In this case the estimator (17.1) is undefined for those values of  $x$ .

To visualize, Figure 17.1 displays a scatter plot of 100 observations on a random pair  $(y_i, x_i)$  generated by simulation<sup>1</sup>. (The observations are displayed as the open circles.) The estimator (17.1) of the CEF  $m(x)$  at  $x = 2$  with  $h = 1/2$  is the average of the  $y_i$  for the observations such that  $x_i$  falls in the interval  $[1.5, 2.5]$ . (Our choice of  $h = 1/2$  is somewhat arbitrary. Selection of  $h$  will be discussed later.) The estimate is  $\hat{m}(2) = 5.16$  and is shown on Figure 17.1 by the first solid square. We repeat the calculation (17.1) for  $x = 3, 4, 5$ , and  $6$ , which is equivalent to partitioning the support of  $x_i$  into the regions  $[1.5, 2.5]$ ,  $[2.5, 3.5]$ ,  $[3.5, 4.5]$ ,  $[4.5, 5.5]$ , and  $[5.5, 6.5]$ . These partitions are shown in Figure 17.1 by the vertical dotted lines, and the estimates (17.1) by the solid squares.

These estimates  $\hat{m}(x)$  can be viewed as estimates of the CEF  $m(x)$ . Sometimes called a binned estimator, this is a step-function approximation to  $m(x)$  and is displayed in Figure 17.1 by the horizontal lines passing through the solid squares. This estimate roughly tracks the central tendency of the scatter of the observations  $(y_i, x_i)$ . However, the huge jumps in the estimated step function at the edges of the partitions are disconcerting, counter-intuitive, and clearly an artifact of the discrete binning.

If we take another look at the estimation formula (17.1) there is no reason why we need to evaluate (17.1) only on a coarse grid. We can evaluate  $\hat{m}(x)$  for any set of values of  $x$ . In particular, we can evaluate (17.1) on a fine grid of values of  $x$  and thereby obtain a smoother estimate of the CEF. This estimator with  $h = 1/2$  is displayed in Figure 17.1 with the solid line. This is a generalization of the binned estimator and by construction passes through the solid squares.

The bandwidth  $h$  determines the degree of smoothing. Larger values of  $h$  increase the width of the bins in Figure 17.1, thereby increasing the smoothness of the estimate  $\hat{m}(x)$  as a function of  $x$ . Smaller values of  $h$  decrease the width of the bins, resulting in less smooth conditional mean estimates.

<sup>1</sup>The distribution is  $x_i \sim N(4, 1)$  and  $y_i | x_i \sim N(m(x_i), 16)$  with  $m(x) = 10 \log(x)$ .



### 17.3 Kernel Regression

One deficiency with the estimator (17.1) is that it is a step function in  $x$ , as it is discontinuous at each observation  $x = x_i$ . That is why its plot in Figure 17.1 is jagged. The source of the discontinuity is that the weights  $w_i(x)$  are constructed from indicator functions, which are themselves discontinuous. If instead the weights are constructed from continuous functions then the CEF estimator will also be continuous in  $x$ .

To generalize (17.1) it is useful to write the weights  $1(|x_i - x| \leq h)$  in terms of the uniform density function on  $[-1, 1]$

$$k_0(u) = \frac{1}{2}1(|u| \leq 1).$$

Then

$$1(|x_i - x| \leq h) = 1\left(\left|\frac{x_i - x}{h}\right| \leq 1\right) = 2k_0\left(\frac{x_i - x}{h}\right).$$

and (17.1) can be written as

$$\hat{m}(x) = \frac{\sum_{i=1}^n k_0\left(\frac{x_i - x}{h}\right) y_i}{\sum_{i=1}^n k_0\left(\frac{x_i - x}{h}\right)}. \quad (17.3)$$

The uniform density  $k_0(u)$  is a special case of what is known as a **kernel function**.

**Definition 17.3.1** A second-order **kernel function**  $k(u)$  satisfies  $0 \leq k(u) < \infty$ ,  $k(u) = k(-u)$ ,  $\int_{-\infty}^{\infty} k(u) du = 1$  and  $\sigma_k^2 = \int_{-\infty}^{\infty} u^2 k(u) du < \infty$ .

Essentially, a kernel function is a probability density function which is bounded and symmetric about zero. A generalization of (17.1) is obtained by replacing the uniform kernel with any other kernel function:

$$\hat{m}(x) = \frac{\sum_{i=1}^n k\left(\frac{x_i - x}{h}\right) y_i}{\sum_{i=1}^n k\left(\frac{x_i - x}{h}\right)}. \quad (17.4)$$

The estimator (17.4) also takes the form (17.2) with

$$w_i(x) = \frac{k\left(\frac{x_i - x}{h}\right)}{\sum_{j=1}^n k\left(\frac{x_j - x}{h}\right)}.$$

The estimator (17.4) is known as the **Nadaraya-Watson** estimator, the **kernel regression** estimator, or the **local constant** estimator.

The bandwidth  $h$  plays the same role in (17.4) as it does in (17.1). Namely, larger values of  $h$  will result in estimates  $\hat{m}(x)$  which are smoother in  $x$ , and smaller values of  $h$  will result in estimates which are more erratic. It might be helpful to consider the two extreme cases  $h \rightarrow 0$  and  $h \rightarrow \infty$ . As  $h \rightarrow 0$  we can see that  $\hat{m}(x_i) \rightarrow y_i$  (if the values of  $x_i$  are unique), so that  $\hat{m}(x)$  is simply the scatter of  $y_i$  on  $x_i$ . In contrast, as  $h \rightarrow \infty$  then for all  $x$ ,  $\hat{m}(x) \rightarrow \bar{y}$ , the sample mean, so that the nonparametric CEF estimate is a constant function. For intermediate values of  $h$ ,  $\hat{m}(x)$  will lie between these two extreme cases.

The uniform density is not a good kernel choice as it produces discontinuous CEF estimates. To obtain a continuous CEF estimate  $\hat{m}(x)$  it is necessary for the kernel  $k(u)$  to be continuous. The two most commonly used choices are the **Epanechnikov kernel**

$$k_1(u) = \frac{3}{4} (1 - u^2) 1(|u| \leq 1)$$

and the **normal** or **Gaussian kernel**

$$k_\phi(u) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right).$$

For computation of the CEF estimate (17.4) the scale of the kernel is not important so long as the bandwidth is selected appropriately. That is, for any  $b > 0$ ,  $k_b(u) = b^{-1}k\left(\frac{u}{b}\right)$  is a valid kernel function with the identical shape as  $k(u)$ . Kernel regression with the kernel  $k(u)$  and bandwidth  $h$  is identical to kernel regression with the kernel  $k_b(u)$  and bandwidth  $h/b$ .

The estimate (17.4) using the Epanechnikov kernel and  $h = 1/2$  is also displayed in Figure 17.1 with the dashed line. As you can see, this estimator appears to be much smoother than that using the uniform kernel.

Two important constants associated with a kernel function  $k(u)$  are its variance  $\sigma_k^2$  and roughness  $R_k$ , which are defined as

$$\sigma_k^2 = \int_{-\infty}^{\infty} u^2 k(u) du \quad (17.5)$$

$$R_k = \int_{-\infty}^{\infty} k(u)^2 du. \quad (17.6)$$

Some common kernels and their roughness and variance values are reported in Table 9.1.

**Table 9.1: Common Second-Order Kernels**

Kernel	Equation	$R_k$	$\sigma_k^2$
Uniform	$k_0(u) = \frac{1}{2} 1( u  \leq 1)$	1/2	1/3
Epanechnikov	$k_1(u) = \frac{3}{4} (1 - u^2) 1( u  \leq 1)$	3/5	1/5
Biweight	$k_2(u) = \frac{15}{16} (1 - u^2)^2 1( u  \leq 1)$	5/7	1/7
Triweight	$k_3(u) = \frac{35}{32} (1 - u^2)^3 1( u  \leq 1)$	350/429	1/9
Gaussian	$k_\phi(u) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right)$	$1/(2\sqrt{\pi})$	1

## 17.4 Local Linear Estimator

The Nadaraya-Watson (NW) estimator is often called a local constant estimator as it locally (about  $x$ ) approximates the CEF  $m(x)$  as a constant function. One way to see this is to observe that  $\hat{m}(x)$  solves the minimization problem

$$\hat{m}(x) = \operatorname{argmin}_{\alpha} \sum_{i=1}^n k\left(\frac{x_i - x}{h}\right) (y_i - \alpha)^2.$$

This is a weighted regression of  $y_i$  on an intercept only. Without the weights, this estimation problem reduces to the sample mean. The NW estimator generalizes this to a local mean.

This interpretation suggests that we can construct alternative nonparametric estimators of the CEF by alternative local approximations. Many such local approximations are possible. A popular choice is the **Local Linear** (LL) approximation. Instead of approximating  $m(x)$  locally

as a constant, LL approximates the CEF locally by a linear function, and estimates this local approximation by locally weighted least squares.

Specifically, for each  $x$  we solve the following minimization problem

$$\{\hat{\alpha}(x), \hat{\beta}(x)\} = \underset{\alpha, \beta}{\operatorname{argmin}} \sum_{i=1}^n k\left(\frac{x_i - x}{h}\right) (y_i - \alpha - \beta(x_i - x))^2.$$

The local linear estimator of  $m(x)$  is the estimated intercept

$$\hat{m}(x) = \hat{\alpha}(x)$$

and the local linear estimator of the regression derivative  $\nabla m(x)$  is the estimated slope coefficient

$$\widehat{\nabla m}(x) = \hat{\beta}(x).$$

Computationally, for each  $x$  set

$$\mathbf{z}_i(x) = \begin{pmatrix} 1 \\ x_i - x \end{pmatrix}$$

and

$$k_i(x) = k\left(\frac{x_i - x}{h}\right).$$

Then

$$\begin{aligned} \begin{pmatrix} \hat{\alpha}(x) \\ \hat{\beta}(x) \end{pmatrix} &= \left( \sum_{i=1}^n k_i(x) \mathbf{z}_i(x) \mathbf{z}_i(x)' \right)^{-1} \sum_{i=1}^n k_i(x) \mathbf{z}_i(x) y_i \\ &= (\mathbf{Z}' \mathbf{K} \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{K} \mathbf{y} \end{aligned}$$

where  $\mathbf{K} = \operatorname{diag}\{k_1(x), \dots, k_n(x)\}$ .

To visualize, Figure 17.2 displays the scatter plot of the same 100 observations from Figure 17.1, divided into three regions depending on the regressor  $x_i$ :  $[1, 3]$ ,  $[3, 5]$ ,  $[5, 7]$ . A linear regression is fit to the observations in each region, with the observations weighted by the Epanechnikov kernel with  $h = 1$ . The three fitted regression lines are displayed by the three straight solid lines. The values of these regression lines at  $x = 2$ ,  $x = 4$  and  $x = 6$ , respectively, are the local linear estimates  $\hat{m}(x)$  at  $x = 2$ , 4, and 6. This estimation is repeated for all  $x$  in the support of the regressors, and plotted as the continuous solid line in Figure 17.2.

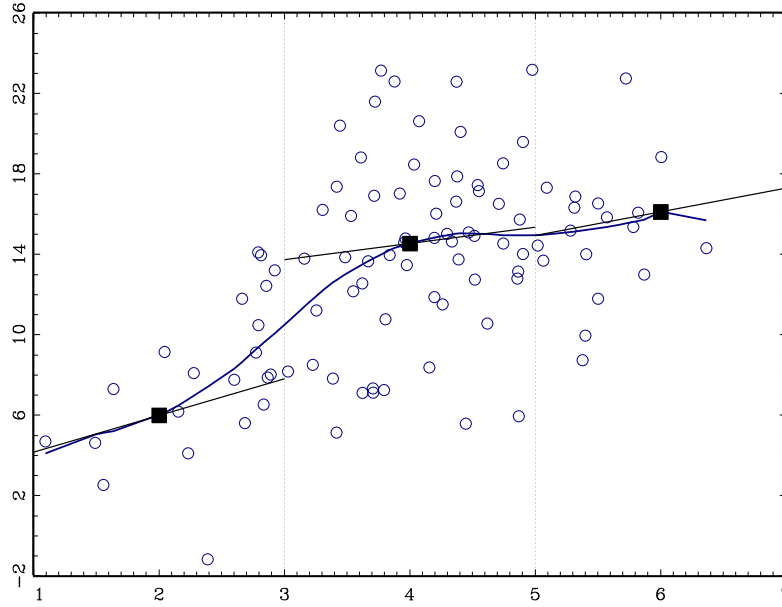
One interesting feature is that as  $h \rightarrow \infty$ , the LL estimator approaches the full-sample linear least-squares estimator  $\hat{m}(x) \rightarrow \hat{\alpha} + \hat{\beta}x$ . That is because as  $h \rightarrow \infty$  all observations receive equal weight regardless of  $x$ . In this sense we can see that the LL estimator is a flexible generalization of the linear OLS estimator.

Which nonparametric estimator should you use in practice: NW or LL? The theoretical literature shows that neither strictly dominates the other, but we can describe contexts where one or the other does better. Roughly speaking, the NW estimator performs better than the LL estimator when  $m(x)$  is close to a flat line, but the LL estimator performs better when  $m(x)$  is meaningfully non-constant. The LL estimator also performs better for values of  $x$  near the boundary of the support of  $x_i$ .

## 17.5 Nonparametric Residuals and Regression Fit

The fitted regression at  $x = x_i$  is  $\hat{m}(x_i)$  and the fitted residual is

$$\hat{e}_i = y_i - \hat{m}(x_i).$$

Figure 17.2: Scatter of  $(y_i, x_i)$  and Local Linear fitted regression

As a general rule, but especially when the bandwidth  $h$  is small, it is hard to view  $\hat{e}_i$  as a good measure of the fit of the regression. As  $h \rightarrow 0$  then  $\hat{m}(x_i) \rightarrow y_i$  and therefore  $\hat{e}_i \rightarrow 0$ . This clearly indicates overfitting as the true error is not zero. In general, since  $\hat{m}(x_i)$  is a local average which includes  $y_i$ , the fitted value will be necessarily close to  $y_i$  and the residual  $\hat{e}_i$  small, and the degree of this overfitting increases as  $h$  decreases.

A standard solution is to measure the fit of the regression at  $x = x_i$  by re-estimating the model excluding the  $i^{th}$  observation. For Nadaraya-Watson regression, the leave-one-out estimator of  $m(x)$  excluding observation  $i$  is

$$\tilde{m}_{-i}(x) = \frac{\sum_{j \neq i} k\left(\frac{x_j - x}{h}\right) y_j}{\sum_{j \neq i} k\left(\frac{x_j - x}{h}\right)}.$$

Notationally, the “ $-i$ ” subscript is used to indicate that the  $i^{th}$  observation is omitted.

The leave-one-out predicted value for  $y_i$  at  $x = x_i$  equals

$$\tilde{y}_i = \tilde{m}_{-i}(x_i) = \frac{\sum_{j \neq i} k\left(\frac{x_j - x_i}{h}\right) y_j}{\sum_{j \neq i} k\left(\frac{x_j - x_i}{h}\right)}.$$

The leave-one-out residuals (or prediction errors) are the difference between the leave-one-out predicted values and the actual observation

$$\tilde{e}_i = y_i - \tilde{y}_i.$$

Since  $\tilde{y}_i$  is not a function of  $y_i$ , there is no tendency for  $\tilde{y}_i$  to overfit for small  $h$ . Consequently,  $\tilde{e}_i$  is a good measure of the fit of the estimated nonparametric regression.

Similarly, the leave-one-out local-linear residual is  $\tilde{e}_i = y_i - \tilde{\alpha}_i$  with

$$\begin{pmatrix} \tilde{\alpha}_i \\ \tilde{\beta}_i \end{pmatrix} = \left( \sum_{j \neq i} k_{ij} \mathbf{z}_{ij} \mathbf{z}_{ij}' \right)^{-1} \sum_{j \neq i} k_{ij} \mathbf{z}_{ij} y_j,$$

$$\mathbf{z}_{ij} = \begin{pmatrix} 1 \\ x_j - x_i \end{pmatrix}$$

and

$$k_{ij} = k\left(\frac{x_j - x_i}{h}\right).$$

## 17.6 Cross-Validation Bandwidth Selection

As we mentioned before, the choice of bandwidth  $h$  is crucial. As  $h$  increases, the kernel regression estimators (both NW and LL) become more smooth, ironing out the bumps and wiggles. This reduces estimation variance but at the cost of increased bias and oversmoothing. As  $h$  decreases the estimators become more wiggly, erratic, and noisy. It is desirable to select  $h$  to trade-off these features. How can this be done systematically?

To be explicit about the dependence of the estimator on the bandwidth, let us write the estimator of  $m(x)$  with a given bandwidth  $h$  as  $\hat{m}(x, h)$ , and our discussion will apply equally to the NW and LL estimators.

Ideally, we would like to select  $h$  to minimize the mean-squared error (MSE) of  $\hat{m}(x, h)$  as a estimate of  $m(x)$ . For a given value of  $x$  the MSE is

$$MSE_n(x, h) = \mathbb{E} \left( (\hat{m}(x, h) - m(x))^2 \right).$$

We are typically interested in estimating  $m(x)$  for all values in the support of  $x$ . A common measure for the average fit is the integrated MSE

$$\begin{aligned} IMSE_n(h) &= \int MSE_n(x, h) f_x(x) dx \\ &= \int \mathbb{E} \left( (\hat{m}(x, h) - m(x))^2 \right) f_x(x) dx \end{aligned}$$

where  $f_x(x)$  is the marginal density of  $x_i$ . Notice that we have defined the IMSE as an integral with respect to the density  $f_x(x)$ . Other weight functions could be used, but it turns out that this is a convenient choice.

The IMSE is closely related with the MSFE of Section 4.11. Let  $(y_{n+1}, x_{n+1})$  be out-of-sample observations (and thus independent of the sample) and consider predicting  $y_{n+1}$  given  $x_{n+1}$  and the nonparametric estimate  $\hat{m}(x, h)$ . The natural point estimate for  $y_{n+1}$  is  $\hat{m}(x_{n+1}, h)$  which has mean-squared forecast error

$$\begin{aligned} MSFE_n(h) &= \mathbb{E} \left( (y_{n+1} - \hat{m}(x_{n+1}, h))^2 \right) \\ &= \mathbb{E} \left( (e_{n+1} + m(x_{n+1}) - \hat{m}(x_{n+1}, h))^2 \right) \\ &= \sigma^2 + \mathbb{E} \left( (m(x_{n+1}) - \hat{m}(x_{n+1}, h))^2 \right) \\ &= \sigma^2 + \int \mathbb{E} \left( (\hat{m}(x, h) - m(x))^2 \right) f_x(x) dx \end{aligned}$$

where the final equality uses the fact that  $x_{n+1}$  is independent of  $\hat{m}(x, h)$ . We thus see that

$$MSFE_n(h) = \sigma^2 + IMSE_n(h).$$

Since  $\sigma^2$  is a constant independent of the bandwidth  $h$ ,  $MSFE_n(h)$  and  $IMSE_n(h)$  are equivalent measures of the fit of the nonparametric regression.

The optimal bandwidth  $h$  is the value which minimizes  $IMSE_n(h)$  (or equivalently  $MSFE_n(h)$ ). While these functions are unknown, we learned in Theorem 4.11.1 that (at least in the case of linear

regression)  $MSFE_n$  can be estimated by the sample mean-squared prediction errors. It turns out that this fact extends to nonparametric regression. The nonparametric leave-one-out residuals are

$$\tilde{e}_i(h) = y_i - \tilde{m}_{-i}(x_i, h)$$

where we are being explicit about the dependence on the bandwidth  $h$ . The mean squared leave-one-out residuals is

$$CV(h) = \frac{1}{n} \sum_{i=1}^n \tilde{e}_i(h)^2.$$

This function of  $h$  is known as the **cross-validation criterion**.

The cross-validation bandwidth  $\hat{h}$  is the value which minimizes  $CV(h)$

$$\hat{h} = \underset{h \geq h_\ell}{\operatorname{argmin}} CV(h) \quad (17.7)$$

for some  $h_\ell > 0$ . The restriction  $h \geq h_\ell$  is imposed so that  $CV(h)$  is not evaluated over unreasonably small bandwidths.

There is not an explicit solution to the minimization problem (17.7), so it must be solved numerically. A typical practical method is to create a grid of values for  $h$ , e.g.  $[h_1, h_2, \dots, h_J]$ , evaluate  $CV(h_j)$  for  $j = 1, \dots, J$ , and set

$$\hat{h} = \underset{h \in [h_1, h_2, \dots, h_J]}{\operatorname{argmin}} CV(h).$$

Evaluation using a coarse grid is typically sufficient for practical application. Plots of  $CV(h)$  against  $h$  are a useful diagnostic tool to verify that the minimum of  $CV(h)$  has been obtained.

We said above that the cross-validation criterion is an estimator of the MSFE. This claim is based on the following result.

**Theorem 17.6.1**

$$\mathbb{E}(CV(h)) = MSFE_{n-1}(h) = IMSE_{n-1}(h) + \sigma^2 \quad (17.8)$$

Theorem 17.6.1 shows that  $CV(h)$  is an unbiased estimator of  $IMSE_{n-1}(h) + \sigma^2$ . The first term,  $IMSE_{n-1}(h)$ , is the integrated MSE of the nonparametric estimator using a sample of size  $n - 1$ . If  $n$  is large,  $IMSE_{n-1}(h)$  and  $IMSE_n(h)$  will be nearly identical, so  $CV(h)$  is essentially unbiased as an estimator of  $IMSE_n(h) + \sigma^2$ . Since the second term ( $\sigma^2$ ) is unaffected by the bandwidth  $h$ , it is irrelevant for the problem of selection of  $h$ . In this sense we can view  $CV(h)$  as an estimator of the IMSE, and more importantly we can view the minimizer of  $CV(h)$  as an estimate of the minimizer of  $IMSE_n(h)$ .

To illustrate, Figure 17.3 displays the cross-validation criteria  $CV(h)$  for the Nadaraya-Watson and Local Linear estimators using the data from Figure 17.1, both using the Epanechnikov kernel. The CV functions are computed on a grid with intervals 0.01. The CV-minimizing bandwidths are  $h = 1.09$  for the Nadaraya-Watson estimator and  $h = 1.59$  for the local linear estimator. Figure 17.3 shows the minimizing bandwidths by the arrows. It is typical to find that the CV criteria recommends a larger bandwidth for the LL estimator than for the NW estimator, which highlights the fact that smoothing parameters such as bandwidths are specific to the particular method.

The CV criterion can also be used to select between different nonparametric estimators. The CV-selected estimator is the one with the lowest minimized CV criterion. For example, in Figure 17.3, the NW estimator has a minimized CV criterion of 16.88, while the LL estimator has a

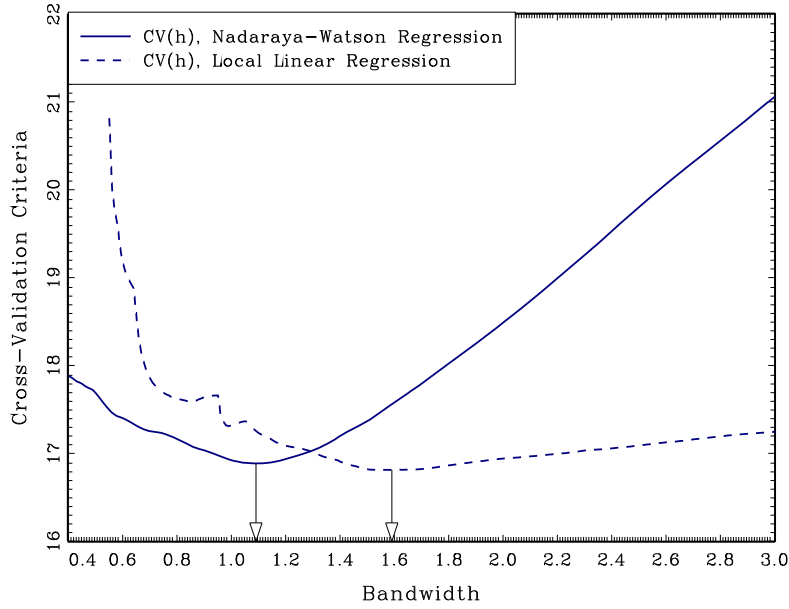


Figure 17.3: Cross-Validation Criteria, Nadaraya-Watson Regression and Local Linear Regression

minimized CV criterion of 16.81. Since the LL estimator achieves a lower value of the CV criterion, LL is the CV-selected estimator. The difference (0.07) is small, suggesting that the two estimators are near equivalent in IMSE.

Figure 17.4 displays the fitted CEF estimates (NW and LL) using the bandwidths selected by cross-validation. Also displayed is the true CEF  $m(x) = 10 \log(x)$ . Notice that the nonparametric estimators with the CV-selected bandwidths (and especially the LL estimator) track the true CEF quite well.

**Proof of Theorem 17.6.1.** Observe that  $m(x_i) - \tilde{m}_{-i}(x_i, h)$  is a function only of  $(x_1, \dots, x_n)$  and  $(e_1, \dots, e_n)$  excluding  $e_i$ , and is thus uncorrelated with  $e_i$ . Since  $\tilde{e}_i(h) = m(x_i) - \tilde{m}_{-i}(x_i, h) + e_i$ , then

$$\begin{aligned}
 \mathbb{E}(CV(h)) &= \mathbb{E}(\tilde{e}_i(h)^2) \\
 &= \mathbb{E}(e_i^2) + \mathbb{E}\left((\tilde{m}_{-i}(x_i, h) - m(x_i))^2\right) \\
 &\quad + 2\mathbb{E}((\tilde{m}_{-i}(x_i, h) - m(x_i))e_i) \\
 &= \sigma^2 + \mathbb{E}\left((\tilde{m}_{-i}(x_i, h) - m(x_i))^2\right). \tag{17.9}
 \end{aligned}$$

The second term is an expectation over the random variables  $x_i$  and  $\tilde{m}_{-i}(x, h)$ , which are independent as the second is not a function of the  $i^{th}$  observation. Thus taking the conditional expectation given the sample excluding the  $i^{th}$  observation, this is the expectation over  $x_i$  only, which is the integral with respect to its density

$$\mathbb{E}_{-i}\left((\tilde{m}_{-i}(x_i, h) - m(x_i))^2\right) = \int (\tilde{m}_{-i}(x, h) - m(x))^2 f_x(x) dx.$$

Taking the unconditional expectation yields

$$\begin{aligned}
 \mathbb{E}\left((\tilde{m}_{-i}(x_i, h) - m(x_i))^2\right) &= \mathbb{E} \int (\tilde{m}_{-i}(x, h) - m(x))^2 f_x(x) dx \\
 &= IMSE_{n-1}(h)
 \end{aligned}$$

where this is the IMSE of a sample of size  $n - 1$  as the estimator  $\tilde{m}_{-i}$  uses  $n - 1$  observations. Combined with (17.9) we obtain (17.8), as desired. ■

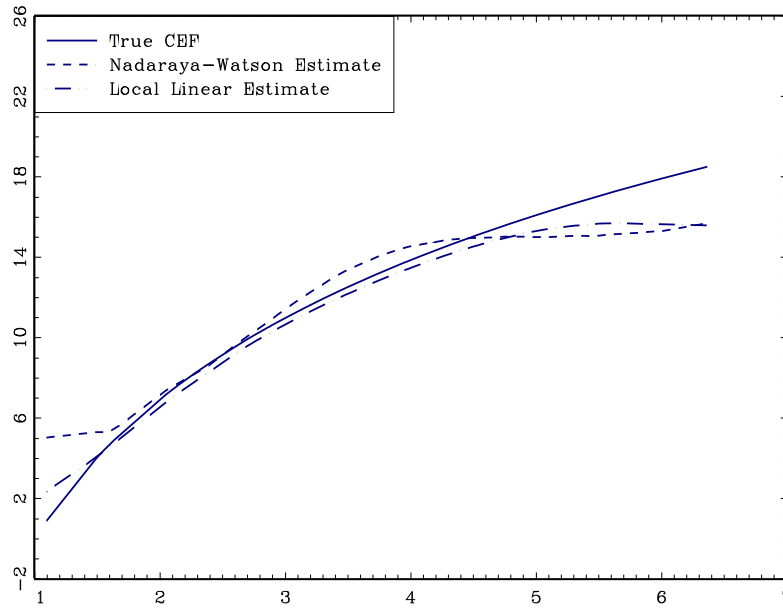


Figure 17.4: Nonparametric Estimates using data-dependent (CV) bandwidths

## 17.7 Asymptotic Distribution

There is no finite sample distribution theory for kernel estimators, but there is a well developed asymptotic distribution theory. The theory is based on the approximation that the bandwidth  $h$  decreases to zero as the sample size  $n$  increases. This means that the smoothing is increasingly localized as the sample size increases. So long as the bandwidth does not decrease to zero too quickly, the estimator can be shown to be asymptotically normal, but with a non-trivial bias.

Let  $f_x(x)$  denote the marginal density of  $x_i$  and  $\sigma^2(x) = \mathbb{E}(e_i^2 \mid x_i = x)$  denote the conditional variance of  $e_i = y_i - m(x_i)$ .

**Theorem 17.7.1** *Let  $\hat{m}(x)$  denote either the Nadarya-Watson or Local Linear estimator of  $m(x)$ . If  $x$  is interior to the support of  $x_i$  and  $f_x(x) > 0$ , then as  $n \rightarrow \infty$  and  $h \rightarrow 0$  such that  $nh \rightarrow \infty$ ,*

$$\sqrt{nh} \left( \hat{m}(x) - m(x) - h^2 \sigma_k^2 B(x) \right) \xrightarrow{d} N \left( 0, \frac{R_k \sigma^2(x)}{f_x(x)} \right) \quad (17.10)$$

where  $\sigma_k^2$  and  $R_k$  are defined in (17.5) and (17.6). For the Nadaraya-Watson estimator

$$B(x) = \frac{1}{2} m''(x) + f_x(x)^{-1} f'_x(x) m'(x)$$

and for the local linear estimator

$$B(x) = \frac{1}{2} f_x(x) m''(x)$$

There are several interesting features about the asymptotic distribution which are noticeably different than for parametric estimators. First, the estimator converges at the rate  $\sqrt{nh}$ , not  $\sqrt{n}$ .



Since  $h \rightarrow 0$ ,  $\sqrt{nh}$  diverges slower than  $\sqrt{n}$ , thus the nonparametric estimator converges more slowly than a parametric estimator. Second, the asymptotic distribution contains a non-negligible bias term  $h^2 \sigma_k^2 B(x)$ . This term asymptotically disappears since  $h \rightarrow 0$ . Third, the assumptions that  $nh \rightarrow \infty$  and  $h \rightarrow 0$  mean that the estimator is consistent for the CEF  $m(x)$ .

The fact that the estimator converges at the rate  $\sqrt{nh}$  has led to the interpretation of  $nh$  as the “effective sample size”. This is because the number of observations being used to construct  $\hat{m}(x)$  is proportional to  $nh$ , not  $n$  as for a parametric estimator.

It is helpful to understand that the nonparametric estimator has a reduced convergence rate because the object being estimated –  $m(x)$  – is nonparametric. This is harder than estimating a finite dimensional parameter, and thus comes at a cost.

Unlike parametric estimation, the asymptotic distribution of the nonparametric estimator includes a term representing the bias of the estimator. The asymptotic distribution (17.10) shows the form of this bias. Not only is it proportional to the squared bandwidth  $h^2$  (the degree of smoothing), it is proportional to the function  $B(x)$  which depends on the slope and curvature of the CEF  $m(x)$ . Interestingly, when  $m(x)$  is constant then  $B(x) = 0$  and the kernel estimator has no asymptotic bias. The bias is essentially increasing in the curvature of the CEF function  $m(x)$ . This is because the local averaging smooths  $m(x)$ , and the smoothing induces more bias when  $m(x)$  is curved.

Theorem 17.7.1 shows that the asymptotic distributions of the NW and LL estimators are similar, with the only difference arising in the bias function  $B(x)$ . The bias term for the NW estimator has an extra component which depends on the first derivative of the CEF  $m(x)$  while the bias term of the LL estimator is invariant to the first derivative. The fact that the bias formula for the LL estimator is simpler and is free of dependence on the first derivative of  $m(x)$  suggests that the LL estimator will generally have smaller bias than the NW estimator (but this is not a precise ranking). Since the asymptotic variances in the two distributions are the same, this means that the LL estimator achieves a reduced bias without an effect on asymptotic variance. This analysis has led to the general preference for the LL estimator over the NW estimator in the nonparametrics literature.

One implication of Theorem 17.7.1 is that we can define the asymptotic MSE (AMSE) of  $\hat{m}(x)$  as the squared bias plus the asymptotic variance

$$AMSE(\hat{m}(x)) = (h^2 \sigma_k^2 B(x))^2 + \frac{R_k \sigma^2(x)}{nh f_x(x)}. \quad (17.11)$$

Focusing on rates, this says

$$AMSE(\hat{m}(x)) \sim h^4 + \frac{1}{nh} \quad (17.12)$$

which means that the AMSE is dominated by the larger of  $h^4$  and  $(nh)^{-1}$ . Notice that the bias is increasing in  $h$  and the variance is decreasing in  $h$ . (More smoothing means more observations are used for local estimation: this increases the bias but decreases estimation variance.) To select  $h$  to minimize the AMSE, these two components should balance each other. Setting  $h^4 \propto (nh)^{-1}$  means setting  $h \propto n^{-1/5}$ . Another way to see this is to pick  $h$  to minimize the right-hand-side of (17.12). The first-order condition for  $h$  is

$$\frac{\partial}{\partial h} \left( h^4 + \frac{1}{nh} \right) = 4h^3 - \frac{1}{nh^2} = 0$$

which when solved for  $h$  yields  $h = n^{-1/5}$ . What this means is that for AMSE-efficient estimation of  $m(x)$ , the optimal rate for the bandwidth is  $h \propto n^{-1/5}$ .

**Theorem 17.7.2** *The bandwidth which minimizes the AMSE (17.12) is of order  $h \propto n^{-1/5}$ . With  $h \propto n^{-1/5}$  then  $AMSE(\hat{m}(x)) = O(n^{-4/5})$  and  $\hat{m}(x) = m(x) + O_p(n^{-2/5})$ .*

This result means that the bandwidth should take the form  $h = cn^{-1/5}$ . The optimal constant  $c$  depends on the kernel  $k$ , the bias function  $B(x)$  and the marginal density  $f_x(x)$ . A common misinterpretation is to set  $h = n^{-1/5}$ , which is equivalent to setting  $c = 1$  and is completely arbitrary. Instead, an empirical bandwidth selection rule such as cross-validation should be used in practice.

When  $h = cn^{-1/5}$  we can rewrite the asymptotic distribution (17.10) as

$$n^{2/5}(\hat{m}(x) - m(x)) \xrightarrow{d} N\left(c^2\sigma_k^2 B(x), \frac{R_k\sigma^2(x)}{cf_x(x)}\right)$$

In this representation, we see that  $\hat{m}(x)$  is asymptotically normal, but with a  $n^{2/5}$  rate of convergence and non-zero mean. The asymptotic distribution depends on the constant  $c$  through the bias (positively) and the variance (inversely).

The asymptotic distribution in Theorem 17.7.1 allows for the optimal rate  $h = cn^{-1/5}$  but this rate is not required. In particular, consider an undersmoothing (smaller than optimal) bandwidth with rate  $h = o(n^{-1/5})$ . For example, we could specify that  $h = cn^{-\alpha}$  for some  $c > 0$  and  $1/5 < \alpha < 1$ . Then  $\sqrt{nh}h^2 = O(n^{(1-5\alpha)/2}) = o(1)$  so the bias term in (17.10) is asymptotically negligible so Theorem 17.7.1 implies

$$\sqrt{nh}(\hat{m}(x) - m(x)) \xrightarrow{d} N\left(0, \frac{R_k\sigma^2(x)}{f_x(x)}\right).$$

That is, the estimator is asymptotically normal without a bias component. Not having an asymptotic bias component is convenient for some theoretical manipulations, so many authors impose the undersmoothing condition  $h = o(n^{-1/5})$  to ensure this situation. This convenience comes at a cost. First, the resulting estimator is inefficient as its convergence rate is  $O_p(n^{-(1-\alpha)/2}) > O_p(n^{-2/5})$  since  $\alpha > 1/5$ . Second, the distribution theory is an inherently misleading approximation as it misses a critically key ingredient of nonparametric estimation – the trade-off between bias and variance. The approximation (17.10) is superior precisely because it contains the asymptotic bias component which is a realistic implication of nonparametric estimation. Undersmoothing assumptions should be avoided when possible.

## 17.8 Conditional Variance Estimation

Let's consider the problem of estimation of the conditional variance

$$\begin{aligned}\sigma^2(x) &= \text{var}(y_i \mid x_i = x) \\ &= \mathbb{E}(e_i^2 \mid x_i = x).\end{aligned}$$

Even if the conditional mean  $m(x)$  is parametrically specified, it is natural to view  $\sigma^2(x)$  as inherently nonparametric as economic models rarely specify the form of the conditional variance. Thus it is quite appropriate to estimate  $\sigma^2(x)$  nonparametrically.

We know that  $\sigma^2(x)$  is the CEF of  $e_i^2$  given  $x_i$ . Therefore if  $e_i^2$  were observed,  $\sigma^2(x)$  could be nonparametrically estimated using NW or LL regression. For example, the ideal NW estimator is

$$\bar{\sigma}^2(x) = \frac{\sum_{i=1}^n k_i(x)e_i^2}{\sum_{i=1}^n k_i(x)}.$$

Since the errors  $e_i$  are not observed, we need to replace them with an empirical residual, such as  $\hat{e}_i = y_i - \hat{m}(x_i)$  where  $\hat{m}(x)$  is the estimated CEF. (The latter could be a nonparametric estimator such as NW or LL, or even a parametric estimator.) Even better, use the leave-one-out prediction errors  $\tilde{e}_i = y_i - \hat{m}_{-i}(x_i)$ , as these are not subject to overfitting.

With this substitution the NW estimator of the conditional variance is

$$\hat{\sigma}^2(x) = \frac{\sum_{i=1}^n k_i(x)\hat{e}_i^2}{\sum_{i=1}^n k_i(x)}. \quad (17.13)$$

This estimator depends on a set of bandwidths  $h_1, \dots, h_q$ , but there is no reason for the bandwidths to be the same as those used to estimate the conditional mean. Cross-validation can be used to select the bandwidths for estimation of  $\hat{\sigma}^2(x)$  separately from cross-validation for estimation of  $\hat{m}(x)$ .

There is one subtle difference between CEF and conditional variance estimation. The conditional variance is inherently non-negative  $\sigma^2(x) \geq 0$  and it is desirable for our estimator to satisfy this property. Interestingly, the NW estimator (17.13) is necessarily non-negative, since it is a smoothed average of the non-negative squared residuals, but the LL estimator is not guaranteed to be non-negative for all  $x$ . For this reason, the NW estimator is preferred for conditional variance estimation.

Fan and Yao (1998, *Biometrika*) derive the asymptotic distribution of the estimator (17.13). They obtain the surprising result that the asymptotic distribution of this two-step estimator is identical to that of the one-step idealized estimator  $\bar{\sigma}^2(x)$ .

## 17.9 Standard Errors

Theorem 17.7.1 shows the asymptotic variances of both the NW and LL nonparametric regression estimators equal

$$V(x) = \frac{R_k \sigma^2(x)}{f_x(x)}.$$

For standard errors we need an estimate of  $V(x)$ . A plug-in estimate replaces the unknowns by estimates. The roughness  $R_k$  can be found from Table 9.1. The conditional variance can be estimated using (17.13). The density of  $x_i$  can be estimated using the methods from Section 22.1. Replacing these estimates into the formula for  $V(x)$  we obtain the asymptotic variance estimate

$$\hat{V}(x) = \frac{R_k \hat{\sigma}^2(x)}{\hat{f}_x(x)}.$$

Then an asymptotic standard error for the kernel estimate  $\hat{m}(x)$  is

$$\hat{s}(x) = \sqrt{\frac{1}{nh} \hat{V}(x)}.$$

Plots of the estimated CEF  $\hat{m}(x)$  can be accompanied by confidence intervals  $\hat{m}(x) \pm 2\hat{s}(x)$ . These are known as **pointwise confidence intervals**, as they are designed to have correct coverage at each  $x$ , not uniformly in  $x$ .

One important caveat about the interpretation of nonparametric confidence intervals is that they are not centered at the true CEF  $m(x)$ , but rather are centered at the biased or pseudo-true value

$$m^*(x) = m(x) + h^2 \sigma_k^2 B(x).$$

Consequently, a correct statement about the confidence interval  $\hat{m}(x) \pm 2\hat{s}(x)$  is that it asymptotically contains  $m^*(x)$  with probability 95%, not that it asymptotically contains  $m(x)$  with probability 95%. The discrepancy is that the confidence interval does not take into account the bias  $h^2 \sigma_k^2 B(x)$ . Unfortunately, nothing constructive can be done about this. The bias is difficult and noisy to estimate, so making a bias-correction only inflates estimation variance and decreases overall precision. A technical “trick” is to assume undersmoothing  $h = o(n^{-1/5})$  but this does not really eliminate the bias, it only assumes it away. The plain fact is that once we honestly acknowledge that the true CEF is nonparametric, it then follows that any finite sample estimate will have finite sample bias, and this bias will be inherently unknown and thus impossible to incorporate into confidence intervals.

## 17.10 Multiple Regressors

Our analysis has focus on the case of real-valued  $x_i$  for simplicity of exposition, but the methods of kernel regression extend easily to the multiple regressor case, at the cost of a reduced rate of convergence. In this section we consider the case of estimation of the conditional expectation function

$$\mathbb{E}(y_i \mid \mathbf{x}_i = \mathbf{x}) = m(\mathbf{x})$$

when

$$\mathbf{x}_i = \begin{pmatrix} x_{1i} \\ \vdots \\ x_{di} \end{pmatrix}$$

is a  $d$ -vector.

For any evaluation point  $\mathbf{x}$  and observation  $i$ , define the kernel weights

$$k_i(\mathbf{x}) = k\left(\frac{x_{1i} - x_1}{h_1}\right) k\left(\frac{x_{2i} - x_2}{h_2}\right) \cdots k\left(\frac{x_{di} - x_d}{h_d}\right),$$

a  $d$ -fold product kernel. The kernel weights  $k_i(\mathbf{x})$  assess if the regressor vector  $\mathbf{x}_i$  is close to the evaluation point  $\mathbf{x}$  in the Euclidean space  $\mathbb{R}^d$ .

These weights depend on a set of  $d$  bandwidths,  $h_j$ , one for each regressor. We can group them together into a single vector for notational convenience:

$$\mathbf{h} = \begin{pmatrix} h_1 \\ \vdots \\ h_d \end{pmatrix}.$$

Given these weights, the Nadaraya-Watson estimator takes the form

$$\hat{m}(\mathbf{x}) = \frac{\sum_{i=1}^n k_i(\mathbf{x}) y_i}{\sum_{i=1}^n k_i(\mathbf{x})}.$$

For the local-linear estimator, define

$$\mathbf{z}_i(\mathbf{x}) = \begin{pmatrix} 1 \\ \mathbf{x}_i - \mathbf{x} \end{pmatrix}$$

and then the local-linear estimator can be written as  $\hat{m}(\mathbf{x}) = \hat{\alpha}(\mathbf{x})$  where

$$\begin{pmatrix} \hat{\alpha}(\mathbf{x}) \\ \hat{\beta}(\mathbf{x}) \end{pmatrix} = \left( \sum_{i=1}^n k_i(\mathbf{x}) \mathbf{z}_i(\mathbf{x}) \mathbf{z}_i(\mathbf{x})' \right)^{-1} \sum_{i=1}^n k_i(\mathbf{x}) \mathbf{z}_i(\mathbf{x}) y_i \\ = (\mathbf{Z}' \mathbf{K} \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{K} \mathbf{y}$$

where  $\mathbf{K} = \text{diag}\{k_1(x), \dots, k_n(x)\}$ .

In multiple regressor kernel regression, cross-validation remains a recommended method for bandwidth selection. The leave-one-out residuals  $\tilde{e}_i$  and cross-validation criterion  $CV(\mathbf{h})$  are defined identically as in the single regressor case. The only difference is that now the CV criterion is a function over the  $d$ -dimensional bandwidth  $\mathbf{h}$ . This is a critical practical difference since finding the bandwidth vector  $\hat{\mathbf{h}}$  which minimizes  $CV(\mathbf{h})$  can be computationally difficult when  $\mathbf{h}$  is high dimensional. Grid search is cumbersome and costly, since  $G$  gridpoints per dimension imply evaluation of  $CV(\mathbf{h})$  at  $G^d$  distinct points, which can be a large number. Furthermore, plots of  $CV(\mathbf{h})$  against  $\mathbf{h}$  are challenging when  $d > 2$ .

The asymptotic distribution of the estimators in the multiple regressor case is an extension of the single regressor case. Let  $f_x(\mathbf{x})$  denote the marginal density of  $\mathbf{x}_i$  and  $\sigma^2(\mathbf{x}) = \mathbb{E}(e_i^2 \mid \mathbf{x}_i = \mathbf{x})$  the conditional variance of  $e_i = y_i - m(\mathbf{x}_i)$ . Let  $|\mathbf{h}| = h_1 h_2 \cdots h_d$ .

**Theorem 17.10.1** Let  $\hat{m}(\mathbf{x})$  denote either the Nadarya-Watson or Local Linear estimator of  $m(\mathbf{x})$ . If  $\mathbf{x}$  is interior to the support of  $\mathbf{x}_i$  and  $f_x(\mathbf{x}) > 0$ , then as  $n \rightarrow \infty$  and  $h_j \rightarrow 0$  such that  $n|\mathbf{h}| \rightarrow \infty$ ,

$$\sqrt{n|\mathbf{h}|} \left( \hat{m}(\mathbf{x}) - m(\mathbf{x}) - \sigma_k^2 \sum_{j=1}^d h_j^2 B_{j,j}(\mathbf{x}) \right) \xrightarrow{d} N \left( 0, \frac{R_k^d \sigma^2(\mathbf{x})}{f_x(\mathbf{x})} \right)$$

where for the Nadaraya-Watson estimator

$$B_j(\mathbf{x}) = \frac{1}{2} \frac{\partial^2}{\partial x_j^2} m(\mathbf{x}) + f_x(\mathbf{x})^{-1} \frac{\partial}{\partial x_j} f_x(\mathbf{x}) \frac{\partial}{\partial x_j} m(\mathbf{x})$$

and for the Local Linear estimator

$$B_j(\mathbf{x}) = \frac{1}{2} \frac{\partial^2}{\partial x_j^2} m(\mathbf{x})$$

For notational simplicity consider the case that there is a single common bandwidth  $h$ . In this case the AMSE takes the form

$$AMSE(\hat{m}(\mathbf{x})) \sim h^4 + \frac{1}{nh^d}$$

That is, the squared bias is of order  $h^4$ , the same as in the single regressor case, but the variance is of larger order  $(nh^d)^{-1}$ . Setting  $h$  to balance these two components requires setting  $h \sim n^{-1/(4+d)}$ .

**Theorem 17.10.2** The bandwidth which minimizes the AMSE is of order  $h \propto n^{-1/(4+d)}$ . With  $h \propto n^{-1/(4+d)}$  then  $AMSE(\hat{m}(\mathbf{x})) = O(n^{-4/(4+d)})$  and  $\hat{m}(\mathbf{x}) = m(\mathbf{x}) + O_p(n^{-2/(4+d)})$

In all estimation problems an increase in the dimension decreases estimation precision. For example, in parametric estimation an increase in dimension typically increases the asymptotic variance. In nonparametric estimation an increase in the dimension typically decreases the convergence rate, which is a more fundamental decrease in precision. For example, in kernel regression the convergence rate  $O_p(n^{-2/(4+d)})$  decreases as  $d$  increases. The reason is the estimator  $\hat{m}(\mathbf{x})$  is a local average of the  $y_i$  for observations such that  $\mathbf{x}_i$  is close to  $\mathbf{x}$ , and when there are multiple regressors the number of such observations is inherently smaller. This phenomenon – that the rate of convergence of nonparametric estimation decreases as the dimension increases – is called the **curse of dimensionality**.

# Chapter 18

## Series Estimation

### 18.1 Approximation by Series

As we mentioned at the beginning of Chapter 17, there are two main methods of nonparametric regression: kernel estimation and series estimation. In this chapter we study series methods.

Series methods approximate an unknown function (e.g. the CEF  $m(\mathbf{x})$ ) with a flexible parametric function, with the number of parameters treated similarly to the bandwidth in kernel regression. A series approximation to  $m(\mathbf{x})$  takes the form  $m_K(\mathbf{x}) = m_K(\mathbf{x}, \boldsymbol{\beta}_K)$  where  $m_K(\mathbf{x}, \boldsymbol{\beta}_K)$  is a known parametric family and  $\boldsymbol{\beta}_K$  is an unknown coefficient. The integer  $K$  is the dimension of  $\boldsymbol{\beta}_K$  and indexes the complexity of the approximation.

A linear series approximation takes the form

$$\begin{aligned} m_K(\mathbf{x}) &= \sum_{j=1}^K z_{jK}(\mathbf{x}) \beta_{jK} \\ &= \mathbf{z}_K(\mathbf{x})' \boldsymbol{\beta}_K \end{aligned} \tag{18.1}$$

where  $z_{jK}(\mathbf{x})$  are (nonlinear) functions of  $\mathbf{x}$ , and are known as **basis functions** or **basis function transformations** of  $\mathbf{x}$ .

For real-valued  $x$ , a well-known linear series approximation is the  $p^{th}$ -order **polynomial**

$$m_K(x) = \sum_{j=0}^p x^j \beta_{jK}$$

where  $K = p + 1$ .

When  $\mathbf{x} \in \mathbb{R}^d$  is vector-valued, a  $p^{th}$ -order polynomial is

$$m_K(\mathbf{x}) = \sum_{j_1=0}^p \cdots \sum_{j_d=0}^p x_1^{j_1} \cdots x_d^{j_d} \beta_{j_1, \dots, j_d K}.$$

This includes all powers and cross-products, and the coefficient vector has dimension  $K = (p+1)^d$ . In general, a common method to create a series approximation for vector-valued  $\mathbf{x}$  is to include all non-redundant cross-products of the basis function transformations of the components of  $\mathbf{x}$ .

### 18.2 Splines

Another common series approximation is a continuous piecewise polynomial function known as a **spline**. While splines can be of any polynomial order (e.g. linear, quadratic, cubic, etc.), a common choice is cubic. To impose smoothness it is common to constrain the spline function to have continuous derivatives up to the order of the spline. Thus a quadratic spline is typically

constrained to have a continuous first derivative, and a cubic spline is typically constrained to have a continuous first and second derivative.

There is more than one way to define a spline series expansion. All are based on the number of **knots** – the join points between the polynomial segments.

To illustrate, a piecewise linear function with two segments and a knot at  $t$  is

$$m_K(x) = \begin{cases} m_1(x) = \beta_{00} + \beta_{01}(x - t) & x < t \\ m_2(x) = \beta_{10} + \beta_{11}(x - t) & x \geq t \end{cases}$$

(For convenience we have written the segments functions as polynomials in  $x - t$ .) The function  $m_K(x)$  equals the linear function  $m_1(x)$  for  $x < t$  and equals  $m_2(t)$  for  $x > t$ . Its left limit at  $x = t$  is  $\beta_{00}$  and its right limit is  $\beta_{10}$ , so is continuous if (and only if)  $\beta_{00} = \beta_{10}$ . Enforcing this constraint is equivalent to writing the function as

$$m_K(x) = \beta_0 + \beta_1(x - t) + \beta_2(x - t)1(x \geq t)$$

or after transforming coefficients, as

$$m_K(x) = \beta_0 + \beta_1x + \beta_2(x - t)1(x \geq t).$$

Notice that this function has  $K = 3$  coefficients, the same as a quadratic polynomial.

A piecewise quadratic function with one knot at  $t$  is

$$m_K(x) = \begin{cases} m_1(x) = \beta_{00} + \beta_{01}(x - t) + \beta_{02}(x - t)^2 & x < t \\ m_2(x) = \beta_{10} + \beta_{11}(x - t) + \beta_{12}(x - t)^2 & x \geq t \end{cases}$$

This function is continuous at  $x = t$  if  $\beta_{00} = \beta_{10}$ , and has a continuous first derivative if  $\beta_{01} = \beta_{11}$ . Imposing these constraints and rewriting, we obtain the function

$$m_K(x) = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3(x - t)^21(x \geq t).$$

Here,  $K = 4$ .

Furthermore, a piecewise cubic function with one knot and a continuous second derivative is

$$m_K(x) = \beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3 + \beta_4(x - t)^31(x \geq t)$$

which has  $K = 5$ .

The polynomial order  $p$  is selected to control the smoothness of the spline, as  $m_K(x)$  has continuous derivatives up to  $p - 1$ .

In general, a  $p^{th}$ -order spline with  $N$  knots at  $t_1, t_2, \dots, t_N$  with  $t_1 < t_2 < \dots < t_N$  is

$$m_K(x) = \sum_{j=0}^p \beta_j x^j + \sum_{k=1}^N \gamma_k (x - t_k)^p 1(x \geq t_k)$$

which has  $K = N + p + 1$  coefficients.

In spline approximation, the typical approach is to treat the polynomial order  $p$  as fixed, and select the number of knots  $N$  to determine the complexity of the approximation. The knots  $t_k$  are typically treated as fixed. A common choice is to set the knots to evenly partition the support  $\mathcal{X}$  of  $\mathbf{x}_i$ .

### 18.3 Partially Linear Model

A common use of a series expansion is to allow the CEF to be nonparametric with respect to one variable, yet linear in the other variables. This allows flexibility in a particular variable of interest. A partially linear CEF with vector-valued regressor  $\mathbf{x}_1$  and real-valued continuous  $x_2$  takes the form

$$m(\mathbf{x}_1, x_2) = \mathbf{x}_1' \boldsymbol{\beta}_1 + m_2(x_2).$$

This model is commonly used when  $\mathbf{x}_1$  are discrete (e.g. binary variables) and  $x_2$  is continuously distributed.

Series methods are particularly convenient for estimation of partially linear models, as we can replace the unknown function  $m_2(x_2)$  with a series expansion to obtain

$$\begin{aligned} m(\mathbf{x}) &\simeq m_K(\mathbf{x}) \\ &= \mathbf{x}_1' \boldsymbol{\beta}_1 + \mathbf{z}_K' \boldsymbol{\beta}_{2K} \\ &= \mathbf{x}_K' \boldsymbol{\beta}_K \end{aligned}$$

where  $\mathbf{z}_K = \mathbf{z}_K(x_2)$  are the basis transformations of  $x_2$  (typically polynomials or splines) and  $\boldsymbol{\beta}_{2K}$  are coefficients. After transformation the regressors are  $\mathbf{x}_K = (\mathbf{x}_1', \mathbf{z}_K')$  and the coefficients are  $\boldsymbol{\beta}_K = (\boldsymbol{\beta}_1', \boldsymbol{\beta}_{2K}')'$ .

### 18.4 Additively Separable Models

When  $\mathbf{x}$  is multivariate a common simplification is to treat the regression function  $m(\mathbf{x})$  as additively separable in the individual regressors, which means that

$$m(\mathbf{x}) = m_1(x_1) + m_2(x_2) + \cdots + m_d(x_d).$$

Series methods are quite convenient for additively separable models, as we simply apply series expansions (polynomials or splines) separately for each component  $m_j(x_j)$ . The advantage of additive separability is the reduction in dimensionality. While an unconstrained  $p^{\text{th}}$  order polynomial has  $(p+1)^d$  coefficients, an additively separable polynomial model has only  $(p+1)d$  coefficients. This can be a major reduction in the number of coefficients. The disadvantage of this simplification is that the interaction effects have been eliminated.

The decision to impose additive separability can be based on an economic model which suggests the absence of interaction effects, or can be a model selection decision similar to the selection of the number of series terms. We will discuss model selection methods below.

### 18.5 Uniform Approximations

A good series approximation  $m_K(\mathbf{x})$  will have the property that it gets close to the true CEF  $m(\mathbf{x})$  as the complexity  $K$  increases. Formal statements can be derived from the theory of functional analysis.

An elegant and famous theorem is the **Stone-Weierstrass theorem**, (Weierstrass, 1885, Stone 1937, 1948) which states that any continuous function can be arbitrarily uniformly well approximated by a polynomial of sufficiently high order. Specifically, the theorem states that for  $\mathbf{x} \in \mathbb{R}^d$ , if  $m(\mathbf{x})$  is continuous on a compact set  $\mathcal{X}$ , then for any  $\varepsilon > 0$  there exists a polynomial  $m_K(\mathbf{x})$  of some order  $K$  which is uniformly within  $\varepsilon$  of  $m(\mathbf{x})$ :

$$\sup_{\mathbf{x} \in \mathcal{X}} |m_K(\mathbf{x}) - m(\mathbf{x})| \leq \varepsilon. \quad (18.2)$$

Thus the true unknown  $m(\mathbf{x})$  can be arbitrarily well approximated by selecting a suitable polynomial.



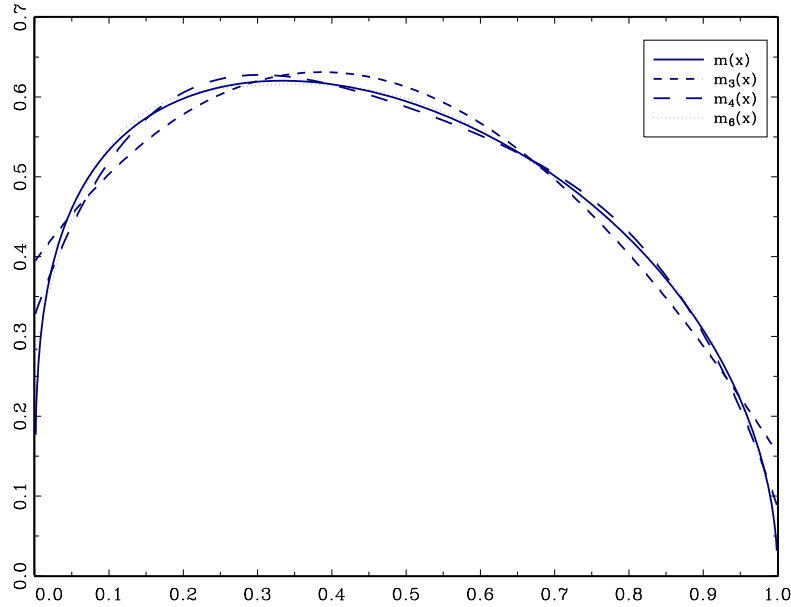


Figure 18.1: True CEF and Best Approximations

The result (18.2) can be strengthened. In particular, if the  $s^{th}$  derivative of  $m(\mathbf{x})$  is continuous then the uniform approximation error satisfies

$$\sup_{\mathbf{x} \in \mathcal{X}} |m_K(\mathbf{x}) - m(\mathbf{x})| = O(K^{-\alpha}) \quad (18.3)$$

as  $K \rightarrow \infty$  where  $\alpha = s/d$ . This result is more useful than (18.2) because it gives a rate at which the approximation  $m_K(\mathbf{x})$  approaches  $m(\mathbf{x})$  as  $K$  increases.

Both (18.2) and (18.3) hold for spline approximations as well.

Intuitively, the number of derivatives  $s$  indexes the smoothness of the function  $m(\mathbf{x})$ . (18.3) says that the best rate at which a polynomial or spline approximates the CEF  $m(\mathbf{x})$  depends on the underlying smoothness of  $m(\mathbf{x})$ . The more smooth is  $m(\mathbf{x})$ , the fewer series terms (polynomial order or spline knots) are needed to obtain a good approximation.

To illustrate polynomial approximation, Figure 18.1 displays the CEF  $m(x) = x^{1/4}(1-x)^{1/2}$  on  $x \in [0, 1]$ . In addition, the best approximations using polynomials of order  $K = 3$ ,  $K = 4$ , and  $K = 6$  are displayed. You can see how the approximation with  $K = 3$  is fairly crude, but improves with  $K = 4$  and especially  $K = 6$ . Approximations obtained with cubic splines are quite similar so not displayed.

As a series approximation can be written as  $m_K(\mathbf{x}) = \mathbf{z}_K(\mathbf{x})'\boldsymbol{\beta}_K$  as in (18.1), then the coefficient of the best uniform approximation (18.3) is then

$$\boldsymbol{\beta}_K^* = \underset{\boldsymbol{\beta}_K}{\operatorname{argmin}} \sup_{\mathbf{x} \in \mathcal{X}} |\mathbf{z}_K(\mathbf{x})'\boldsymbol{\beta}_K - m(\mathbf{x})|. \quad (18.4)$$

The approximation error is

$$r_K^*(\mathbf{x}) = m(\mathbf{x}) - \mathbf{z}_K(\mathbf{x})'\boldsymbol{\beta}_K^*.$$

We can write this as

$$m(\mathbf{x}) = \mathbf{z}_K(\mathbf{x})'\boldsymbol{\beta}_K^* + r_K^*(\mathbf{x}) \quad (18.5)$$

to emphasize that the true conditional mean can be written as the linear approximation plus error. A useful consequence of equation (18.3) is

$$\sup_{\mathbf{x} \in \mathcal{X}} |r_K^*(\mathbf{x})| \leq O(K^{-\alpha}). \quad (18.6)$$

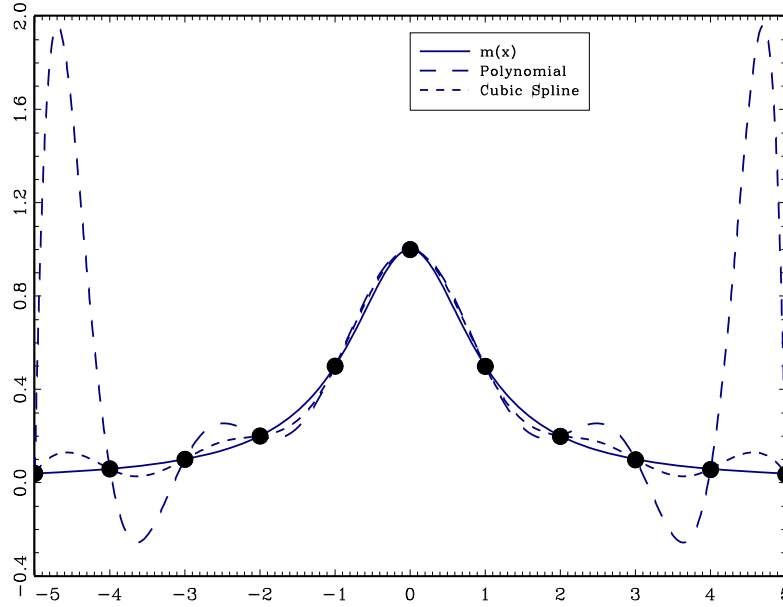


Figure 18.2: True CEF, polynomial interpolation, and spline interpolation

## 18.6 Runge's Phenomenon

Despite the excellent approximation implied by the Stone-Weierstrass theorem, polynomials have the troubling disadvantage that they are very poor at simple interpolation. The problem is known as **Runge's phenomenon**, and is illustrated in Figure 18.2. The solid line is the CEF  $m(x) = (1 + x^2)^{-1}$  displayed on  $[-5, 5]$ . The circles display the function at the  $K = 11$  integers in this interval. The long dashes display the 10<sup>th</sup> order polynomial fit through these points. Notice that the polynomial approximation is erratic and far from the smooth CEF. This discrepancy gets worse as the number of evaluation points increases, as Runge (1901) showed that the discrepancy increases to infinity with  $K$ .

In contrast, splines do not exhibit Runge's phenomenon. In Figure 18.2 the short dashes display a cubic spline with seven knots fit through the same points as the polynomial. While the fitted spline displays some oscillation relative to the true CEF, they are relatively moderate.

Because of Runge's phenomenon, high-order polynomials are not used for interpolation, and are not popular choices for high-order series approximations. Instead, splines are widely used.

## 18.7 Approximating Regression

For each observation  $i$  we observe  $(y_i, \mathbf{x}_i)$  and then construct the regressor vector  $\mathbf{z}_{Ki} = \mathbf{z}_K(\mathbf{x}_i)$  using the series transformations. Stacking the observations in the matrices  $\mathbf{y}$  and  $\mathbf{Z}_K$ , the least squares estimate of the coefficient  $\beta_K$  in the series approximation  $\mathbf{z}_K(\mathbf{x})'\beta_K$  is

$$\hat{\beta}_K = (\mathbf{Z}_K' \mathbf{Z}_K)^{-1} \mathbf{Z}_K' \mathbf{y},$$

and the least squares estimate of the regression function is

$$\hat{m}_K(\mathbf{x}) = \mathbf{z}_K(\mathbf{x})' \hat{\beta}_K. \quad (18.7)$$

As we learned in Chapter 2, the least-squares coefficient is estimating the best linear predictor of  $y_i$  given  $\mathbf{z}_{Ki}$ . This is

$$\beta_K = \mathbb{E} (\mathbf{z}_{Ki} \mathbf{z}_{Ki}')^{-1} \mathbb{E} (\mathbf{z}_{Ki} y_i).$$

Given this coefficient, the series approximation is  $\mathbf{z}_K(\mathbf{x})'\boldsymbol{\beta}_K$  with approximation error

$$r_K(\mathbf{x}) = m(\mathbf{x}) - \mathbf{z}_K(\mathbf{x})'\boldsymbol{\beta}_K. \quad (18.8)$$

The true CEF equation for  $y_i$  is

$$y_i = m(\mathbf{x}_i) + e_i \quad (18.9)$$

with  $e_i$  the CEF error. Defining  $r_{Ki} = r_K(\mathbf{x}_i)$ , we find

$$y_i = \mathbf{z}'_{Ki}\boldsymbol{\beta}_K + e_{Ki}$$

where the equation error is

$$e_{Ki} = r_{Ki} + e_i.$$

Observe that the error  $e_{Ki}$  includes the approximation error and thus does not have the properties of a CEF error.

In matrix notation we can write these equations as

$$\begin{aligned} \mathbf{y} &= \mathbf{Z}_K\boldsymbol{\beta}_K + \mathbf{r}_K + \mathbf{e} \\ &= \mathbf{Z}_K\boldsymbol{\beta}_K + \mathbf{e}_K. \end{aligned} \quad (18.10)$$

We now impose some regularity conditions on the regression model to facilitate the theory. Define the  $K \times K$  expected design matrix

$$\mathbf{Q}_K = \mathbb{E}(\mathbf{z}_{Ki}\mathbf{z}'_{Ki}),$$

let  $\mathcal{X}$  denote the support of  $\mathbf{x}_i$ , and define the largest normalized length of the regressor vector in the support of  $\mathbf{x}_i$

$$\zeta_K = \sup_{\mathbf{x} \in \mathcal{X}} (\mathbf{z}_K(\mathbf{x})'\mathbf{Q}_K^{-1}\mathbf{z}_K(\mathbf{x}))^{1/2}. \quad (18.11)$$

$\zeta_K$  will increase with  $K$ . For example, if the support of the variables  $\mathbf{z}_K(\mathbf{x}_i)$  is the unit cube  $[0, 1]^K$ , then you can compute that  $\zeta_K = \sqrt{K}$ . As discussed in Newey (1997) and Li and Racine (2007, Corollary 15.1) if the support of  $\mathbf{x}_i$  is compact then  $\zeta_K = O(K)$  for polynomials and  $\zeta_K = O(K^{1/2})$  for splines.

#### Assumption 18.7.1

1. For some  $\alpha > 0$  the series approximation satisfies (18.3).
2.  $\mathbb{E}(e_i^2 \mid \mathbf{x}_i) \leq \bar{\sigma}^2 < \infty$ .
3.  $\lambda_{\min}(\mathbf{Q}_K) \geq \underline{\lambda} > 0$ .
4.  $K = K(n)$  is a function of  $n$  which satisfies  $K/n \rightarrow 0$  and  $\zeta_K^2 K/n \rightarrow 0$  as  $n \rightarrow \infty$ .

Assumptions 18.7.1.1 through 18.7.1.3 concern properties of the regression model. Assumption 18.7.1.1 holds with  $\alpha = s/d$  if  $\mathcal{X}$  is compact and the  $s$ 'th derivative of  $m(\mathbf{x})$  is continuous. Assumption 18.7.1.2 allows for conditional heteroskedasticity, but requires the conditional variance to be bounded. Assumption 18.7.1.3 excludes near-singular designs. Since estimates of the conditional mean are unchanged if we replace  $\mathbf{z}_{Ki}$  with  $\mathbf{z}_{Ki}^* = \mathbf{B}_K\mathbf{z}_{Ki}$  for any non-singular  $\mathbf{B}_K$ , Assumption 18.7.1.3 can be viewed as holding after transformation by an appropriate non-singular  $\mathbf{B}_K$ .

Assumption 18.7.1.4 concerns the choice of the number of series terms, which is under the control of the user. It specifies that  $K$  can increase with sample size, but at a controlled rate of growth. Since  $\zeta_K = O(K)$  for polynomials and  $\zeta_K = O(K^{1/2})$  for splines, Assumption 18.7.1.4 is satisfied if  $K^3/n \rightarrow 0$  for polynomials and  $K^2/n \rightarrow 0$  for splines. This means that while the number of series terms  $K$  can increase with the sample size,  $K$  must increase at a much slower rate.

In Section 18.5 we introduced the best uniform approximation, and in this section we introduced the best linear predictor. What is the relationship? They may be similar in practice, but they are not the same and we should be careful to maintain the distinction. Note that from (18.5) we can write  $m(\mathbf{x}_i) = \mathbf{z}'_{Ki} \boldsymbol{\beta}_K^* + r_{Ki}^*$  where  $r_{Ki}^* = r_K^*(\mathbf{x}_i)$  satisfies  $\sup_i |r_{Ki}^*| = O(K^{-\alpha})$  from (18.6). Then the best linear predictor equals

$$\begin{aligned} \boldsymbol{\beta}_K &= \mathbb{E}(\mathbf{z}_{Ki} \mathbf{z}'_{Ki})^{-1} \mathbb{E}(\mathbf{z}_{Ki} y_i) \\ &= \mathbb{E}(\mathbf{z}_{Ki} \mathbf{z}'_{Ki})^{-1} \mathbb{E}(\mathbf{z}_{Ki} m(\mathbf{x}_i)) \\ &= \mathbb{E}(\mathbf{z}_{Ki} \mathbf{z}'_{Ki})^{-1} \mathbb{E}(\mathbf{z}_{Ki} (\mathbf{z}'_{Ki} \boldsymbol{\beta}_K^* + r_{Ki}^*)) \\ &= \boldsymbol{\beta}_K^* + \mathbb{E}(\mathbf{z}_{Ki} \mathbf{z}'_{Ki})^{-1} \mathbb{E}(\mathbf{z}_{Ki} r_{Ki}^*). \end{aligned}$$

Thus the difference between the two approximations is

$$\begin{aligned} r_K(\mathbf{x}) - r_K^*(\mathbf{x}) &= \mathbf{z}_K(\mathbf{x})' (\boldsymbol{\beta}_K^* - \boldsymbol{\beta}_K) \\ &= \mathbf{z}_K(\mathbf{x})' \mathbb{E}(\mathbf{z}_{Ki} \mathbf{z}'_{Ki})^{-1} \mathbb{E}(\mathbf{z}_{Ki} r_{Ki}^*). \end{aligned} \quad (18.12)$$

Observe that by the properties of projection

$$\mathbb{E}(\mathbf{r}_{Ki}^{*2}) - \mathbb{E}(\mathbf{r}_{Ki}^* \mathbf{z}_{Ki})' \mathbb{E}(\mathbf{z}_{Ki} \mathbf{z}'_{Ki})^{-1} \mathbb{E}(\mathbf{z}_{Ki} r_{Ki}^*) \geq 0 \quad (18.13)$$

and by (18.6)

$$\mathbb{E}(\mathbf{r}_{Ki}^{*2}) = \int r_K^*(\mathbf{x})^2 f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \leq O(K^{-2\alpha}). \quad (18.14)$$

Then applying the Schwarz inequality to (18.12), Definition (18.11), (18.13) and (18.14), we find

$$\begin{aligned} |r_K(\mathbf{x}) - r_K^*(\mathbf{x})| &\leq \left( \mathbf{z}_K(\mathbf{x})' \mathbb{E}(\mathbf{z}_{Ki} \mathbf{z}'_{Ki})^{-1} \mathbf{z}_K(\mathbf{x}) \right)^{1/2} \\ &\quad \left( \mathbb{E}(\mathbf{r}_{Ki}^* \mathbf{z}_{Ki})' \mathbb{E}(\mathbf{z}_{Ki} \mathbf{z}'_{Ki})^{-1} \mathbb{E}(\mathbf{z}_{Ki} r_{Ki}^*) \right)^{1/2} \\ &\leq O(\zeta_K K^{-\alpha}). \end{aligned} \quad (18.15)$$

It follows that the best linear predictor approximation error satisfies

$$\sup_{\mathbf{x} \in \mathcal{X}} |r_K(\mathbf{x})| \leq O(\zeta_K K^{-\alpha}). \quad (18.16)$$

The bound (18.16) is probably not the best possible, but it shows that the best linear predictor satisfies a uniform approximation bound. Relative to (18.6), the rate is slower by the factor  $\zeta_K$ . The bound (18.16) term is  $o(1)$  as  $K \rightarrow \infty$  if  $\zeta_K K^{-\alpha} \rightarrow 0$ . A sufficient condition is that  $\alpha > 1$  ( $s > d$ ) for polynomials and  $\alpha > 1/2$  ( $s > d/2$ ) for splines, where  $d = \dim(\mathbf{x})$  and  $s$  is the number of continuous derivatives of  $m(\mathbf{x})$ .

It is also useful to observe that since  $\boldsymbol{\beta}_K$  is the best linear approximation to  $m(\mathbf{x}_i)$  in mean-square (see Section 2.24), then

$$\begin{aligned} \mathbb{E}(r_{Ki}^2) &= \mathbb{E}\left((m(\mathbf{x}_i) - \mathbf{z}'_{Ki} \boldsymbol{\beta}_K)^2\right) \\ &\leq \mathbb{E}\left((m(\mathbf{x}_i) - \mathbf{z}'_{Ki} \boldsymbol{\beta}_K^*)^2\right) \\ &\leq O(K^{-2\alpha}) \end{aligned} \quad (18.17)$$

the final inequality by (18.14).

## 18.8 Residuals and Regression Fit

The fitted regression at  $\mathbf{x} = \mathbf{x}_i$  is  $\hat{m}_K(\mathbf{x}_i) = \mathbf{z}'_{Ki} \hat{\boldsymbol{\beta}}_K$  and the fitted residual is

$$\hat{e}_{iK} = y_i - \hat{m}_K(\mathbf{x}_i).$$

The leave-one-out prediction errors are

$$\begin{aligned} \tilde{e}_{iK} &= y_i - \hat{m}_{K,-i}(\mathbf{x}_i) \\ &= y_i - \mathbf{z}'_{Ki} \hat{\boldsymbol{\beta}}_{K,-i} \end{aligned}$$

where  $\hat{\boldsymbol{\beta}}_{K,-i}$  is the least-squares coefficient with the  $i$ 'th observation omitted. Using (3.44) we can also write

$$\tilde{e}_{iK} = \hat{e}_{iK}(1 - h_{Kii})^{-1}$$

where  $h_{Kii} = \mathbf{z}'_{Ki} (\mathbf{Z}'_K \mathbf{Z}_K)^{-1} \mathbf{z}_{Ki}$ .

As for kernel regression, the prediction errors  $\tilde{e}_{iK}$  are better estimates of the errors than the fitted residuals  $\hat{e}_{iK}$ , as they do not have the tendency to over-fit when the number of series terms is large.

To assess the fit of the nonparametric regression, the estimate of the mean-square prediction error is

$$\tilde{\sigma}_K^2 = \frac{1}{n} \sum_{i=1}^n \tilde{e}_{iK}^2 = \frac{1}{n} \sum_{i=1}^n \hat{e}_{iK}^2 (1 - h_{Kii})^{-2}$$

and the prediction  $R^2$  is

$$\tilde{R}_K^2 = 1 - \frac{\sum_{i=1}^n \tilde{e}_{iK}^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

## 18.9 Cross-Validation Model Selection

The cross-validation criterion for selection of the number of series terms is the MSPE

$$CV(K) = \tilde{\sigma}_K^2 = \frac{1}{n} \sum_{i=1}^n \tilde{e}_{iK}^2 (1 - h_{Kii})^{-2}.$$

By selecting the series terms to minimize  $CV(K)$ , or equivalently maximize  $\tilde{R}_K^2$ , we have a data-dependent rule which is designed to produce estimates with low integrated mean-squared error (IMSE) and mean-squared forecast error (MSFE). As shown in Theorem 17.6.1,  $CV(K)$  is an approximately unbiased estimated of the MSFE and IMSE, so finding the model which produces the smallest value of  $CV(K)$  is a good indicator that the estimated model has small MSFE and IMSE. The proof of the result is the same for all nonparametric estimators (series as well as kernels) so does not need to be repeated here.

As a practical matter, an estimator corresponds to a set of regressors  $\mathbf{z}_{Ki}$ , that is, a set of transformations of the original variables  $\mathbf{x}_i$ . For each set of regressions, the regression is estimated and  $CV(K)$  calculated, and the estimator is selected which has the smallest value of  $CV(K)$ . If there are  $p$  ordered regressors, then there are  $p$  possible estimators. Typically, this calculation is simple even if  $p$  is large. However, if the  $p$  regressors are unordered (and this is typical) then there are  $2^p$  possible subsets of conceivable models. If  $p$  is even moderately large,  $2^p$  can be immensely large so brute-force computation of all models may be computationally demanding.

## 18.10 Convergence in Mean-Square

The series estimate  $\hat{\beta}_K$  are indexed by  $K$ . The point of nonparametric estimation is to let  $K$  be flexible so as to incorporate greater complexity when the data are sufficiently informative. This means that  $K$  will typically be increasing with sample size  $n$ . This invalidates conventional asymptotic distribution theory. However, we can develop extensions which use appropriate matrix norms, and by focusing on real-valued functions of the parameters including the estimated regression function itself.

The asymptotic theory we present in this and the next several sections is largely taken from Newey (1997).

Our first main result shows that the least-squares estimate converges to  $\beta_K$  in mean-square distance.

**Theorem 18.10.1** *Under Assumption 18.7.1, as  $n \rightarrow \infty$ ,*

$$\left(\hat{\beta}_K - \beta_K\right)' \mathbf{Q}_K \left(\hat{\beta}_K - \beta_K\right) = O_p\left(\frac{K}{n}\right) + o_p(K^{-2\alpha}) \quad (18.18)$$

The proof of Theorem 18.10.1 is rather technical and deferred to Section 18.16.

The rate of convergence in (18.18) has two terms. The  $O_p(K/n)$  term is due to estimation variance. Note in contrast that the corresponding rate would be  $O_p(1/n)$  in the parametric case. The difference is that in the parametric case we assume that the number of regressors  $K$  is fixed as  $n$  increases, while in the nonparametric case we allow the number of regressors  $K$  to be flexible. As  $K$  increases, the estimation variance increases. The  $o_p(K^{-2\alpha})$  term in (18.18) is due to the series approximation error.

Using Theorem 18.10.1 we can establish the following convergence rate for the estimated regression function.

**Theorem 18.10.2** *Under Assumption 18.7.1, as  $n \rightarrow \infty$ ,*

$$\int (\hat{m}_K(\mathbf{x}) - m(\mathbf{x}))^2 f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} = O_p\left(\frac{K}{n}\right) + O_p(K^{-2\alpha}) \quad (18.19)$$

Theorem 18.10.2 shows that the integrated squared difference between the fitted regression and the true CEF converges in probability to zero if  $K \rightarrow \infty$  as  $n \rightarrow \infty$ . The convergence results of Theorem 18.10.2 show that the number of series terms  $K$  involves a trade-off similar to the role of the bandwidth  $h$  in kernel regression. Larger  $K$  implies smaller approximation error but increased estimation variance.

The optimal rate which minimizes the average squared error in (18.19) is  $K = O(n^{1/(1+2\alpha)})$ , yielding an optimal rate of convergence in (18.19) of  $O_p(n^{-2\alpha/(1+2\alpha)})$ . This rate depends on the unknown smoothness  $\alpha$  of the true CEF (the number of derivatives  $s$ ) and so does not directly suggest a practical rule for determining  $K$ . Still, the implication is that when the function being estimated is less smooth ( $\alpha$  is small) then it is necessary to use a larger number of series terms  $K$  to reduce the bias. In contrast, when the function is more smooth then it is better to use a smaller number of series terms  $K$  to reduce the variance.

To establish (18.19), using (18.7) and (18.8) we can write

$$\hat{m}_K(\mathbf{x}) - m(\mathbf{x}) = \mathbf{z}_K(\mathbf{x})' \left(\hat{\beta}_K - \beta_K\right) - r_K(\mathbf{x}). \quad (18.20)$$

Since  $e_{Ki}$  are projection errors, they satisfy  $\mathbb{E}(\mathbf{z}_{Ki}e_{Ki}) = 0$  and thus  $\mathbb{E}(\mathbf{z}_{Ki}r_{Ki}) = 0$ . This means  $\int \mathbf{z}_K(\mathbf{x})r_K(\mathbf{x})f_x(\mathbf{x})d\mathbf{x} = 0$ . Also observe that  $\mathbf{Q}_K = \int \mathbf{z}_K(\mathbf{x})\mathbf{z}_K(\mathbf{x})'f_x(\mathbf{x})d\mathbf{x}$  and  $\mathbb{E}(r_{Ki}^2) = \int r_K(\mathbf{x})^2f_x(\mathbf{x})d\mathbf{x}$ . Then

$$\begin{aligned} & \int (\hat{m}_K(\mathbf{x}) - m(\mathbf{x}))^2 f_x(\mathbf{x})d\mathbf{x} \\ &= (\hat{\beta}_K - \beta_K)' \mathbf{Q}_K (\hat{\beta}_K - \beta_K) + \mathbb{E}(r_{Ki}^2) \\ &\leq O_p\left(\frac{K}{n}\right) + O_p(K^{-2\alpha}) \end{aligned}$$

by (18.18) and (18.17), establishing (18.19).

## 18.11 Uniform Convergence

Theorem 18.10.2 established conditions under which  $\hat{m}_K(\mathbf{x})$  is consistent in a squared error norm. It is also of interest to know the rate at which the largest deviation converges to zero. We have the following rate.

**Theorem 18.11.1** *Under Assumption 18.7.1, then as  $n \rightarrow \infty$ ,*

$$\sup_{\mathbf{x} \in \mathcal{X}} |\hat{m}_K(\mathbf{x}) - m(\mathbf{x})| = O_p\left(\sqrt{\frac{\zeta_K^2 K}{n}}\right) + O_p(\zeta_K K^{-\alpha}). \quad (18.21)$$

Relative to Theorem 18.10.2, the error has been increased multiplicatively by  $\zeta_K$ . This slower convergence rate is a penalty for the stronger uniform convergence, though it is probably not the best possible rate. Examining the bound in (18.21) notice that the first term is  $o_p(1)$  under Assumption 18.7.1.4. The second term is  $o_p(1)$  if  $\zeta_K K^{-\alpha} \rightarrow 0$ , which requires that  $K \rightarrow \infty$  and that  $\alpha$  be sufficiently large. A sufficient condition is that  $s > d$  for polynomials and  $s > d/2$  for splines, where  $d = \dim(\mathbf{x})$  and  $s$  is the number of continuous derivatives of  $m(\mathbf{x})$ . Thus higher dimensional  $\mathbf{x}$  require a smoother CEF  $m(\mathbf{x})$  to ensure that the series estimate  $\hat{m}_K(\mathbf{x})$  is uniformly consistent.

The convergence (18.21) is straightforward to show using (18.18). Using (18.20), the Triangle Inequality, the Schwarz inequality (A.20), Definition (18.11), (18.18) and (18.16),

$$\begin{aligned} & \sup_{\mathbf{x} \in \mathcal{X}} |\hat{m}_K(\mathbf{x}) - m(\mathbf{x})| \\ &\leq \sup_{\mathbf{x} \in \mathcal{X}} \left| \mathbf{z}_K(\mathbf{x})' (\hat{\beta}_K - \beta_K) \right| + \sup_{\mathbf{x} \in \mathcal{X}} |r_K(\mathbf{x})| \\ &\leq \sup_{\mathbf{x} \in \mathcal{X}} \left( \mathbf{z}_K(\mathbf{x})' \mathbf{Q}_K^{-1} \mathbf{z}_K(\mathbf{x}) \right)^{1/2} \left( (\hat{\beta}_K - \beta_K)' \mathbf{Q}_K (\hat{\beta}_K - \beta_K) \right)^{1/2} \\ &\quad + O(\zeta_K K^{-\alpha}) \\ &\leq \zeta_K \left( O_p\left(\frac{K}{n}\right) + O_p(K^{-2\alpha}) \right)^{1/2} + O(\zeta_K K^{-\alpha}), \\ &= O_p\left(\sqrt{\frac{\zeta_K^2 K}{n}}\right) + O_p(\zeta_K K^{-\alpha}). \end{aligned} \quad (18.22)$$

This is (18.21).

## 18.12 Asymptotic Normality

One advantage of series methods is that the estimators are (in finite samples) equivalent to parametric estimators, so it is easy to calculate covariance matrix estimates. We now show that we can also justify normal asymptotic approximations.

The theory we present in this section will apply to any linear function of the regression function. That is, we allow the parameter of interest to be any non-trivial real-valued linear function of the entire regression function  $m(\cdot)$

$$\theta = a(m).$$

This includes the regression function  $m(\mathbf{x})$  at a given point  $\mathbf{x}$ , derivatives of  $m(\mathbf{x})$ , and integrals over  $m(\mathbf{x})$ . Given  $\hat{m}_K(\mathbf{x}) = \mathbf{z}_K(\mathbf{x})'\hat{\boldsymbol{\beta}}_K$  as an estimator for  $m(\mathbf{x})$ , the estimator for  $\theta$  is

$$\hat{\theta}_K = a(\hat{m}_K) = \mathbf{a}'_K \hat{\boldsymbol{\beta}}_K$$

for some  $K \times 1$  vector of constants  $\mathbf{a}_K \neq \mathbf{0}$ . (The relationship  $a(\hat{m}_K) = \mathbf{a}'_K \hat{\boldsymbol{\beta}}_K$  follows since  $a$  is linear in  $m$  and  $\hat{m}_K$  is linear in  $\hat{\boldsymbol{\beta}}_K$ .)

If  $K$  were fixed as  $n \rightarrow \infty$ , then by standard asymptotic theory we would expect  $\hat{\theta}_K$  to be asymptotically normal with variance

$$v_K = \mathbf{a}'_K \mathbf{Q}_K^{-1} \boldsymbol{\Omega}_K \mathbf{Q}_K^{-1} \mathbf{a}_K$$

where

$$\boldsymbol{\Omega}_K = \mathbb{E}(\mathbf{z}_{Ki} \mathbf{z}'_{Ki} e_{Ki}^2).$$

The standard justification, however, is not valid in the nonparametric case, in part because  $v_K$  may diverge as  $K \rightarrow \infty$ , and in part due to the finite sample bias due to the approximation error. Therefore a new theory is required. Interestingly, it turns out that in the nonparametric case  $\hat{\theta}_K$  is still asymptotically normal, and  $v_K$  is still the appropriate variance for  $\hat{\theta}_K$ . The proof is different than the parametric case as the dimensions of the matrices are increasing with  $K$ , and we need to be attentive to the estimator's bias due to the series approximation.

**Theorem 18.12.1** *Under Assumption 18.7.1, if in addition  $\mathbb{E}(e_i^4 | \mathbf{x}_i) \leq \kappa_4 < \infty$ ,  $\mathbb{E}(e_i^2 | \mathbf{x}_i) \geq \underline{\sigma}^2 > 0$ , and  $\zeta_K K^{-\alpha} = O(1)$ , then as  $n \rightarrow \infty$ ,*

$$\frac{\sqrt{n}(\hat{\theta}_K - \theta + a(r_K))}{v_K^{1/2}} \xrightarrow{d} N(0, 1) \quad (18.23)$$

The proof of Theorem 18.12.1 can be found in Section 18.16.

Theorem 18.12.1 shows that the estimator  $\hat{\theta}_K$  is approximately normal with bias  $-a(r_K)$  and variance  $v_K/n$ . The variance is the same as in the parametric case, but the asymptotic distribution contains an asymptotic bias, similar as is found in kernel regression. We discuss the bias in more detail below.

Notice that Theorem 18.12.1 requires  $\zeta_K K^{-\alpha} = O(1)$ , which is similar to that found in Theorem 18.11.1 to establish uniform convergence. The bound  $\zeta_K K^{-\alpha} = O(1)$  allows  $K$  to be constant with  $n$  or to increase with  $n$ . However, when  $K$  is increasing the bound requires that  $\alpha$  be sufficient large so that  $K^\alpha$  grows faster than  $\zeta_K$ . A sufficient condition is that  $s = d$  for polynomials and  $s = d/2$  for splines. The fact that the condition allows for  $K$  to be constant means that Theorem 18.12.1 includes parametric least-squares as a special case with explicit attention to estimation bias.



One useful message from Theorem 18.12.1 is that the classic variance formula  $v_K$  for  $\hat{\theta}_K$  still applies for series regression. Indeed, we can estimate the asymptotic variance using the standard White formula

$$\begin{aligned}\hat{v}_K &= \mathbf{a}'_K \hat{\mathbf{Q}}_K^{-1} \hat{\mathbf{\Omega}}_K \hat{\mathbf{Q}}_K^{-1} \mathbf{a}_K \\ \hat{\mathbf{\Omega}}_K &= \frac{1}{n} \sum_{i=1}^n \mathbf{z}_{Ki} \mathbf{z}'_{Ki} \hat{e}_{iK}^2 \\ \hat{\mathbf{Q}}_K &= \frac{1}{n} \sum_{i=1}^n \mathbf{z}_{Ki} \mathbf{z}'_{Ki}.\end{aligned}$$

Hence a standard error for  $\hat{\theta}_K$  is

$$\hat{s}(\theta_K) = \sqrt{\frac{1}{n} \mathbf{a}'_K \hat{\mathbf{Q}}_K^{-1} \hat{\mathbf{\Omega}}_K \hat{\mathbf{Q}}_K^{-1} \mathbf{a}_K}.$$

It can be shown (Newey, 1997) that  $\hat{v}_K/v_K \xrightarrow{p} 1$  as  $n \rightarrow \infty$  and thus the distribution in (18.23) is unchanged if  $v_K$  is replaced with  $\hat{v}_K$ .

Theorem 18.12.1 shows that the estimator  $\hat{\theta}_K$  has a bias term  $a(r_K)$ . What is this? It is the same transformation of the function  $r_K(\mathbf{x})$  as  $\theta = a(m)$  is of the regression function  $m(\mathbf{x})$ . For example, if  $\theta = m(\mathbf{x})$  is the regression at a fixed point  $\mathbf{x}$ , then  $a(r_K) = r_K(\mathbf{x})$ , the approximation error at the same point. If  $\theta = \frac{d}{dx}m(x)$  is the regression derivative, then  $a(r_K) = \frac{d}{dx}r_K(\mathbf{x})$  is the derivative of the approximation error.

This means that the bias in the estimator  $\hat{\theta}_K$  for  $\theta$  shown in Theorem 18.12.1 is simply the approximation error, transformed by the functional of interest. If we are estimating the regression function then the bias is the error in approximating the regression function; if we are estimating the regression derivative then the bias is the error in the derivative in the approximation error for the regression function.

### 18.13 Asymptotic Normality with Undersmoothing

An unpleasant aspect about Theorem 18.12.1 is the bias term. An interesting trick is that this bias term can be made asymptotically negligible if we assume that  $K$  increases with  $n$  at a sufficiently fast rate.

**Theorem 18.13.1** *Under Assumption 18.7.1, if in addition  $\mathbb{E}(e_i^4|\mathbf{x}_i) \leq \kappa_4 < \infty$ ,  $\mathbb{E}(e_i^2|\mathbf{x}_i) \geq \underline{\sigma}^2 > 0$ ,  $a(r_K^*) \leq O(K^{-\alpha})$ ,  $nK^{-2\alpha} \rightarrow 0$ , and  $\mathbf{a}'_K \mathbf{Q}_K^{-1} \mathbf{a}_K$  is bounded away from zero, then*

$$\frac{\sqrt{n}(\hat{\theta}_K - \theta)}{v_K^{1/2}} \xrightarrow{d} N(0, 1). \quad (18.24)$$

The condition  $a(r_K^*) \leq O(K^{-\alpha})$  states that the function of interest (for example, the regression function, its derivative, or its integral) applied to the uniform approximation error converges to zero as the number of terms  $K$  in the series approximation increases. If  $a(m) = m(\mathbf{x})$  then this condition holds by (18.6).

The condition that  $\mathbf{a}'_K \mathbf{Q}_K^{-1} \mathbf{a}_K$  is bounded away from zero is simply a technical requirement to exclude degeneracy.

The critical condition is the assumption that  $nK^{-2\alpha} \rightarrow 0$ . This requires that  $K \rightarrow \infty$  at a rate *faster* than  $n^{1/2\alpha}$ . This is a troubling condition. The optimal rate for estimation of  $m(\mathbf{x})$  is  $K = O(n^{1/(1+2\alpha)})$ . If we set  $K = n^{1/(1+2\alpha)}$  by this rule then  $nK^{-2\alpha} = n^{1/(1+2\alpha)} \rightarrow \infty$ , not zero. Thus this assumption is equivalent to assuming that  $K$  is much larger than optimal. The reason why this trick works (that is, why the bias is negligible) is that by increasing  $K$ , the asymptotic bias decreases and the asymptotic variance increases and thus the variance dominates. Because  $K$  is larger than optimal, we typically say that  $\hat{m}_K(\mathbf{x})$  is **undersmoothed** relative to the optimal series estimator.

Many authors like to focus their asymptotic theory on the assumptions in Theorem 18.13.1, as the distribution (18.24) appears cleaner. However, it is a poor use of asymptotic theory. There are three problems with the assumption  $nK^{-2\alpha} \rightarrow 0$  and the approximation (18.24). First, it says that if we intentionally pick  $K$  to be larger than optimal, we can increase the estimation variance relative to the bias so the variance will dominate the bias. But why would we want to intentionally use an estimator which is sub-optimal? Second, the assumption  $nK^{-2\alpha} \rightarrow 0$  does not eliminate the asymptotic bias, it only makes it of lower order than the variance. So the approximation (18.24) is technically valid, but the missing asymptotic bias term is just slightly smaller in asymptotic order, and thus still relevant in finite samples. Third, the condition  $nK^{-2\alpha} \rightarrow 0$  is just an assumption, it has nothing to do with actual empirical practice. Thus the difference between (18.23) and (18.24) is in the assumptions, not in the actual reality or in the actual empirical practice. Eliminating a nuisance (the asymptotic bias) through an assumption is a trick, not a substantive use of theory. My strong view is that the result (18.23) is more informative than (18.24). It shows that the asymptotic distribution is normal but has a non-trivial finite sample bias.

## 18.14 Regression Estimation

A special yet important example of a linear estimator of the regression function is the regression function at a fixed point  $\mathbf{x}$ . In the notation of the previous section,  $a(m) = m(\mathbf{x})$  and  $\mathbf{a}_K = \mathbf{z}_K(\mathbf{x})$ . The series estimator of  $m(\mathbf{x})$  is  $\hat{\theta}_K = \hat{m}_K(\mathbf{x}) = \mathbf{z}_K(\mathbf{x})' \hat{\beta}_K$ . As this is a key problem of interest, we restate the asymptotic results of Theorems 18.12.1 and 18.13.1 for this estimator.

**Theorem 18.14.1** *Under Assumption 18.7.1, if in addition  $\mathbb{E}(e_i^4 | \mathbf{x}_i) \leq \kappa_4 < \infty$ ,  $\mathbb{E}(e_i^2 | \mathbf{x}_i) \geq \underline{\sigma}^2 > 0$ , and  $\zeta_K K^{-\alpha} = O(1)$ , then as  $n \rightarrow \infty$ ,*

$$\frac{\sqrt{n}(\hat{m}_K(\mathbf{x}) - m(\mathbf{x}) + \mathbf{r}_K(\mathbf{x}))}{v_K^{1/2}(\mathbf{x})} \xrightarrow{d} N(0, 1) \quad (18.25)$$

where

$$v_K(\mathbf{x}) = \mathbf{z}_K(\mathbf{x})' \mathbf{Q}_K^{-1} \mathbf{\Omega}_K \mathbf{Q}_K^{-1} \mathbf{z}_K(\mathbf{x}).$$

If  $\zeta_K K^{-\alpha} = O(1)$  is replaced by  $nK^{-2\alpha} \rightarrow 0$ , and  $\mathbf{z}_K(\mathbf{x})' \mathbf{Q}_K^{-1} \mathbf{z}_K(\mathbf{x})$  is bounded away from zero, then

$$\frac{\sqrt{n}(\hat{m}_K(\mathbf{x}) - m(\mathbf{x}))}{v_K^{1/2}(\mathbf{x})} \xrightarrow{d} N(0, 1) \quad (18.26)$$

There are two important features about the asymptotic distribution (18.25).

First, as mentioned in the previous section, it shows how to construct asymptotic standard errors for the CEF  $m(\mathbf{x})$ . These are

$$\hat{s}(\mathbf{x}) = \sqrt{\frac{1}{n} \mathbf{z}_K(\mathbf{x})' \hat{\mathbf{Q}}_K^{-1} \hat{\mathbf{\Omega}}_K \hat{\mathbf{Q}}_K^{-1} \mathbf{z}_K(\mathbf{x})}.$$

Second, (18.25) shows that the estimator has the asymptotic bias component  $\mathbf{r}_K(\mathbf{x})$ . This is due to the fact that the finite order series is an approximation to the unknown CEF  $m(\mathbf{x})$ , and this results in finite sample bias.

The asymptotic distribution (18.26) shows that the bias term is negligible if  $K$  diverges fast enough so that  $nK^{-2\alpha} \rightarrow 0$ . As discussed in the previous section, this means that  $K$  is larger than optimal.

The assumption that  $\mathbf{z}_K(\mathbf{x})' \mathbf{Q}_K^{-1} \mathbf{z}_K(\mathbf{x})$  is bounded away from zero is a technical condition to exclude degenerate cases, and is automatically satisfied if  $\mathbf{z}_K(\mathbf{x})$  includes an intercept.

Plots of the CEF estimate  $\hat{m}_K(\mathbf{x})$  can be accompanied by 95% confidence intervals  $\hat{m}_K(\mathbf{x}) \pm 2\hat{s}(\mathbf{x})$ . As we discussed in the chapter on kernel regression, this can be viewed as a confidence interval for the pseudo-true CEF  $m_K^*(\mathbf{x}) = m(\mathbf{x}) - \mathbf{r}_K(\mathbf{x})$ , not for the true  $m(\mathbf{x})$ . As for kernel regression, the difference is the unavoidable consequence of nonparametric estimation.

## 18.15 Kernel Versus Series Regression

In this and the previous chapter we have presented two distinct methods of nonparametric regression based on kernel methods and series methods. Which should be used in practice? Both methods have advantages and disadvantages and there is no clear overall winner.

First, while the asymptotic theory of the two estimators appear quite different, they are actually rather closely related. When the regression function  $m(\mathbf{x})$  is twice differentiable ( $s = 2$ ) then the rate of convergence of both the MSE of the kernel regression estimator with optimal bandwidth  $h$  and the series estimator with optimal  $K$  is  $n^{-2/(d+4)}$ . There is no difference. If the regression function is smoother than twice differentiable ( $s > 2$ ) then the rate of the convergence of the series estimator improves. This may appear to be an advantage for series methods, but kernel regression can also take advantage of the higher smoothness by using so-called higher-order kernels or local polynomial regression, so perhaps this advantage is not too large.

Both estimators are asymptotically normal and have straightforward asymptotic standard error formulae. The series estimators are a bit more convenient for this purpose, as classic parametric standard error formula work without amendment.

An advantage of kernel methods is that their distributional theory is easier to derive. The theory is all based on local averages which is relatively straightforward. In contrast, series theory is more challenging, dealing with increasing parameter spaces. An important difference in the theory is that for kernel estimators we have explicit representations for the bias while we only have rates for series methods. This means that plug-in methods can be used for bandwidth selection in kernel regression. However, typically we rely on cross-validation, which is equally applicable in both kernel and series regression.

Kernel methods are also relatively easy to implement when the dimension  $d$  is large. There is not a major change in the methodology as  $d$  increases. In contrast, series methods become quite cumbersome as  $d$  increases as the number of cross-terms increases exponentially.

A major advantage of series methods is that it has inherently a high degree of flexibility, and the user is able to implement shape restrictions quite easily. For example, in series estimation it is relatively simple to implement a partial linear CEF, an additively separable CEF, monotonicity, concavity or convexity. These restrictions are harder to implement in kernel regression.

## 18.16 Technical Proofs

Define  $\mathbf{z}_{Ki} = \mathbf{z}_K(\mathbf{x}_i)$  and let  $\mathbf{Q}_K^{1/2}$  denote the positive definite square root of  $\mathbf{Q}_K$ . As mentioned before Theorem 18.10.1, the regression problem is unchanged if we replace  $\mathbf{z}_{Ki}$  with a rotated regressor such as  $\mathbf{z}_{Ki}^* = \mathbf{Q}_K^{-1/2} \mathbf{z}_{Ki}$ . This is a convenient choice for then  $\mathbb{E}(\mathbf{z}_{Ki}^* \mathbf{z}_{Ki}^{*'}) = \mathbf{I}_K$ . For notational convenience we will simply write the transformed regressors as  $\mathbf{z}_{Ki}$  and set  $\mathbf{Q}_K = \mathbf{I}_K$ .

We start with some convergence results for the sample design matrix

$$\hat{\mathbf{Q}}_K = \frac{1}{n} \mathbf{Z}'_K \mathbf{Z}_K = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_{Ki} \mathbf{z}'_{Ki}.$$

**Theorem 18.16.1** *Under Assumption 18.7.1 and  $\mathbf{Q}_K = \mathbf{I}_K$ , as  $n \rightarrow \infty$ ,*

$$\left\| \hat{\mathbf{Q}}_K - \mathbf{I}_K \right\| = o_p(1) \quad (18.27)$$

and

$$\lambda_{\min}(\hat{\mathbf{Q}}_K) \xrightarrow{p} 1. \quad (18.28)$$

**Proof.** Since

$$\left\| \hat{\mathbf{Q}}_K - \mathbf{I}_K \right\|^2 = \sum_{j=1}^K \sum_{\ell=1}^K \left( \frac{1}{n} \sum_{i=1}^n (\mathbf{z}_{jKi} \mathbf{z}_{\ell Ki} - \mathbb{E}(\mathbf{z}_{jKi} \mathbf{z}_{\ell Ki})) \right)^2$$

then

$$\begin{aligned} \mathbb{E} \left( \left\| \hat{\mathbf{Q}}_K - \mathbf{I}_K \right\|^2 \right) &= \sum_{j=1}^K \sum_{\ell=1}^K \text{var} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{z}_{jKi} \mathbf{z}_{\ell Ki} \right) \\ &= n^{-1} \sum_{j=1}^K \sum_{\ell=1}^K \text{var}(\mathbf{z}_{jKi} \mathbf{z}_{\ell Ki}) \\ &\leq n^{-1} \mathbb{E} \left( \sum_{j=1}^K \mathbf{z}_{jKi}^2 \sum_{\ell=1}^K \mathbf{z}_{\ell Ki}^2 \right) \\ &= n^{-1} \mathbb{E} \left( (\mathbf{z}'_{Ki} \mathbf{z}_{Ki})^2 \right). \end{aligned} \quad (18.29)$$

Since  $\mathbf{z}'_{Ki} \mathbf{z}_{Ki} \leq \zeta_K^2$  by definition (18.11) and using (A.1) we find

$$\mathbb{E}(\mathbf{z}'_{Ki} \mathbf{z}_{Ki}) = \text{tr}(\mathbb{E}(\mathbf{z}_{Ki} \mathbf{z}'_{Ki})) = \text{tr} \mathbf{I}_K = K, \quad (18.30)$$

so that

$$\mathbb{E} \left( (\mathbf{z}'_{Ki} \mathbf{z}_{Ki})^2 \right) \leq \zeta_K^2 K \quad (18.31)$$

and hence (18.29) is  $o(1)$  under Assumption 18.7.1.4. Theorem 6.13.1 shows that this implies (18.27).

Let  $\lambda_1, \lambda_2, \dots, \lambda_K$  be the eigenvalues of  $\hat{\mathbf{Q}}_K - \mathbf{I}_K$  which are real as  $\hat{\mathbf{Q}}_K - \mathbf{I}_K$  is symmetric. Then

$$\left| \lambda_{\min}(\hat{\mathbf{Q}}_K) - 1 \right| = \left| \lambda_{\min}(\hat{\mathbf{Q}}_K - \mathbf{I}_K) \right| \leq \left( \sum_{\ell=1}^K \lambda_{\ell}^2 \right)^{1/2} = \left\| \hat{\mathbf{Q}}_K - \mathbf{I}_K \right\|$$

where the second equality is (A.22). This is  $o_p(1)$  by (18.27), establishing (18.28).  $\blacksquare$

**Proof of Theorem 18.10.1.** As above, assume that the regressors have been transformed so that  $\mathbf{Q}_K = \mathbf{I}_K$ .

From expression (18.10) we can substitute to find

$$\begin{aligned}\widehat{\beta}_K - \beta_K &= (\mathbf{Z}'_K \mathbf{Z}_K)^{-1} \mathbf{Z}'_K \mathbf{e}_K. \\ &= \widehat{\mathbf{Q}}_K^{-1} \left( \frac{1}{n} \mathbf{Z}'_K \mathbf{e}_K \right)\end{aligned}\tag{18.32}$$

Using (18.32) and the Quadratic Inequality (A.28),

$$\begin{aligned}& (\widehat{\beta}_K - \beta_K)' (\widehat{\beta}_K - \beta_K) \\ &= n^{-2} (\mathbf{e}'_K \mathbf{Z}_K) \widehat{\mathbf{Q}}_K^{-1} \widehat{\mathbf{Q}}_K^{-1} (\mathbf{Z}'_K \mathbf{e}_K) \\ &\leq \left( \lambda_{\max} \left( \widehat{\mathbf{Q}}_K^{-1} \right) \right)^2 n^{-2} (\mathbf{e}'_K \mathbf{Z}_K \mathbf{Z}'_K \mathbf{e}_K).\end{aligned}\tag{18.33}$$

Observe that (18.28) implies

$$\lambda_{\max} \left( \widehat{\mathbf{Q}}_K^{-1} \right) = \left( \lambda_{\min} \left( \widehat{\mathbf{Q}}_K \right) \right)^{-1} = O_p(1).\tag{18.34}$$

Since  $e_{Ki} = e_i + r_{Ki}$ , and using Assumption 18.7.1.2 and (18.16), then

$$\sup_i \mathbb{E} (e_{Ki}^2 | \mathbf{x}_i) = \bar{\sigma}^2 + \sup_i r_{Ki}^2 \leq \bar{\sigma}^2 + O(\zeta_K^2 K^{-2\alpha}).\tag{18.35}$$

As  $e_{Ki}$  are projection errors, they satisfy  $\mathbb{E}(\mathbf{z}_{Ki} e_{Ki}) = 0$ . Since the observations are independent, using (18.30) and (18.35), then

$$\begin{aligned}n^{-2} \mathbb{E} (\mathbf{e}'_K \mathbf{Z}_K \mathbf{Z}'_K \mathbf{e}_K) &= n^{-2} \mathbb{E} \left( \sum_{i=1}^n e_{Ki} \mathbf{z}'_{Ki} \sum_{j=1}^n \mathbf{z}_{Kj} e_{Kj} \right) \\ &= n^{-2} \sum_{i=1}^n \mathbb{E} (\mathbf{z}'_{Ki} \mathbf{z}_{Ki} e_{Ki}^2) \\ &\leq n^{-1} \mathbb{E} (\mathbf{z}'_{Ki} \mathbf{z}_{Ki}) \sup_i \mathbb{E} (e_{Ki}^2 | \mathbf{x}_i) \\ &\leq \bar{\sigma}^2 \frac{K}{n} + O \left( \frac{\zeta_K^2 K^{1-2\alpha}}{n} \right) \\ &= \bar{\sigma}^2 \frac{K}{n} + o(K^{-2\alpha})\end{aligned}\tag{18.36}$$

since  $\zeta_K^2 K/n = o(1)$  by Assumption 18.7.1.4. Theorem 6.13.1 shows that this implies

$$n^{-2} \mathbf{e}'_K \mathbf{Z}_K \mathbf{Z}'_K \mathbf{e}_K = O_p(n^{-2}) + o_p(K^{-2\alpha}).\tag{18.37}$$

Together, (18.33), (18.34) and (18.37) imply (18.18).  $\blacksquare$

**Proof of Theorem 18.12.1.** As above, assume that the regressors have been transformed so that  $\mathbf{Q}_K = \mathbf{I}_K$ .

Using  $m(\mathbf{x}) = \mathbf{z}_K(\mathbf{x})' \beta_K + r_K(\mathbf{x})$  and linearity

$$\begin{aligned}\theta &= a(m) \\ &= a(\mathbf{z}_K(\mathbf{x})' \beta_K) + a(r_K) \\ &= \mathbf{a}'_K \beta_K + a(r_K)\end{aligned}$$

Combined with (18.32) we find

$$\begin{aligned}\widehat{\theta}_K - \theta + a(r_K) &= \mathbf{a}'_K (\widehat{\beta}_K - \beta_K) \\ &= \frac{1}{n} \mathbf{a}'_K \widehat{\mathbf{Q}}_K^{-1} \mathbf{Z}'_K \mathbf{e}_K\end{aligned}$$

and thus

$$\begin{aligned}\sqrt{\frac{n}{v_k}} (\widehat{\theta}_K - \theta_K + a(r_K)) &= \sqrt{\frac{n}{v_k}} \mathbf{a}'_K (\widehat{\beta}_K - \beta_K) \\ &= \sqrt{\frac{1}{nv_k}} \mathbf{a}'_K \widehat{\mathbf{Q}}_K^{-1} \mathbf{Z}'_K \mathbf{e}_K \\ &= \frac{1}{\sqrt{nv_K}} \mathbf{a}'_K \mathbf{Z}'_K \mathbf{e}_K\end{aligned}\tag{18.38}$$

$$+ \frac{1}{\sqrt{nv_K}} \mathbf{a}'_K (\widehat{\mathbf{Q}}_K^{-1} - \mathbf{I}_K) \mathbf{Z}'_K \mathbf{e}\tag{18.39}$$

$$+ \frac{1}{\sqrt{nv_K}} \mathbf{a}'_K (\widehat{\mathbf{Q}}_K^{-1} - \mathbf{I}_K) \mathbf{Z}'_K \mathbf{r}_K.\tag{18.40}$$

where we have used  $\mathbf{e}_K = \mathbf{e} + \mathbf{r}_K$ . We now take the terms in (18.38)-(18.40) separately.

First, take (18.38). We can write

$$\frac{1}{\sqrt{nv_K}} \mathbf{a}'_K \mathbf{Z}'_K \mathbf{e}_K = \frac{1}{\sqrt{nv_K}} \sum_{i=1}^n \mathbf{a}'_K \mathbf{z}_{Ki} e_{Ki}.\tag{18.41}$$

Observe that  $\mathbf{a}'_K \mathbf{z}_{Ki} e_{Ki}$  are independent across  $i$ , mean zero, and have variance

$$\mathbb{E} \left( (\mathbf{a}'_K \mathbf{z}_{Ki} e_{Ki})^2 \right) = \mathbf{a}'_K \mathbb{E} (\mathbf{z}_{Ki} \mathbf{z}'_{Ki} e_{Ki}^2) \mathbf{a}_K = v_K.$$

We will apply the Lindeberg CLT 6.8.2, for which it is sufficient to verify Lyapunov's condition (6.6):

$$\frac{1}{n^2 v_K^2} \sum_{i=1}^n \mathbb{E} \left( (\mathbf{a}'_K \mathbf{z}_{Ki} e_{Ki})^4 \right) = \frac{1}{nv_K^2} \mathbb{E} \left( (\mathbf{a}'_K \mathbf{z}_{Ki})^4 e_{Ki}^4 \right) \rightarrow 0.\tag{18.42}$$

The assumption that  $\zeta_K K^{-\alpha} = O(1)$  means  $\zeta_K K^{-\alpha} \leq \kappa_1$  for some  $\kappa_1 < \infty$ . Then by the  $c_r$  inequality and  $\mathbb{E}(e_i^4 | \mathbf{x}_i) \leq \kappa$

$$\sup_i \mathbb{E}(e_{Ki}^4 | \mathbf{x}_i) \leq 8 \sup_i (\mathbb{E}(e_i^4 | \mathbf{x}_i) + r_{Ki}^4) \leq 8(\kappa + \kappa_1).\tag{18.43}$$

Using (18.43), the Schwarz Inequality, and (18.31)

$$\begin{aligned}\mathbb{E} \left( (\mathbf{a}'_K \mathbf{z}_{Ki})^4 e_{Ki}^4 \right) &= \mathbb{E} \left( (\mathbf{a}'_K \mathbf{z}_{Ki})^4 \mathbb{E}(e_{Ki}^4 | \mathbf{x}_i) \right) \\ &\leq 8(\kappa + \kappa_1) \mathbb{E} \left( (\mathbf{a}'_K \mathbf{z}_{Ki})^4 \right) \\ &\leq 8(\kappa + \kappa_1) (\mathbf{a}'_K \mathbf{a}_K)^2 \mathbb{E} \left( (\mathbf{z}'_{Ki} \mathbf{z}_{Ki})^2 \right) \\ &= 8(\kappa + \kappa_1) (\mathbf{a}'_K \mathbf{a}_K)^2 \zeta_K^2 K.\end{aligned}\tag{18.44}$$

Since  $\mathbb{E}(e_{Ki}^2 | \mathbf{x}_i) = \mathbb{E}(e_i^2 | \mathbf{x}_i) + r_{Ki}^2 \geq \underline{\sigma}^2$ ,

$$\begin{aligned}v_K &= \mathbf{a}'_K \mathbb{E} (\mathbf{z}_{Ki} \mathbf{z}'_{Ki} e_{Ki}^2) \mathbf{a}_K \\ &\geq \underline{\sigma}^2 \mathbf{a}'_K \mathbb{E} (\mathbf{z}_{Ki} \mathbf{z}'_{Ki}) \mathbf{a}_K \\ &= \underline{\sigma}^2 \mathbf{a}'_K \mathbf{a}_K.\end{aligned}\tag{18.45}$$

Equation (18.44) and (18.45) combine to show that

$$\frac{1}{nv_K^2} \mathbb{E} \left( (a'_K z_{Ki})^4 e_{Ki}^4 \right) \leq \frac{8(\kappa + \kappa_1) \zeta_K^2 K}{\underline{\sigma}^4 n} = o(1)$$

under Assumption 18.7.1.4. This establishes Lyapunov's condition (18.42). Hence the Lindeberg CLT applies to (18.41) and we conclude

$$\frac{1}{\sqrt{nv_K}} a'_K \mathbf{Z}'_K \mathbf{e}_K \xrightarrow{d} N(0, 1). \quad (18.46)$$

Second, take (18.39). Since  $\mathbb{E}(\mathbf{e} \mid \mathbf{X}) = 0$ , then applying  $\mathbb{E}(e_i^2 \mid \mathbf{x}_i) \leq \bar{\sigma}^2$ , the Schwarz and Norm Inequalities, (18.45), (18.34) and (18.27),

$$\begin{aligned} & \mathbb{E} \left( \left( \frac{1}{\sqrt{nv_K}} a'_K (\hat{\mathbf{Q}}_K^{-1} - \mathbf{I}_K) \mathbf{Z}'_K \mathbf{e} \right)^2 \mid \mathbf{X} \right) \\ &= \frac{1}{nv_K} a'_K (\hat{\mathbf{Q}}_K^{-1} - \mathbf{I}_K) \mathbf{Z}'_K \mathbb{E}(\mathbf{e} \mathbf{e}' \mid \mathbf{X}) \mathbf{Z}_K (\hat{\mathbf{Q}}_K^{-1} - \mathbf{I}_K) a_K \\ &\leq \frac{\bar{\sigma}^2}{v_K} a'_K (\hat{\mathbf{Q}}_K^{-1} - \mathbf{I}_K) \hat{\mathbf{Q}}_K (\hat{\mathbf{Q}}_K^{-1} - \mathbf{I}_K) a_K \\ &= \frac{\bar{\sigma}^2}{v_K} a'_K (\hat{\mathbf{Q}}_K - \mathbf{I}_K) \hat{\mathbf{Q}}_K^{-1} (\hat{\mathbf{Q}}_K - \mathbf{I}_K) a_K \\ &\leq \frac{\bar{\sigma}^2 a'_K a_K}{v_K} \lambda_{\max}(\hat{\mathbf{Q}}_K^{-1}) \|\hat{\mathbf{Q}}_K - \mathbf{I}_K\|^2 \\ &\leq \frac{\bar{\sigma}^2}{\underline{\sigma}^2} o_p(1). \end{aligned}$$

This establishes

$$\frac{1}{\sqrt{nv_K}} a'_K (\hat{\mathbf{Q}}_K^{-1} - \mathbf{I}_K) \mathbf{Z}'_K \mathbf{e} \xrightarrow{p} 0. \quad (18.47)$$

Third, take (18.40). By the Cauchy-Schwarz inequality, (18.45), and the Quadratic Inequality,

$$\begin{aligned} & \left( \frac{1}{\sqrt{nv_K}} a'_K (\hat{\mathbf{Q}}_K^{-1} - \mathbf{I}_K) \mathbf{Z}'_K \mathbf{r}_K \right)^2 \\ &\leq \frac{a'_K a_K}{nv_K} \mathbf{r}'_K \mathbf{Z}_K (\hat{\mathbf{Q}}_K^{-1} - \mathbf{I}_K) (\hat{\mathbf{Q}}_K^{-1} - \mathbf{I}_K) \mathbf{Z}'_K \mathbf{r}_K \\ &\leq \frac{1}{\underline{\sigma}^2} \lambda_{\max}(\hat{\mathbf{Q}}_K^{-1} - \mathbf{I}_K)^2 \frac{1}{n} \mathbf{r}'_K \mathbf{Z}_K \mathbf{Z}'_K \mathbf{r}_K. \end{aligned} \quad (18.48)$$

Observe that since the observations are independent and  $\mathbb{E} \mathbf{z}_{Ki} \mathbf{r}_{Ki} = 0$ ,  $\mathbf{z}'_{Ki} \mathbf{z}_{Ki} \leq \zeta_K^2$ , and (18.17)

$$\begin{aligned} \mathbb{E} \left( \frac{1}{n} \mathbf{r}'_K \mathbf{Z}_K \mathbf{Z}'_K \mathbf{r}_K \right) &= \mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n r_{Ki} \mathbf{z}'_{Ki} \sum_{ij=1}^n \mathbf{z}_{Kj} r_{Kj} \right) \\ &= \mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{z}'_{Ki} \mathbf{z}_{Ki} r_{Ki}^2 \right) \\ &\leq \zeta_K^2 \mathbb{E}(r_{Ki}^2) \\ &= O(\zeta_K^2 K^{-2\alpha}) \\ &= O(1) \end{aligned}$$

since  $\zeta_K K^{-2} = O(1)$ . Thus  $\frac{1}{n} \mathbf{r}'_K \mathbf{Z}_K \mathbf{Z}'_K \mathbf{r}_K = O_p(1)$ . This means that (18.48) is  $o_p(1)$  since (18.28) implies

$$\lambda_{\max} \left( \widehat{\mathbf{Q}}_K^{-1} - \mathbf{I}_K \right) = \lambda_{\max} \left( \widehat{\mathbf{Q}}_K^{-1} \right) - 1 = o_p(1). \quad (18.49)$$

Equivalently,

$$\frac{1}{\sqrt{nv_K}} \mathbf{a}'_K \left( \widehat{\mathbf{Q}}_K^{-1} - \mathbf{I}_K \right) \mathbf{Z}'_K \mathbf{r}_K \xrightarrow{p} 0. \quad (18.50)$$

Equations (18.46), (18.47) and (18.50) applied to (18.38)-(18.40) show that

$$\sqrt{\frac{n}{v_k}} \left( \widehat{\theta}_K - \theta_K + a(r_K) \right) \xrightarrow{d} N(0, 1)$$

completing the proof. ■

**Proof of Theorem 18.13.1.** The assumption that  $nK^{-2\alpha} = o(1)$  implies  $K^{-\alpha} = o(n^{-1/2})$ . Thus

$$\zeta_K K^{-\alpha} \leq o \left( \left( \frac{\zeta_K^2}{n} \right)^{1/2} \right) \leq o \left( \left( \frac{\zeta_K^2 K}{n} \right)^{1/2} \right) = o(1)$$

so the conditions of Theorem 18.12.1 are satisfied. It is thus sufficient to show that

$$\sqrt{\frac{n}{v_k}} a(r_K) = o(1).$$

From (18.12)

$$\begin{aligned} r_K(\mathbf{x}) &= r_K^*(\mathbf{x}) + \mathbf{z}_K(\mathbf{x})' \gamma_K \\ \gamma_K &= \mathbb{E} \left( \mathbf{z}_{Ki} \mathbf{z}'_{Ki} \right)^{-1} \mathbb{E} \left( \mathbf{z}_{Ki} r_{Ki}^* \right). \end{aligned}$$

Thus by linearity, applying (18.45), and the Schwarz inequality

$$\begin{aligned} \sqrt{\frac{n}{v_k}} a(r_K) &= \sqrt{\frac{n}{v_k}} \left( a(r_K^*) + \mathbf{a}'_K \gamma_K \right) \\ &\leq \frac{n^{1/2}}{\underline{\sigma}^2 (\mathbf{a}'_K \mathbf{a}_K)^{1/2}} a(r_K^*) \end{aligned} \quad (18.51)$$

$$+ \frac{(n \gamma'_K \gamma_K)^{1/2}}{\underline{\sigma}}. \quad (18.52)$$

By assumption,  $n^{1/2} a(r_K^*) = O(n^{1/2} K^{-\alpha}) = o(1)$ . By (18.14) and  $nK^{-2\alpha} = o(1)$

$$\begin{aligned} n \gamma'_K \gamma_K &= n \mathbb{E} \left( r_{Ki}^* \mathbf{z}'_{Ki} \right) \mathbb{E} \left( \mathbf{z}_{Ki} \mathbf{z}'_{Ki} \right)^{-1} \mathbb{E} \left( \mathbf{z}_{Ki} r_{Ki}^* \right) \\ &\leq n O(K^{-2\alpha}) \\ &= o(1). \end{aligned}$$

Together, both (18.51) and (18.52) are  $o(1)$ , as required. ■



## Exercises

**Exercise 18.1** You have a friend who wants to estimate  $\beta$  in the model

$$y_i = x_i\beta + e_i$$
$$\mathbb{E}(e_i \mid z_i) = 0$$

with both  $x_i \in \mathbb{R}$  and  $z_i \in \mathbb{R}$ , and  $z_i$  is continuously distributed. Your friend wants to treat the reduced form equation for  $x_i$  as nonparametric

$$x_i = g(z_i) + u_i$$
$$\mathbb{E}(u_i \mid z_i) = 0$$

Your friend asks you for advice and help to construct an estimator  $\hat{\beta}$  of  $\beta$ . Describe an appropriate estimator. You do not have to develop the distribution theory, but try to be sufficiently complete with your advice so your friend can compute  $\hat{\beta}$ .

# Chapter 19

## Empirical Likelihood

### 19.1 Non-Parametric Likelihood

An alternative to GMM is **empirical likelihood**. The idea is due to Art Owen (1988, 2001) and has been extended to moment condition models by Qin and Lawless (1994). It is a non-parametric analog of likelihood estimation.

The idea is to construct a multinomial distribution  $F(p_1, \dots, p_n)$  which places probability  $p_i$  at each observation. To be a valid multinomial distribution, these probabilities must satisfy the requirements that  $p_i \geq 0$  and

$$\sum_{i=1}^n p_i = 1. \quad (19.1)$$

Since each observation is observed once in the sample, the log-likelihood function for this multinomial distribution is

$$\log L(p_1, \dots, p_n) = \sum_{i=1}^n \log(p_i). \quad (19.2)$$

First let us consider a just-identified model. In this case the moment condition places no additional restrictions on the multinomial distribution. The maximum likelihood estimators of the probabilities  $(p_1, \dots, p_n)$  are those which maximize the log-likelihood subject to the constraint (19.1). This is equivalent to maximizing

$$\sum_{i=1}^n \log(p_i) - \mu \left( \sum_{i=1}^n p_i - 1 \right)$$

where  $\mu$  is a Lagrange multiplier. The  $n$  first order conditions are  $0 = p_i^{-1} - \mu$ . Combined with the constraint (19.1) we find that the MLE is  $p_i = n^{-1}$  yielding the log-likelihood  $-n \log(n)$ .

Now consider the case of an overidentified model with moment condition

$$\mathbb{E}(\mathbf{g}_i(\boldsymbol{\beta})) = \mathbf{0}$$

where  $\mathbf{g}$  is  $\ell \times 1$  and  $\boldsymbol{\beta}$  is  $k \times 1$  and for simplicity we write  $\mathbf{g}_i(\boldsymbol{\beta}) = \mathbf{g}(y_i, \mathbf{z}_i, \mathbf{x}_i, \boldsymbol{\beta})$ . The multinomial distribution which places probability  $p_i$  at each observation  $(y_i, \mathbf{x}_i, \mathbf{z}_i)$  will satisfy this condition if and only if

$$\sum_{i=1}^n p_i \mathbf{g}_i(\boldsymbol{\beta}) = \mathbf{0} \quad (19.3)$$

The **empirical likelihood estimator** is the value of  $\boldsymbol{\beta}$  which maximizes the multinomial log-likelihood (19.2) subject to the restrictions (19.1) and (19.3).

The Lagrangian for this maximization problem is

$$\mathcal{L}(\boldsymbol{\beta}, p_1, \dots, p_n, \boldsymbol{\lambda}, \mu) = \sum_{i=1}^n \log(p_i) - \mu \left( \sum_{i=1}^n p_i - 1 \right) - n \boldsymbol{\lambda}' \sum_{i=1}^n p_i \mathbf{g}_i(\boldsymbol{\beta})$$

where  $\boldsymbol{\lambda}$  and  $\mu$  are Lagrange multipliers. The first-order-conditions of  $\mathcal{L}$  with respect to  $p_i$ ,  $\mu$ , and  $\boldsymbol{\lambda}$  are

$$\begin{aligned} \frac{1}{p_i} &= \mu + n \boldsymbol{\lambda}' \mathbf{g}_i(\boldsymbol{\beta}) \\ \sum_{i=1}^n p_i &= 1 \\ \sum_{i=1}^n p_i \mathbf{g}_i(\boldsymbol{\beta}) &= \mathbf{0}. \end{aligned}$$

Multiplying the first equation by  $p_i$ , summing over  $i$ , and using the second and third equations, we find  $\mu = n$  and

$$p_i = \frac{1}{n(1 + \boldsymbol{\lambda}' \mathbf{g}_i(\boldsymbol{\beta}))}.$$

Substituting into  $\mathcal{L}$  we find

$$R(\boldsymbol{\beta}, \boldsymbol{\lambda}) = -n \log(n) - \sum_{i=1}^n \log(1 + \boldsymbol{\lambda}' \mathbf{g}_i(\boldsymbol{\beta})). \quad (19.4)$$

For given  $\boldsymbol{\beta}$ , the Lagrange multiplier  $\boldsymbol{\lambda}(\boldsymbol{\beta})$  minimizes  $R(\boldsymbol{\beta}, \boldsymbol{\lambda})$ :

$$\boldsymbol{\lambda}(\boldsymbol{\beta}) = \underset{\boldsymbol{\lambda}}{\operatorname{argmin}} R(\boldsymbol{\beta}, \boldsymbol{\lambda}). \quad (19.5)$$

This minimization problem is the dual of the constrained maximization problem. The solution (when it exists) is well defined since  $R(\boldsymbol{\beta}, \boldsymbol{\lambda})$  is a convex function of  $\boldsymbol{\lambda}$ . The solution cannot be obtained explicitly, but must be obtained numerically (see section 6.5). This yields the (profile) empirical log-likelihood function for  $\boldsymbol{\beta}$ .

$$\begin{aligned} R(\boldsymbol{\beta}) &= R(\boldsymbol{\beta}, \boldsymbol{\lambda}(\boldsymbol{\beta})) \\ &= -n \log(n) - \sum_{i=1}^n \log(1 + \boldsymbol{\lambda}(\boldsymbol{\beta})' \mathbf{g}_i(\boldsymbol{\beta})) \end{aligned}$$

The EL estimate  $\hat{\boldsymbol{\beta}}$  is the value which maximizes  $R(\boldsymbol{\beta})$ , or equivalently minimizes its negative

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} [-R(\boldsymbol{\beta})] \quad (19.6)$$

Numerical methods are required for calculation of  $\hat{\boldsymbol{\beta}}$  (see Section 19.5).

As a by-product of estimation, we also obtain the Lagrange multiplier  $\hat{\boldsymbol{\lambda}} = \boldsymbol{\lambda}(\hat{\boldsymbol{\beta}})$ , probabilities

$$\hat{p}_i = \frac{1}{n(1 + \hat{\boldsymbol{\lambda}}' \mathbf{g}_i(\hat{\boldsymbol{\beta}}))}.$$

and maximized empirical likelihood

$$R(\hat{\boldsymbol{\beta}}) = \sum_{i=1}^n \log(\hat{p}_i). \quad (19.7)$$

## 19.2 Asymptotic Distribution of EL Estimator

Define

$$\mathbf{G}_i(\boldsymbol{\beta}) = \frac{\partial}{\partial \boldsymbol{\beta}'} \mathbf{g}_i(\boldsymbol{\beta}) \quad (19.8)$$

$$\mathbf{G} = \mathbb{E}(\mathbf{G}_i(\boldsymbol{\beta}))$$

$$\boldsymbol{\Omega} = \mathbb{E}(\mathbf{g}_i(\boldsymbol{\beta}) \mathbf{g}_i(\boldsymbol{\beta})')$$

and

$$\mathbf{V} = (\mathbf{G}'\boldsymbol{\Omega}^{-1}\mathbf{G})^{-1} \quad (19.9)$$

$$\mathbf{V}_\lambda = \boldsymbol{\Omega} - \mathbf{G}(\mathbf{G}'\boldsymbol{\Omega}^{-1}\mathbf{G})^{-1}\mathbf{G}' \quad (19.10)$$

For example, in the linear model,  $\mathbf{G}_i(\boldsymbol{\beta}) = -\mathbf{z}_i \mathbf{x}_i'$ ,  $\mathbf{G} = -\mathbb{E}(\mathbf{z}_i \mathbf{x}_i')$ , and  $\boldsymbol{\Omega} = \mathbb{E}(\mathbf{z}_i \mathbf{z}_i' e_i^2)$ .

**Theorem 19.2.1** *Under regularity conditions,*

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{V}_\beta)$$

$$\sqrt{n}\hat{\boldsymbol{\lambda}} \xrightarrow{d} \boldsymbol{\Omega}^{-1}\mathbf{N}(\mathbf{0}, \mathbf{V}_\lambda)$$

where  $\mathbf{V}$  and  $\mathbf{V}_\lambda$  are defined in (19.9) and (19.10), and  $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$  and  $\sqrt{n}\hat{\boldsymbol{\lambda}}$  are asymptotically independent.

The theorem shows that the asymptotic variance  $\mathbf{V}_\beta$  for  $\hat{\boldsymbol{\beta}}$  is the same as for efficient GMM. Thus the EL estimator is asymptotically efficient.

Chamberlain (1987) showed that  $\mathbf{V}_\beta$  is the semiparametric efficiency bound for  $\boldsymbol{\beta}$  in the overidentified moment condition model. This means that no consistent estimator for this class of models can have a lower asymptotic variance than  $\mathbf{V}_\beta$ . Since the EL estimator achieves this bound, it is an asymptotically efficient estimator for  $\boldsymbol{\beta}$ .

---

**Proof of Theorem 19.2.1.**  $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\lambda}})$  jointly solve

$$\mathbf{0} = \frac{\partial}{\partial \boldsymbol{\lambda}} R(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\lambda}}) = - \sum_{i=1}^n \frac{\mathbf{g}_i(\hat{\boldsymbol{\beta}})}{(1 + \hat{\boldsymbol{\lambda}}' \mathbf{g}_i(\hat{\boldsymbol{\beta}}))} \quad (19.11)$$

$$\mathbf{0} = \frac{\partial}{\partial \boldsymbol{\beta}} R(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\lambda}}) = - \sum_{i=1}^n \frac{\mathbf{G}_i(\hat{\boldsymbol{\beta}})' \boldsymbol{\lambda}}{1 + \hat{\boldsymbol{\lambda}}' \mathbf{g}_i(\hat{\boldsymbol{\beta}})}. \quad (19.12)$$

Let  $\mathbf{G}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{G}_i(\boldsymbol{\beta})$ ,  $\bar{\mathbf{g}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{g}_i(\boldsymbol{\beta})$  and  $\boldsymbol{\Omega}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{g}_i(\boldsymbol{\beta}) \mathbf{g}_i(\boldsymbol{\beta})'$ .

Expanding (19.12) around  $\boldsymbol{\beta}$  and  $\boldsymbol{\lambda} = \mathbf{0}$  yields

$$\mathbf{0} \simeq \mathbf{G}_n' \hat{\boldsymbol{\lambda}}. \quad (19.13)$$

Expanding (19.11) around  $\boldsymbol{\beta} = \boldsymbol{\beta}_0$  and  $\boldsymbol{\lambda} = \boldsymbol{\lambda}_0 = \mathbf{0}$  yields

$$\mathbf{0} \simeq -\bar{\mathbf{g}}_n - \mathbf{G}_n(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) + \boldsymbol{\Omega}_n \hat{\boldsymbol{\lambda}} \quad (19.14)$$

Premultiplying by  $\mathbf{G}'_n \boldsymbol{\Omega}_n^{-1}$  and using (19.13) yields

$$\begin{aligned} \mathbf{0} &\simeq -\mathbf{G}'_n \boldsymbol{\Omega}_n^{-1} \bar{\mathbf{g}}_n - \mathbf{G}'_n \boldsymbol{\Omega}_n^{-1} \mathbf{G}_n (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) + \mathbf{G}'_n \boldsymbol{\Omega}_n^{-1} \boldsymbol{\Omega}_n \hat{\boldsymbol{\lambda}} \\ &= -\mathbf{G}'_n \boldsymbol{\Omega}_n^{-1} \bar{\mathbf{g}}_n - \mathbf{G}'_n \boldsymbol{\Omega}_n^{-1} \mathbf{G}_n (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \end{aligned}$$

Solving for  $\hat{\boldsymbol{\beta}}$  and using the WLLN and CLT yields

$$\begin{aligned} \sqrt{n} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) &\simeq -(\mathbf{G}'_n \boldsymbol{\Omega}_n^{-1} \mathbf{G}_n)^{-1} \mathbf{G}'_n \boldsymbol{\Omega}_n^{-1} \sqrt{n} \bar{\mathbf{g}}_n \\ &\xrightarrow{d} (\mathbf{G}' \boldsymbol{\Omega}^{-1} \mathbf{G})^{-1} \mathbf{G}' \boldsymbol{\Omega}^{-1} \mathbf{N}(\mathbf{0}, \boldsymbol{\Omega}) \\ &= \mathbf{N}(\mathbf{0}, \mathbf{V}_\beta) \end{aligned} \tag{19.15}$$

Solving (19.14) for  $\hat{\boldsymbol{\lambda}}$  and using (19.15) yields

$$\begin{aligned} \sqrt{n} \hat{\boldsymbol{\lambda}} &\simeq \boldsymbol{\Omega}_n^{-1} \left( \mathbf{I} - \mathbf{G}_n (\mathbf{G}'_n \boldsymbol{\Omega}_n^{-1} \mathbf{G}_n)^{-1} \mathbf{G}'_n \boldsymbol{\Omega}_n^{-1} \right) \sqrt{n} \bar{\mathbf{g}}_n \\ &\xrightarrow{d} \boldsymbol{\Omega}^{-1} \left( \mathbf{I} - \mathbf{G} (\mathbf{G}' \boldsymbol{\Omega}^{-1} \mathbf{G})^{-1} \mathbf{G}' \boldsymbol{\Omega}^{-1} \right) \mathbf{N}(\mathbf{0}, \boldsymbol{\Omega}) \\ &= \boldsymbol{\Omega}^{-1} \mathbf{N}(\mathbf{0}, \mathbf{V}_\lambda) \end{aligned} \tag{19.16}$$

Furthermore, since

$$\mathbf{G}' \left( \mathbf{I} - \boldsymbol{\Omega}^{-1} \mathbf{G} (\mathbf{G}' \boldsymbol{\Omega}^{-1} \mathbf{G})^{-1} \mathbf{G}' \right) = \mathbf{0}$$

$\sqrt{n} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$  and  $\sqrt{n} \hat{\boldsymbol{\lambda}}$  are asymptotically uncorrelated and hence independent.

---

### 19.3 Overidentifying Restrictions

In a parametric likelihood context, tests are based on the difference in the log likelihood functions. The same statistic can be constructed for empirical likelihood. Twice the difference between the unrestricted empirical log-likelihood  $-n \log(n)$  and the maximized empirical log-likelihood for the model (19.7) is

$$LR_n = \sum_{i=1}^n 2 \log \left( 1 + \hat{\boldsymbol{\lambda}}' \mathbf{g}_i(\hat{\boldsymbol{\beta}}) \right). \tag{19.17}$$

**Theorem 19.3.1** *If  $\mathbb{E}(\mathbf{g}_i(\boldsymbol{\beta})) = \mathbf{0}$  then  $LR_n \xrightarrow{d} \chi^2_{\ell-k}$ .*

The EL overidentification test is similar to the GMM overidentification test. They are asymptotically first-order equivalent, and have the same interpretation. The overidentification test is a very useful by-product of EL estimation, and it is advisable to report the statistic  $LR_n$  whenever EL is the estimation method.

---

**Proof of Theorem 19.3.1.** First, by a Taylor expansion, (19.15), and (19.16),

$$\begin{aligned} \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{g}_i(\hat{\boldsymbol{\beta}}) &\simeq \sqrt{n} \left( \bar{\mathbf{g}}_n + \mathbf{G}_n (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \right) \\ &\simeq \left( \mathbf{I} - \mathbf{G}_n (\mathbf{G}'_n \boldsymbol{\Omega}_n^{-1} \mathbf{G}_n)^{-1} \mathbf{G}'_n \boldsymbol{\Omega}_n^{-1} \right) \sqrt{n} \bar{\mathbf{g}}_n \\ &\simeq \boldsymbol{\Omega}_n \sqrt{n} \hat{\boldsymbol{\lambda}}. \end{aligned}$$

Second, since  $\log(1 + u) \simeq u - u^2/2$  for  $u$  small,

$$\begin{aligned}
 LR_n &= \sum_{i=1}^n 2 \log \left( 1 + \hat{\lambda}' g_i(\hat{\beta}) \right) \\
 &\simeq 2 \hat{\lambda}' \sum_{i=1}^n g_i(\hat{\beta}) - \hat{\lambda}' \sum_{i=1}^n g_i(\hat{\beta}) g_i(\hat{\beta})' \hat{\lambda} \\
 &\simeq n \hat{\lambda}' \Omega_n \hat{\lambda} \\
 &\xrightarrow{d} N(\mathbf{0}, \mathbf{V}_\lambda)' \Omega^{-1} N(\mathbf{0}, \mathbf{V}_\lambda) \\
 &= \chi_{\ell-k}^2
 \end{aligned}$$

where the proof of the final equality is left as an exercise.

---

## 19.4 Testing

Let the maintained model be

$$\mathbb{E}(g_i(\beta)) = \mathbf{0} \quad (19.18)$$

where  $g$  is  $\ell \times 1$  and  $\beta$  is  $k \times 1$ . By “maintained” we mean that the overidentifying restrictions contained in (19.18) are assumed to hold and are not being challenged (at least for the test discussed in this section). The hypothesis of interest is

$$h(\beta) = \mathbf{0}.$$

where  $h : \mathbb{R}^k \rightarrow \mathbb{R}^a$ . The restricted EL estimator and likelihood are the values which solve

$$\begin{aligned}
 \tilde{\beta} &= \underset{h(\beta)=\mathbf{0}}{\operatorname{argmax}} R(\beta) \\
 R(\tilde{\beta}) &= \max_{h(\beta)=\mathbf{0}} R(\beta).
 \end{aligned}$$

Fundamentally, the restricted EL estimator  $\tilde{\beta}$  is simply an EL estimator with  $\ell - k + a$  overidentifying restrictions, so there is no fundamental change in the distribution theory for  $\tilde{\beta}$  relative to  $\hat{\beta}$ . To test the hypothesis  $h(\beta)$  while maintaining (19.18), the simple overidentifying restrictions test (19.17) is not appropriate. Instead we use the difference in log-likelihoods:

$$LR_n = 2 \left( R(\hat{\beta}) - R(\tilde{\beta}) \right).$$

This test statistic is a natural analog of the GMM distance statistic.

**Theorem 19.4.1** Under (19.18) and  $\mathbb{H}_0 : h(\beta) = \mathbf{0}$ ,  $LR_n \xrightarrow{d} \chi_a^2$ .

The proof of this result is more challenging and is omitted.

## 19.5 Numerical Computation

### Derivatives

The numerical calculations depend on derivatives of the dual likelihood function (19.4). Define

$$\begin{aligned} g_i^*(\beta, \lambda) &= \frac{g_i(\beta)}{(1 + \lambda' g_i(\beta))} \\ G_i^*(\beta, \lambda) &= \frac{G_i(\beta)' \lambda}{1 + \lambda' g_i(\beta)} \end{aligned}$$

The first derivatives of (19.4) are

$$\begin{aligned} R_\lambda &= \frac{\partial}{\partial \lambda} R(\beta, \lambda) = - \sum_{i=1}^n g_i^*(\beta, \lambda) \\ R_\beta &= \frac{\partial}{\partial \beta} R(\beta, \lambda) = - \sum_{i=1}^n G_i^*(\beta, \lambda). \end{aligned}$$

The second derivatives are

$$\begin{aligned} R_{\lambda\lambda} &= \frac{\partial^2}{\partial \lambda \partial \lambda'} R(\beta, \lambda) = \sum_{i=1}^n g_i^*(\beta, \lambda) g_i^*(\beta, \lambda)' \\ R_{\lambda\beta} &= \frac{\partial^2}{\partial \lambda \partial \beta'} R(\beta, \lambda) = \sum_{i=1}^n \left( g_i^*(\beta, \lambda) G_i^*(\beta, \lambda)' - \frac{G_i(\beta)}{1 + \lambda' g_i(\beta)} \right) \\ R_{\beta\beta} &= \frac{\partial^2}{\partial \beta \partial \beta'} R(\beta, \lambda) = \sum_{i=1}^n \left( G_i^*(\beta, \lambda) G_i^*(\beta, \lambda)' - \frac{\frac{\partial^2}{\partial \beta \partial \beta'} (g_i(\beta)' \lambda)}{1 + \lambda' g_i(\beta)} \right) \end{aligned}$$

### Inner Loop

The so-called “inner loop” solves (19.5) for given  $\beta$ . The modified Newton method takes a quadratic approximation to  $R_n(\beta, \lambda)$  yielding the iteration rule

$$\lambda_{j+1} = \lambda_j - \delta (R_{\lambda\lambda}(\beta, \lambda_j))^{-1} R_\lambda(\beta, \lambda_j). \quad (19.19)$$

where  $\delta > 0$  is a scalar steplength (to be discussed next). The starting value  $\lambda_1$  can be set to the zero vector. The iteration (19.19) is continued until the gradient  $R_\lambda(\beta, \lambda_j)$  is smaller than some prespecified tolerance.

Efficient convergence requires a good choice of steplength  $\delta$ . One method uses the following quadratic approximation. Set  $\delta_0 = 0$ ,  $\delta_1 = \frac{1}{2}$  and  $\delta_2 = 1$ . For  $p = 0, 1, 2$ , set

$$\begin{aligned} \lambda_p &= \lambda_j - \delta_p (R_{\lambda\lambda}(\beta, \lambda_j))^{-1} R_\lambda(\beta, \lambda_j) \\ R_p &= R(\beta, \lambda_p) \end{aligned}$$

A quadratic function can be fit exactly through these three points. The value of  $\delta$  which minimizes this quadratic is

$$\hat{\delta} = \frac{R_2 + 3R_0 - 4R_1}{4R_2 + 4R_0 - 8R_1}.$$

yielding the steplength to be plugged into (19.19).

A complication is that  $\lambda$  must be constrained so that  $0 \leq p_i \leq 1$  which holds if

$$n(1 + \lambda' g_i(\beta)) \geq 1 \quad (19.20)$$

for all  $i$ . If (19.20) fails, the stepsize  $\delta$  needs to be decreased.

**Outer Loop**

The outer loop is the minimization (19.6). This can be done by the modified Newton method described in the previous section. The gradient for (19.6) is

$$\mathbf{R}_\beta = \frac{\partial}{\partial \beta} R(\beta) = \frac{\partial}{\partial \beta} R(\beta, \lambda) = \mathbf{R}_\beta + \lambda'_\beta \mathbf{R}_\lambda = \mathbf{R}_\beta$$

since  $\mathbf{R}_\lambda(\beta, \lambda) = 0$  at  $\lambda = \lambda(\beta)$ , where

$$\lambda_\beta = \frac{\partial}{\partial \beta'} \lambda(\beta) = -\mathbf{R}_{\lambda\lambda}^{-1} \mathbf{R}_{\lambda\beta},$$

the second equality following from the implicit function theorem applied to  $\mathbf{R}_\lambda(\beta, \lambda(\beta)) = 0$ .

The Hessian for (19.6) is

$$\begin{aligned} \mathbf{R}_{\beta\beta} &= -\frac{\partial^2}{\partial \beta \partial \beta'} R(\beta) \\ &= -\frac{\partial}{\partial \beta'} [\mathbf{R}_\beta(\beta, \lambda(\beta)) + \lambda'_\beta \mathbf{R}_\lambda(\beta, \lambda(\beta))] \\ &= -(\mathbf{R}_{\beta\beta}(\beta, \lambda(\beta)) + \mathbf{R}'_{\lambda\beta} \lambda_\beta + \lambda'_\beta \mathbf{R}_{\lambda\beta} + \lambda'_\beta \mathbf{R}_{\lambda\lambda} \lambda_\beta) \\ &= \mathbf{R}'_{\lambda\beta} \mathbf{R}_{\lambda\lambda}^{-1} \mathbf{R}_{\lambda\beta} - \mathbf{R}_{\beta\beta}. \end{aligned}$$

It is not guaranteed that  $\mathbf{R}_{\beta\beta} > 0$ . If not, the eigenvalues of  $\mathbf{R}_{\beta\beta}$  should be adjusted so that all are positive. The Newton iteration rule is

$$\beta_{j+1} = \beta_j - \delta \mathbf{R}_{\beta\beta}^{-1} \mathbf{R}_\beta$$

where  $\delta$  is a scalar stepsize, and the rule is iterated until convergence.



# Chapter 20

## Regression Extensions

### 20.1 Nonlinear Least Squares

In some cases we might use a parametric regression function  $m(\mathbf{x}, \boldsymbol{\theta}) = \mathbb{E}(y_i \mid \mathbf{x}_i = \mathbf{x})$  which is a non-linear function of the parameters  $\boldsymbol{\theta}$ . We describe this setting as **nonlinear regression**.

#### Example 20.1.1 Exponential Link Regression

$$m(\mathbf{x}, \boldsymbol{\theta}) = \exp(\mathbf{x}'\boldsymbol{\theta})$$

The exponential link function is strictly positive, so this choice can be useful when it is desired to constrain the mean to be strictly positive.

#### Example 20.1.2 Logistic Link Regression

$$m(\mathbf{x}, \boldsymbol{\theta}) = \Lambda(\mathbf{x}'\boldsymbol{\theta})$$

where

$$\Lambda(u) = (1 + \exp(-u))^{-1} \quad (20.1)$$

is the Logistic distribution function. Since the logistic link function lies in  $[0, 1]$ , this choice can be useful when the conditional mean is bounded between 0 and 1.

#### Example 20.1.3 Exponentially Transformed Regressors

$$m(x, \boldsymbol{\theta}) = \theta_1 + \theta_2 \exp(\theta_3 x)$$

#### Example 20.1.4 Power Transformation

$$m(x, \boldsymbol{\theta}) = \theta_1 + \theta_2 x^{\theta_3}$$

with  $x > 0$ .

#### Example 20.1.5 Box-Cox Transformed Regressors

$$m(x, \boldsymbol{\theta}) = \theta_1 + \theta_2 x^{(\theta_3)}$$

where

$$x^{(\lambda)} = \begin{cases} \frac{x^\lambda - 1}{\lambda}, & \text{if } \lambda > 0 \\ \log(x), & \text{if } \lambda = 0 \end{cases} \quad (20.2)$$

and  $x > 0$ . The function (20.2) is called the Box-Cox Transformation and was introduced by Box and Cox (1964). The function nests linearity ( $\lambda = 1$ ) and logarithmic ( $\lambda = 0$ ) transformations continuously.

**Example 20.1.6** *Continuous Threshold Regression*

$$m(x, \theta) = \theta_1 + \theta_2 x + \theta_3 (x - \theta_4) 1(x > \theta_4)$$

**Example 20.1.7** *Threshold Regression*

$$m(x, \theta) = (\theta'_1 x_1) 1(x_2 < \theta_3) + (\theta'_2 x_1) 1(x_2 \geq \theta_3)$$

**Example 20.1.8** *Smooth Transition*

$$m(x, \theta) = \theta'_1 x_1 + (\theta'_2 x_1) \Lambda\left(\frac{x_2 - \theta_3}{\theta_4}\right)$$

where  $\Lambda(u)$  is the logit function (20.1).

What differentiates these examples from the linear regression model is that the conditional mean cannot be written as a linear function of the parameter vector  $\theta$ .

Nonlinear regression is sometimes adopted because the functional form  $m(x, \theta)$  is suggested by an economic model. In other cases, it is adopted as a flexible approximation to an unknown regression function.

The least squares estimator  $\hat{\theta}$  minimizes the normalized sum-of-squared-errors

$$\hat{S}(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - m(x_i, \theta))^2.$$

When the regression function is nonlinear, we call  $\hat{\theta}$  the **nonlinear least squares** (NLLS) estimator. The NLLS residuals are  $\hat{e}_i = y_i - m(x_i, \hat{\theta})$ .

One motivation for the choice of NLLS as the estimation method is that the parameter  $\theta$  is the solution to the population problem  $\min_{\theta} \mathbb{E} (y_i - m(x_i, \theta))^2$ .

Since the criterion  $\hat{S}(\theta)$  is not quadratic,  $\hat{\theta}$  must be found by numerical methods. See Appendix E. When  $m(x, \theta)$  is differentiable, then the FOC for minimization are

$$0 = \sum_{i=1}^n m_{\theta}(x_i, \hat{\theta}) \hat{e}_i \quad (20.3)$$

where

$$m_{\theta}(x, \theta) = \frac{\partial}{\partial \theta} m(x, \theta).$$

**Theorem 20.1.1** *Asymptotic Distribution of NLLS Estimator*

If the model is identified and  $m(x, \theta)$  is differentiable with respect to  $\theta$ ,

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, V_{\theta})$$

$$V_{\theta} = (\mathbb{E}(m_{\theta i} m'_{\theta i}))^{-1} (\mathbb{E}(m_{\theta i} m'_{\theta i} e_i^2)) (\mathbb{E}(m_{\theta i} m'_{\theta i}))^{-1}$$

where  $m_{\theta i} = m_{\theta}(x_i, \theta_0)$ .

Based on Theorem 20.1.1, an estimate of the asymptotic variance  $\mathbf{V}_\theta$  is

$$\widehat{\mathbf{V}}_\theta = \left( \frac{1}{n} \sum_{i=1}^n \widehat{\mathbf{m}}_{\theta_i} \widehat{\mathbf{m}}_{\theta_i}' \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \widehat{\mathbf{m}}_{\theta_i} \widehat{\mathbf{m}}_{\theta_i}' \widehat{e}_i^2 \right) \left( \frac{1}{n} \sum_{i=1}^n \widehat{\mathbf{m}}_{\theta_i} \widehat{\mathbf{m}}_{\theta_i}' \right)^{-1}$$

where  $\widehat{\mathbf{m}}_{\theta_i} = \mathbf{m}_\theta(\mathbf{x}_i, \widehat{\boldsymbol{\theta}})$  and  $\widehat{e}_i = y_i - m(\mathbf{x}_i, \widehat{\boldsymbol{\theta}})$ .

Identification is often tricky in nonlinear regression models. Suppose that

$$m(\mathbf{x}_i, \boldsymbol{\theta}) = \beta_1' \mathbf{z}_i + \beta_2' \mathbf{x}_i(\gamma)$$

where  $\mathbf{x}_i(\gamma)$  is a function of  $\mathbf{x}_i$  and the unknown parameter  $\gamma$ . Examples include  $x_i(\gamma) = x_i^\gamma$ ,  $x_i(\gamma) = \exp(\gamma x_i)$ , and  $x_i(\gamma) = x_i 1(g(x_i) > \gamma)$ . The model is linear when  $\beta_2 = \mathbf{0}$ , and this is often a useful hypothesis (sub-model) to consider. Thus we want to test

$$\mathbb{H}_0 : \beta_2 = \mathbf{0}.$$

However, under  $\mathbb{H}_0$ , the model is

$$y_i = \beta_1' \mathbf{z}_i + e_i$$

and both  $\beta_2$  and  $\gamma$  have dropped out. This means that under  $\mathbb{H}_0$ ,  $\gamma$  is not identified. This renders the distribution theory presented in the previous section invalid. Thus when the truth is that  $\beta_2 = \mathbf{0}$ , the parameter estimates are not asymptotically normally distributed. Furthermore, tests of  $\mathbb{H}_0$  do not have asymptotic normal or chi-square distributions.

The asymptotic theory of such tests have been worked out by Andrews and Ploberger (1994) and B. E. Hansen (1996). In particular, Hansen shows how to use simulation (similar to the bootstrap) to construct the asymptotic critical values (or p-values) in a given application.

**Proof of Theorem 20.1.1 (Sketch).** NLLS estimation falls in the class of optimization estimators. For this theory, it is useful to denote the true value of the parameter  $\boldsymbol{\theta}$  as  $\boldsymbol{\theta}_0$ .

The first step is to show that  $\widehat{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}_0$ . Proving that nonlinear estimators are consistent is more challenging than for linear estimators. We sketch the main argument. The idea is that  $\widehat{\boldsymbol{\theta}}$  minimizes the sample criterion function  $\widehat{S}(\boldsymbol{\theta})$ , which (for any  $\boldsymbol{\theta}$ ) converges in probability to the mean-squared error function  $\mathbb{E}((y_i - m(\mathbf{x}_i, \boldsymbol{\theta}))^2)$ . Thus it seems reasonable that the minimizer  $\widehat{\boldsymbol{\theta}}$  will converge in probability to  $\boldsymbol{\theta}_0$ , the minimizer of  $\mathbb{E}((y_i - m(\mathbf{x}_i, \boldsymbol{\theta}))^2)$ . It turns out that to show this rigorously, we need to show that  $\widehat{S}(\boldsymbol{\theta})$  converges *uniformly* to its expectation  $\mathbb{E}((y_i - m(\mathbf{x}_i, \boldsymbol{\theta}))^2)$ , which means that the maximum discrepancy must converge in probability to zero, to exclude the possibility that  $\widehat{S}(\boldsymbol{\theta})$  is excessively wiggly in  $\boldsymbol{\theta}$ . Proving uniform convergence is technically challenging, but it can be shown to hold broadly for relevant nonlinear regression models, especially if the regression function  $m(\mathbf{x}_i, \boldsymbol{\theta})$  is differentiable in  $\boldsymbol{\theta}$ . For a complete treatment of the theory of optimization estimators see Newey and McFadden (1994).

Since  $\widehat{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}_0$ ,  $\widehat{\boldsymbol{\theta}}$  is close to  $\boldsymbol{\theta}_0$  for  $n$  large, so the minimization of  $\widehat{S}(\boldsymbol{\theta})$  only needs to be examined for  $\boldsymbol{\theta}$  close to  $\boldsymbol{\theta}_0$ . Let

$$y_i^0 = e_i + \mathbf{m}_{\theta_i}' \boldsymbol{\theta}_0.$$

For  $\boldsymbol{\theta}$  close to the true value  $\boldsymbol{\theta}_0$ , by a first-order Taylor series approximation,

$$m(\mathbf{x}_i, \boldsymbol{\theta}) \simeq m(\mathbf{x}_i, \boldsymbol{\theta}_0) + \mathbf{m}_{\theta_i}' (\boldsymbol{\theta} - \boldsymbol{\theta}_0).$$

Thus

$$\begin{aligned} y_i - m(\mathbf{x}_i, \boldsymbol{\theta}) &\simeq (e_i + m(\mathbf{x}_i, \boldsymbol{\theta}_0)) - (m(\mathbf{x}_i, \boldsymbol{\theta}_0) + \mathbf{m}_{\theta_i}' (\boldsymbol{\theta} - \boldsymbol{\theta}_0)) \\ &= e_i - \mathbf{m}_{\theta_i}' (\boldsymbol{\theta} - \boldsymbol{\theta}_0) \\ &= y_i^0 - \mathbf{m}_{\theta_i}' \boldsymbol{\theta}. \end{aligned}$$

Hence the normalized sum of squared errors function is

$$\widehat{S}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (y_i - m(\mathbf{x}_i, \boldsymbol{\theta}))^2 \simeq \frac{1}{n} \sum_{i=1}^n (y_i^0 - \mathbf{m}'_{\boldsymbol{\theta}i} \boldsymbol{\theta})^2$$

and the right-hand-side is the criterion function for a linear regression of  $y_i^0$  on  $\mathbf{m}_{\boldsymbol{\theta}i}$ . Thus the NLLS estimator  $\widehat{\boldsymbol{\theta}}$  has the same asymptotic distribution as the (infeasible) OLS regression of  $y_i^0$  on  $\mathbf{m}_{\boldsymbol{\theta}i}$ , which is that stated in the theorem.

---

## 20.2 Generalized Least Squares

In the projection model, we know that the least-squares estimator is semi-parametrically efficient for the projection coefficient. However, in the linear regression model

$$\begin{aligned} y_i &= \mathbf{x}'_i \boldsymbol{\beta} + e_i \\ \mathbb{E}(e_i \mid \mathbf{x}_i) &= 0, \end{aligned}$$

the least-squares estimator is inefficient. The theory of Chamberlain (1987) can be used to show that in this model the semiparametric efficiency bound is obtained by the **Generalized Least Squares** (GLS) estimator (4.19) introduced in Section 4.7.1. The GLS estimator is sometimes called the Aitken estimator. The GLS estimator (20.2) is infeasible since the matrix  $\mathbf{D}$  is unknown. A feasible GLS (FGLS) estimator replaces the unknown  $\mathbf{D}$  with an estimate  $\widehat{\mathbf{D}} = \text{diag}\{\widehat{\sigma}_1^2, \dots, \widehat{\sigma}_n^2\}$ . We now discuss this estimation problem.

Suppose that we model the conditional variance using the parametric form

$$\begin{aligned} \sigma_i^2 &= \alpha_0 + \mathbf{z}'_{1i} \boldsymbol{\alpha}_1 \\ &= \boldsymbol{\alpha}' \mathbf{z}_i, \end{aligned}$$

where  $\mathbf{z}_{1i}$  is some  $q \times 1$  function of  $\mathbf{x}_i$ . Typically,  $\mathbf{z}_{1i}$  are squares (and perhaps levels) of some (or all) elements of  $\mathbf{x}_i$ . Often the functional form is kept simple for parsimony.

Let  $\eta_i = e_i^2$ . Then

$$\mathbb{E}(\eta_i \mid \mathbf{x}_i) = \alpha_0 + \mathbf{z}'_{1i} \boldsymbol{\alpha}_1$$

and we have the regression equation

$$\begin{aligned} \eta_i &= \alpha_0 + \mathbf{z}'_{1i} \boldsymbol{\alpha}_1 + \xi_i \\ \mathbb{E}(\xi_i \mid \mathbf{x}_i) &= 0. \end{aligned} \tag{20.4}$$

This regression error  $\xi_i$  is generally heteroskedastic and has the conditional variance

$$\begin{aligned} \text{var}(\xi_i \mid \mathbf{x}_i) &= \text{var}(e_i^2 \mid \mathbf{x}_i) \\ &= \mathbb{E}\left((e_i^2 - \mathbb{E}(e_i^2 \mid \mathbf{x}_i))^2 \mid \mathbf{x}_i\right) \\ &= \mathbb{E}(e_i^4 \mid \mathbf{x}_i) - (\mathbb{E}(e_i^2 \mid \mathbf{x}_i))^2. \end{aligned}$$

Suppose  $e_i$  (and thus  $\eta_i$ ) were observed. Then we could estimate  $\boldsymbol{\alpha}$  by OLS:

$$\widehat{\boldsymbol{\alpha}} = (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\boldsymbol{\eta} \xrightarrow{p} \boldsymbol{\alpha}$$

and

$$\sqrt{n}(\widehat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}) \xrightarrow{d} \text{N}(\mathbf{0}, \mathbf{V}_{\boldsymbol{\alpha}})$$

where

$$\mathbf{V}_\alpha = (\mathbb{E}(\mathbf{z}_i \mathbf{z}_i'))^{-1} \mathbb{E}(\mathbf{z}_i \mathbf{z}_i' \xi_i^2) (\mathbb{E}(\mathbf{z}_i \mathbf{z}_i'))^{-1}. \quad (20.5)$$

While  $e_i$  is not observed, we have the OLS residual  $\hat{e}_i = y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}} = e_i - \mathbf{x}_i' (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$ . Thus

$$\begin{aligned} \phi_i &\equiv \hat{\eta}_i - \eta_i \\ &= \hat{e}_i^2 - e_i^2 \\ &= -2e_i \mathbf{x}_i' (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) + (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{x}_i \mathbf{x}_i' (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}). \end{aligned}$$

And then

$$\begin{aligned} \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{z}_i \phi_i &= \frac{-2}{n} \sum_{i=1}^n \mathbf{z}_i e_i \mathbf{x}_i' \sqrt{n} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) + \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{x}_i \mathbf{x}_i' (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \sqrt{n} \\ &\xrightarrow{p} \mathbf{0} \end{aligned}$$

Let

$$\tilde{\boldsymbol{\alpha}} = (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \hat{\boldsymbol{\eta}} \quad (20.6)$$

be from OLS regression of  $\hat{\eta}_i$  on  $\mathbf{z}_i$ . Then

$$\begin{aligned} \sqrt{n} (\tilde{\boldsymbol{\alpha}} - \boldsymbol{\alpha}) &= \sqrt{n} (\hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}) + (n^{-1} \mathbf{Z}' \mathbf{Z})^{-1} n^{-1/2} \mathbf{Z}' \boldsymbol{\phi} \\ &\xrightarrow{d} \mathbf{N}(\mathbf{0}, \mathbf{V}_\alpha) \end{aligned} \quad (20.7)$$

Thus the fact that  $\eta_i$  is replaced with  $\hat{\eta}_i$  is asymptotically irrelevant. We call (20.6) the *skedastic* regression, as it is estimating the conditional variance of the regression of  $y_i$  on  $\mathbf{x}_i$ . We have shown that  $\boldsymbol{\alpha}$  is consistently estimated by a simple procedure, and hence we can estimate  $\sigma_i^2 = \mathbf{z}_i' \boldsymbol{\alpha}$  by

$$\tilde{\sigma}_i^2 = \tilde{\boldsymbol{\alpha}}' \mathbf{z}_i. \quad (20.8)$$

Suppose that  $\tilde{\sigma}_i^2 > 0$  for all  $i$ . Then set

$$\tilde{\mathbf{D}} = \text{diag}\{\tilde{\sigma}_1^2, \dots, \tilde{\sigma}_n^2\}$$

and

$$\tilde{\boldsymbol{\beta}} = (\mathbf{X}' \tilde{\mathbf{D}}^{-1} \mathbf{X})^{-1} \mathbf{X}' \tilde{\mathbf{D}}^{-1} \mathbf{y}.$$

This is the feasible GLS, or FGLS, estimator of  $\boldsymbol{\beta}$ . Since there is not a unique specification for the conditional variance the FGLS estimator is not unique, and will depend on the model (and estimation method) for the skedastic regression.

One typical problem with implementation of FGLS estimation is that in the linear specification (20.4), there is no guarantee that  $\tilde{\sigma}_i^2 > 0$  for all  $i$ . If  $\tilde{\sigma}_i^2 < 0$  for some  $i$ , then the FGLS estimator is not well defined. Furthermore, if  $\tilde{\sigma}_i^2 \approx 0$  for some  $i$  then the FGLS estimator will force the regression equation through the point  $(y_i, \mathbf{x}_i)$ , which is undesirable. This suggests that there is a need to bound the estimated variances away from zero. A trimming rule takes the form

$$\bar{\sigma}_i^2 = \max[\tilde{\sigma}_i^2, c\hat{\sigma}^2]$$

for some  $c > 0$ . For example, setting  $c = 1/4$  means that the conditional variance function is constrained to exceed one-fourth of the unconditional variance. As there is no clear method to select  $c$ , this introduces a degree of arbitrariness. In this context it is useful to re-estimate the model with several choices for the trimming parameter. If the estimates turn out to be sensitive to its choice, the estimation method should probably be reconsidered.

It is possible to show that if the skedastic regression is correctly specified, then FGLS is asymptotically equivalent to GLS. As the proof is tricky, we just state the result without proof.

**Theorem 20.2.1** *If the skedastic regression is correctly specified,*

$$\sqrt{n} \left( \tilde{\beta}_{GLS} - \tilde{\beta}_{FGLS} \right) \xrightarrow{p} \mathbf{0},$$

*and thus*

$$\sqrt{n} \left( \tilde{\beta}_{FGLS} - \beta \right) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_\beta),$$

*where*

$$\mathbf{V}_\beta = \left( \mathbb{E} \left( \sigma_i^{-2} \mathbf{x}_i \mathbf{x}_i' \right) \right)^{-1}.$$

Examining the asymptotic distribution of Theorem 20.2.1, the natural estimator of the asymptotic variance of  $\tilde{\beta}$  is

$$\tilde{\mathbf{V}}_\beta^0 = \left( \frac{1}{n} \sum_{i=1}^n \tilde{\sigma}_i^{-2} \mathbf{x}_i \mathbf{x}_i' \right)^{-1} = \left( \frac{1}{n} \mathbf{X}' \tilde{\mathbf{D}}^{-1} \mathbf{X} \right)^{-1}.$$

which is consistent for  $\mathbf{V}_\beta$  as  $n \rightarrow \infty$ . This estimator  $\tilde{\mathbf{V}}_\beta^0$  is appropriate when the skedastic regression (20.4) is correctly specified.

It may be the case that  $\alpha' \mathbf{z}_i$  is only an approximation to the true conditional variance  $\sigma_i^2 = \mathbb{E}(e_i^2 | \mathbf{x}_i)$ . In this case we interpret  $\alpha' \mathbf{z}_i$  as a linear projection of  $e_i^2$  on  $\mathbf{z}_i$ .  $\tilde{\beta}$  should perhaps be called a quasi-FGLS estimator of  $\beta$ . Its asymptotic variance is not that given in Theorem 20.2.1. Instead,

$$\mathbf{V}_\beta = \left( \mathbb{E} \left( (\alpha' \mathbf{z}_i)^{-1} \mathbf{x}_i \mathbf{x}_i' \right) \right)^{-1} \left( \mathbb{E} \left( (\alpha' \mathbf{z}_i)^{-2} \sigma_i^2 \mathbf{x}_i \mathbf{x}_i' \right) \right) \left( \mathbb{E} \left( (\alpha' \mathbf{z}_i)^{-1} \mathbf{x}_i \mathbf{x}_i' \right) \right)^{-1}.$$

$\mathbf{V}_\beta$  takes a sandwich form similar to the covariance matrix of the OLS estimator. Unless  $\sigma_i^2 = \alpha' \mathbf{z}_i$ ,  $\tilde{\mathbf{V}}_\beta^0$  is inconsistent for  $\mathbf{V}_\beta$ .

An appropriate solution is to use a White-type estimator in place of  $\tilde{\mathbf{V}}_\beta^0$ . This may be written as

$$\begin{aligned} \tilde{\mathbf{V}}_\beta &= \left( \frac{1}{n} \sum_{i=1}^n \tilde{\sigma}_i^{-2} \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \tilde{\sigma}_i^{-4} \hat{e}_i^2 \mathbf{x}_i \mathbf{x}_i' \right) \left( \frac{1}{n} \sum_{i=1}^n \tilde{\sigma}_i^{-2} \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \\ &= \left( \frac{1}{n} \mathbf{X}' \tilde{\mathbf{D}}^{-1} \mathbf{X} \right)^{-1} \left( \frac{1}{n} \mathbf{X}' \tilde{\mathbf{D}}^{-1} \hat{\mathbf{D}} \tilde{\mathbf{D}}^{-1} \mathbf{X} \right) \left( \frac{1}{n} \mathbf{X}' \tilde{\mathbf{D}}^{-1} \mathbf{X} \right)^{-1} \end{aligned}$$

where  $\hat{\mathbf{D}} = \text{diag}\{\hat{e}_1^2, \dots, \hat{e}_n^2\}$ . This estimator is robust to misspecification of the conditional variance, and was proposed by Cragg (1992).

In the linear regression model, FGLS is asymptotically superior to OLS. Why then do we not exclusively estimate regression models by FGLS? This is a good question. There are three reasons.

First, FGLS estimation depends on specification and estimation of the skedastic regression. Since the form of the skedastic regression is unknown, and it may be estimated with considerable error, the estimated conditional variances may contain more noise than information about the true conditional variances. In this case, FGLS can do worse than OLS in practice.

Second, individual estimated conditional variances may be negative, and this requires trimming to solve. This introduces an element of arbitrariness which is unsettling to empirical researchers.

Third, and probably most importantly, OLS is a robust estimator of the parameter vector. It is consistent not only in the regression model, but also under the assumptions of linear projection. The GLS and FGLS estimators, on the other hand, require the assumption of a correct conditional

mean. If the equation of interest is a linear projection and not a conditional mean, then the OLS and FGLS estimators will converge in probability to different limits as they will be estimating two different projections. The FGLS probability limit will depend on the particular function selected for the skedastic regression. The point is that the efficiency gains from FGLS are built on the stronger assumption of a correct conditional mean, and the cost is a loss of robustness to misspecification.

### 20.3 Testing for Heteroskedasticity

The hypothesis of homoskedasticity is that  $\mathbb{E}(e_i^2 | \mathbf{x}_i) = \sigma^2$ , or equivalently that

$$\mathbb{H}_0 : \boldsymbol{\alpha}_1 = 0$$

in the regression (20.4). We may therefore test this hypothesis by the estimation (20.6) and constructing a Wald statistic. In the classic literature it is typical to impose the stronger assumption that  $e_i$  is independent of  $\mathbf{x}_i$ , in which case  $\xi_i$  is independent of  $\mathbf{x}_i$  and the asymptotic variance (20.5) for  $\tilde{\boldsymbol{\alpha}}$  simplifies to

$$V_{\boldsymbol{\alpha}} = (\mathbb{E}(\mathbf{z}_i \mathbf{z}_i'))^{-1} \mathbb{E}(\xi_i^2). \quad (20.9)$$

Hence the standard test of  $\mathbb{H}_0$  is a classic  $F$  (or Wald) test for exclusion of all regressors from the skedastic regression (20.6). The asymptotic distribution (20.7) and the asymptotic variance (20.9) under independence show that this test has an asymptotic chi-square distribution.

**Theorem 20.3.1** *Under  $\mathbb{H}_0$  and  $e_i$  independent of  $\mathbf{x}_i$ , the Wald test of  $\mathbb{H}_0$  is asymptotically  $\chi_q^2$ .*

Most tests for heteroskedasticity take this basic form. The main differences between popular tests are which transformations of  $\mathbf{x}_i$  enter  $\mathbf{z}_i$ . Motivated by the form of the asymptotic variance of the OLS estimator  $\hat{\boldsymbol{\beta}}$ , White (1980) proposed that the test for heteroskedasticity be based on setting  $\mathbf{z}_i$  to equal all non-redundant elements of  $\mathbf{x}_i$ , its squares, and all cross-products. Breusch-Pagan (1979) proposed what might appear to be a distinct test, but the only difference is that they allowed for general choice of  $\mathbf{z}_i$ , and replaced  $\mathbb{E}(\xi_i^2)$  with  $2\sigma^4$  which holds when  $e_i$  is  $N(\mathbf{0}, \sigma^2)$ . If this simplification is replaced by the standard formula (under independence of the error), the two tests coincide.

It is important not to misuse tests for heteroskedasticity. It should not be used to determine whether to estimate a regression equation by OLS or FGLS, nor to determine whether classic or White standard errors should be reported. Hypothesis tests are not designed for these purposes. Rather, tests for heteroskedasticity should be used to answer the scientific question of whether or not the conditional variance is a function of the regressors. If this question is not of economic interest, then there is no value in conducting a test for heteroskedasticity.

### 20.4 Testing for Omitted Nonlinearity

If the goal is to estimate the conditional expectation  $\mathbb{E}(y_i | \mathbf{x}_i)$ , it is useful to have a general test of the adequacy of the specification.

One simple test for neglected nonlinearity is to add nonlinear functions of the regressors to the regression, and test their significance using a Wald test. Thus, if the model  $y_i = \mathbf{x}_i' \hat{\boldsymbol{\beta}} + \hat{e}_i$  has been fit by OLS, let  $\mathbf{z}_i = \mathbf{h}(\mathbf{x}_i)$  denote functions of  $\mathbf{x}_i$  which are not linear functions of  $\mathbf{x}_i$  (perhaps squares of non-binary regressors) and then fit  $y_i = \mathbf{x}_i' \hat{\boldsymbol{\beta}} + \mathbf{z}_i' \tilde{\boldsymbol{\gamma}} + \tilde{e}_i$  by OLS, and form a Wald statistic for  $\boldsymbol{\gamma} = \mathbf{0}$ .

Another popular approach is the RESET test proposed by Ramsey (1969). The null model is

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

which is estimated by OLS, yielding predicted values  $\hat{y}_i = \mathbf{x}_i' \hat{\boldsymbol{\beta}}$ . Now let

$$\mathbf{z}_i = \begin{pmatrix} \hat{y}_i^2 \\ \vdots \\ \hat{y}_i^m \end{pmatrix}$$

be a  $(m-1)$ -vector of powers of  $\hat{y}_i$ . Then run the auxiliary regression

$$y_i = \mathbf{x}_i' \tilde{\boldsymbol{\beta}} + \mathbf{z}_i' \tilde{\boldsymbol{\gamma}} + \tilde{e}_i \quad (20.10)$$

by OLS, and form the Wald statistic  $W$  for  $\boldsymbol{\gamma} = \mathbf{0}$ . It is easy (although somewhat tedious) to show that under the null hypothesis,  $W \xrightarrow{d} \chi_{m-1}^2$ . Thus the null is rejected at the  $\alpha\%$  level if  $W$  exceeds the upper  $1 - \alpha$  critical value of the  $\chi_{m-1}^2$  distribution.

To implement the test,  $m$  must be selected in advance. Typically, small values such as  $m = 2, 3$ , or  $4$  seem to work best.

The RESET test appears to work well as a test of functional form against a wide range of smooth alternatives. It is particularly powerful at detecting *single-index* models of the form

$$y_i = G(\mathbf{x}_i' \boldsymbol{\beta}) + e_i$$

where  $G(\cdot)$  is a smooth “link” function. To see why this is the case, note that (20.10) may be written as

$$y_i = \mathbf{x}_i' \tilde{\boldsymbol{\beta}} + \left(\mathbf{x}_i' \hat{\boldsymbol{\beta}}\right)^2 \tilde{\gamma}_1 + \left(\mathbf{x}_i' \hat{\boldsymbol{\beta}}\right)^3 \tilde{\gamma}_2 + \cdots + \left(\mathbf{x}_i' \hat{\boldsymbol{\beta}}\right)^m \tilde{\gamma}_{m-1} + \tilde{e}_i$$

which has essentially approximated  $G(\cdot)$  by a  $m$ 'th order polynomial

## 20.5 Least Absolute Deviations

We stated that a conventional goal in econometrics is estimation of impact of variation in  $\mathbf{x}_i$  on the central tendency of  $y_i$ . We have discussed projections and conditional means, but these are not the only measures of central tendency. An alternative good measure is the conditional median.

To recall the definition and properties of the median, let  $y$  be a continuous random variable. The median  $\theta = \text{med}(y)$  is the value such that  $\Pr(y \leq \theta) = \Pr(y \geq \theta) = 0.5$ . Two useful facts about the median are that

$$\theta = \underset{\theta}{\operatorname{argmin}} \mathbb{E} |y - \theta| \quad (20.11)$$

and

$$\mathbb{E} (\operatorname{sgn}(y - \theta)) = 0$$

where

$$\operatorname{sgn}(u) = \begin{cases} 1 & \text{if } u \geq 0 \\ -1 & \text{if } u < 0 \end{cases}$$

is the sign function.

These facts and definitions motivate three estimators of  $\theta$ . The first definition is the 50th empirical quantile. The second is the value which minimizes  $\frac{1}{n} \sum_{i=1}^n |y_i - \theta|$ , and the third definition is the solution to the moment equation  $\frac{1}{n} \sum_{i=1}^n \operatorname{sgn}(y_i - \theta) = 0$ . These distinctions are illusory, however, as these estimators are indeed identical.

Now let's consider the conditional median of  $y$  given a random vector  $\mathbf{x}$ . Let  $m(\mathbf{x}) = \text{med}(y | \mathbf{x})$  denote the conditional median of  $y$  given  $\mathbf{x}$ . The linear median regression model takes the form

$$\begin{aligned} y_i &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ \text{med}(e_i | \mathbf{x}_i) &= 0 \end{aligned}$$

In this model, the linear function  $\text{med}(y_i | \mathbf{x}_i = \mathbf{x}) = \mathbf{x}' \boldsymbol{\beta}$  is the conditional median function, and the substantive assumption is that the median function is linear in  $\mathbf{x}$ .

Conditional analogs of the facts about the median are



- $\Pr(y_i \leq \mathbf{x}'_i \boldsymbol{\beta} \mid \mathbf{x}_i = \mathbf{x}) = \Pr(y_i > \mathbf{x}'_i \boldsymbol{\beta} \mid \mathbf{x}_i = \mathbf{x}) = .5$
- $\mathbb{E}(\text{sgn}(e_i) \mid \mathbf{x}_i) = 0$
- $\mathbb{E}(\mathbf{x}_i \text{sgn}(e_i)) = 0$
- $\boldsymbol{\beta} = \min_{\boldsymbol{\beta}} \mathbb{E} |y_i - \mathbf{x}'_i \boldsymbol{\beta}|$

These facts motivate the following estimator. Let

$$LAD(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n |y_i - \mathbf{x}'_i \boldsymbol{\beta}|$$

be the average of absolute deviations. The **least absolute deviations** (LAD) estimator of  $\boldsymbol{\beta}$  minimizes this function

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\text{argmin}} LAD(\boldsymbol{\beta})$$

Equivalently, it is a solution to the moment condition

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \text{sgn}(y_i - \mathbf{x}'_i \hat{\boldsymbol{\beta}}) = 0. \quad (20.12)$$

The LAD estimator has an asymptotic normal distribution.

**Theorem 20.5.1 Asymptotic Distribution of LAD Estimator**

*When the conditional median is linear in  $\mathbf{x}$*

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} N(\mathbf{0}, \mathbf{V})$$

where

$$\mathbf{V} = \frac{1}{4} (\mathbb{E}(\mathbf{x}_i \mathbf{x}'_i f(0 \mid \mathbf{x}_i)))^{-1} (\mathbb{E}(\mathbf{x}_i \mathbf{x}'_i)) (\mathbb{E}(\mathbf{x}_i \mathbf{x}'_i f(0 \mid \mathbf{x}_i)))^{-1}$$

and  $f(e \mid \mathbf{x})$  is the conditional density of  $e_i$  given  $\mathbf{x}_i = \mathbf{x}$ .

The variance of the asymptotic distribution inversely depends on  $f(0 \mid \mathbf{x})$ , the conditional density of the error at its median. When  $f(0 \mid \mathbf{x})$  is large, then there are many innovations near to the median, and this improves estimation of the median. In the special case where the error is independent of  $\mathbf{x}_i$ , then  $f(0 \mid \mathbf{x}) = f(0)$  and the asymptotic variance simplifies

$$\mathbf{V} = \frac{(\mathbb{E}(\mathbf{x}_i \mathbf{x}'_i))^{-1}}{4f(0)^2} \quad (20.13)$$

This simplification is similar to the simplification of the asymptotic covariance of the OLS estimator under homoskedasticity.

Computation of standard error for LAD estimates typically is based on equation (20.13). The main difficulty is the estimation of  $f(0)$ , the height of the error density at its median. This can be done with kernel estimation techniques. See Chapter 22. While a complete proof of Theorem 20.5.1 is advanced, we provide a sketch here for completeness.

---

**Proof of Theorem 20.5.1:** Similar to NLLS, LAD is an optimization estimator. Let  $\boldsymbol{\beta}_0$  denote the true value of  $\boldsymbol{\beta}_0$ .

The first step is to show that  $\hat{\beta} \xrightarrow{p} \beta_0$ . The general nature of the proof is similar to that for the NLLS estimator, and is sketched here. For any fixed  $\beta$ , by the WLLN,  $LAD(\beta) \xrightarrow{p} \mathbb{E} |y_i - \mathbf{x}'_i \beta|$ . Furthermore, it can be shown that this convergence is uniform in  $\beta$ . (Proving uniform convergence is more challenging than for the NLLS criterion since the LAD criterion is not differentiable in  $\beta$ .) It follows that  $\hat{\beta}$ , the minimizer of  $LAD(\beta)$ , converges in probability to  $\beta_0$ , the minimizer of  $\mathbb{E} |y_i - \mathbf{x}'_i \beta|$ .

Since  $\text{sgn}(a) = 1 - 2 \cdot 1(a \leq 0)$ , (20.12) is equivalent to  $\bar{\mathbf{g}}_n(\hat{\beta}) = 0$ , where  $\bar{\mathbf{g}}_n(\beta) = n^{-1} \sum_{i=1}^n \mathbf{g}_i(\beta)$  and  $\mathbf{g}_i(\beta) = \mathbf{x}_i (1 - 2 \cdot 1(y_i \leq \mathbf{x}'_i \beta))$ . Let  $\mathbf{g}(\beta) = \mathbb{E}(\mathbf{g}_i(\beta))$ . We need three preliminary results. First, since  $\mathbb{E}(\mathbf{g}_i(\beta_0)) = 0$  and  $\mathbb{E}(\mathbf{g}_i(\beta_0) \mathbf{g}_i(\beta_0)') = \mathbb{E}(\mathbf{x}_i \mathbf{x}_i')$ , we can apply the central limit theorem (Theorem 6.8.1) and find that

$$\sqrt{n} \bar{\mathbf{g}}_n(\beta_0) = n^{-1/2} \sum_{i=1}^n \mathbf{g}_i(\beta_0) \xrightarrow{d} N(\mathbf{0}, \mathbb{E}(\mathbf{x}_i \mathbf{x}_i')).$$

Second using the law of iterated expectations and the chain rule of differentiation,

$$\begin{aligned} \frac{\partial}{\partial \beta'} \mathbf{g}(\beta) &= \frac{\partial}{\partial \beta'} \mathbb{E} \mathbf{x}_i (1 - 2 \cdot 1(y_i \leq \mathbf{x}'_i \beta)) \\ &= -2 \frac{\partial}{\partial \beta'} \mathbb{E} (\mathbf{x}_i \mathbb{E}(1(e_i \leq \mathbf{x}'_i \beta - \mathbf{x}'_i \beta_0) | \mathbf{x}_i)) \\ &= -2 \frac{\partial}{\partial \beta'} \mathbb{E} \left( \mathbf{x}_i \int_{-\infty}^{\mathbf{x}'_i \beta - \mathbf{x}'_i \beta_0} f(e | \mathbf{x}_i) de \right) \\ &= -2 \mathbb{E} (\mathbf{x}_i \mathbf{x}'_i f(\mathbf{x}'_i \beta - \mathbf{x}'_i \beta_0 | \mathbf{x}_i)) \end{aligned}$$

so

$$\frac{\partial}{\partial \beta'} \mathbf{g}(\beta) = -2 \mathbb{E} (\mathbf{x}_i \mathbf{x}'_i f(0 | \mathbf{x}_i)).$$

Third, by a Taylor series expansion and the fact  $\mathbf{g}(\beta) = 0$

$$\mathbf{g}(\hat{\beta}) \simeq \frac{\partial}{\partial \beta'} \mathbf{g}(\beta) (\hat{\beta} - \beta).$$

Together

$$\begin{aligned} \sqrt{n} (\hat{\beta} - \beta_0) &\simeq \left( \frac{\partial}{\partial \beta'} \mathbf{g}(\beta_0) \right)^{-1} \sqrt{n} \mathbf{g}(\hat{\beta}) \\ &= (-2 \mathbb{E} (\mathbf{x}_i \mathbf{x}'_i f(0 | \mathbf{x}_i)))^{-1} \sqrt{n} (\mathbf{g}(\hat{\beta}) - \bar{\mathbf{g}}_n(\hat{\beta})) \\ &\simeq \frac{1}{2} (\mathbb{E} (\mathbf{x}_i \mathbf{x}'_i f(0 | \mathbf{x}_i)))^{-1} \sqrt{n} (\bar{\mathbf{g}}_n(\beta_0) - \mathbf{g}(\beta_0)) \\ &\xrightarrow{d} \frac{1}{2} (\mathbb{E} [\mathbf{x}_i \mathbf{x}'_i f(0 | \mathbf{x}_i)])^{-1} N(\mathbf{0}, \mathbb{E}(\mathbf{x}_i \mathbf{x}_i')) \\ &= N(\mathbf{0}, \mathbf{V}). \end{aligned}$$

The third line follows from an asymptotic empirical process argument and the fact that  $\hat{\beta} \xrightarrow{p} \beta_0$ .

## 20.6 Quantile Regression

Quantile regression has become quite popular in recent econometric practice. For  $\tau \in [0, 1]$  the  $\tau^{th}$  quantile  $Q_\tau$  of a random variable with distribution function  $F(u)$  is defined as

$$Q_\tau = \inf \{u : F(u) \geq \tau\}$$

When  $F(u)$  is continuous and strictly monotonic, then  $F(Q_\tau) = \tau$ , so you can think of the quantile as the inverse of the distribution function. The quantile  $Q_\tau$  is the value such that  $\tau$  (percent) of the mass of the distribution is less than  $Q_\tau$ . The median is the special case  $\tau = .5$ .

The following alternative representation is useful. If the random variable  $U$  has  $\tau^{th}$  quantile  $Q_\tau$ , then

$$Q_\tau = \underset{\theta}{\operatorname{argmin}} \mathbb{E}(\rho_\tau(U - \theta)). \quad (20.14)$$

where  $\rho_\tau(q)$  is the piecewise linear function

$$\begin{aligned} \rho_\tau(q) &= \begin{cases} -q(1 - \tau) & q < 0 \\ q\tau & q \geq 0 \end{cases} \\ &= q(\tau - 1(q < 0)). \end{aligned} \quad (20.15)$$

This generalizes representation (20.11) for the median to all quantiles.

For the random variables  $(y_i, \mathbf{x}_i)$  with conditional distribution function  $F(y | \mathbf{x})$  the conditional quantile function  $q_\tau(\mathbf{x})$  is

$$Q_\tau(\mathbf{x}) = \inf \{y : F(y | \mathbf{x}) \geq \tau\}.$$

Again, when  $F(y | \mathbf{x})$  is continuous and strictly monotonic in  $y$ , then  $F(Q_\tau(\mathbf{x}) | \mathbf{x}) = \tau$ . For fixed  $\tau$ , the quantile regression function  $q_\tau(\mathbf{x})$  describes how the  $\tau^{th}$  quantile of the conditional distribution varies with the regressors.

As functions of  $\mathbf{x}$ , the quantile regression functions can take any shape. However for computational convenience it is typical to assume that they are (approximately) linear in  $\mathbf{x}$  (after suitable transformations). This linear specification assumes that  $Q_\tau(\mathbf{x}) = \beta'_\tau \mathbf{x}$  where the coefficients  $\beta_\tau$  vary across the quantiles  $\tau$ . We then have the linear quantile regression model

$$y_i = \mathbf{x}'_i \beta_\tau + e_i$$

where  $e_i$  is the error defined to be the difference between  $y_i$  and its  $\tau^{th}$  conditional quantile  $\mathbf{x}'_i \beta_\tau$ . By construction, the  $\tau^{th}$  conditional quantile of  $e_i$  is zero, otherwise its properties are unspecified without further restrictions.

Given the representation (20.14), the quantile regression estimator  $\hat{\beta}_\tau$  for  $\beta_\tau$  solves the minimization problem

$$\hat{\beta}_\tau = \underset{\beta}{\operatorname{argmin}} S^\tau(\beta)$$

where

$$S^\tau(\beta) = \frac{1}{n} \sum_{i=1}^n \rho_\tau(y_i - \mathbf{x}'_i \beta)$$

and  $\rho_\tau(q)$  is defined in (20.15).

Since the quantile regression criterion function  $S^\tau(\beta)$  does not have an algebraic solution, numerical methods are necessary for its minimization. Furthermore, since it has discontinuous derivatives, conventional Newton-type optimization methods are inappropriate. Fortunately, fast linear programming methods have been developed for this problem, and are widely available.

An asymptotic distribution theory for the quantile regression estimator can be derived using similar arguments as those for the LAD estimator in Theorem 20.5.1.

**Theorem 20.6.1** *Asymptotic Distribution of the Quantile Regression Estimator*

When the  $\tau^{\text{th}}$  conditional quantile is linear in  $\mathbf{x}$

$$\sqrt{n} \left( \hat{\beta}_\tau - \beta_\tau \right) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_\tau),$$

where

$$\mathbf{V}_\tau = \tau(1-\tau) \left( \mathbb{E}(\mathbf{x}_i \mathbf{x}_i' f(0 | \mathbf{x}_i)) \right)^{-1} \left( \mathbb{E}(\mathbf{x}_i \mathbf{x}_i') \right) \left( \mathbb{E}(\mathbf{x}_i \mathbf{x}_i' f(0 | \mathbf{x}_i)) \right)^{-1}$$

and  $f(e | \mathbf{x})$  is the conditional density of  $e_i$  given  $\mathbf{x}_i = \mathbf{x}$ .

In general, the asymptotic variance depends on the conditional density of the quantile regression error. When the error  $e_i$  is independent of  $\mathbf{x}_i$ , then  $f(0 | \mathbf{x}_i) = f(0)$ , the unconditional density of  $e_i$  at 0, and we have the simplification

$$\mathbf{V}_\tau = \frac{\tau(1-\tau)}{f(0)^2} \left( \mathbb{E}(\mathbf{x}_i \mathbf{x}_i') \right)^{-1}.$$

A recent monograph on the details of quantile regression is Koenker (2005).

## Exercises

**Exercise 20.1** Suppose that  $y_i = g(\mathbf{x}_i, \boldsymbol{\theta}) + e_i$  with  $\mathbb{E}(e_i | \mathbf{x}_i) = 0$ ,  $\hat{\boldsymbol{\theta}}$  is the NLLS estimator, and  $\hat{\mathbf{V}}$  is the estimate of  $\text{var}(\hat{\boldsymbol{\theta}})$ . You are interested in the conditional mean function  $\mathbb{E}(y_i | \mathbf{x}_i = \mathbf{x}) = g(\mathbf{x})$  at some  $\mathbf{x}$ . Find an asymptotic 95% confidence interval for  $g(\mathbf{x})$ .

**Exercise 20.2** In Exercise 9.26, you estimated a cost function on a cross-section of electric companies. The equation you estimated was

$$\log TC_i = \beta_1 + \beta_2 \log Q_i + \beta_3 \log PL_i + \beta_4 \log PK_i + \beta_5 \log PF_i + e_i. \quad (20.16)$$

- (a) Following Nerlove, add the variable  $(\log Q_i)^2$  to the regression. Do so. Assess the merits of this new specification using a hypothesis test. Do you agree with this modification?
- (b) Now try a non-linear specification. Consider model (20.16) plus the extra term  $\beta_6 z_i$ , where

$$z_i = \log Q_i (1 + \exp(-(\log Q_i - \beta_7)))^{-1}.$$

In addition, impose the restriction  $\beta_3 + \beta_4 + \beta_5 = 1$ . This model is called a smooth threshold model. For values of  $\log Q_i$  much below  $\beta_7$ , the variable  $\log Q_i$  has a regression slope of  $\beta_2$ . For values much above  $\beta_7$ , the regression slope is  $\beta_2 + \beta_6$ , and the model imposes a smooth transition between these regimes. The model is non-linear because of the parameter  $\beta_7$ .

The model works best when  $\beta_7$  is selected so that several values (in this example, at least 10 to 15) of  $\log Q_i$  are both below and above  $\beta_7$ . Examine the data and pick an appropriate range for  $\beta_7$ .

- (c) Estimate the model by non-linear least squares. I recommend the concentration method: Pick 10 (or more if you like) values of  $\beta_7$  in this range. For each value of  $\beta_7$ , calculate  $z_i$  and estimate the model by OLS. Record the sum of squared errors, and find the value of  $\beta_7$  for which the sum of squared errors is minimized.
- (d) Calculate standard errors for all the parameters  $(\beta_1, \dots, \beta_7)$ .

**Exercise 20.3** Using the CPS data set, return to the linear regression model reported in Table 4.1

- (a) Re-estimate the model by least-squares. You do not need to report the estimates, but confirm that you obtain the same results.
- (b) Test whether the error variance is different for men and women. Interpret.
- (c) Test whether the error variance is different across the race groups (White, Black, American Indian, Asian, Mixed Race). Interpret.
- (d) Construct a model for the conditional variance. Estimate such a model, test for general heteroskedasticity and report the results.
- (e) Using this model for the conditional variance, re-estimate the model from part (c) using FGLS. Report the results.
- (f) Do the OLS and FGLS estimates differ greatly? Note any interesting differences.
- (g) Compare the estimated standard errors. Note any interesting differences.

**Exercise 20.4** For any predictor  $g(\mathbf{x}_i)$  for  $y_i$ , the mean absolute error (MAE) is

$$\mathbb{E} |y_i - g(\mathbf{x}_i)|.$$

Show that the function  $g(\mathbf{x})$  which minimizes the MAE is the conditional median  $m(\mathbf{x}) = \text{med}(y_i | \mathbf{x}_i)$ .

**Exercise 20.5** Define

$$g(u) = \tau - 1(u < 0)$$

where  $1(\cdot)$  is the indicator function (takes the value 1 if the argument is true, else equals zero). Let  $\theta$  satisfy  $\mathbb{E}(g(y_i - \theta)) = 0$ . Is  $\theta$  a quantile of the distribution of  $y_i$ ?

**Exercise 20.6** Verify equation (20.14)

**Exercise 20.7** You are interested in estimating the equation  $y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$ . You believe the regressors are exogenous, but you are uncertain about the properties of the error. You estimate the equation both by least absolute deviations (LAD) and OLS. A colleague suggests that you should prefer the OLS estimate, because it produces a higher  $R^2$  than the LAD estimate. Is your colleague correct?

# Chapter 21

## Limited Dependent Variables

$y$  is a **limited dependent variable** if it takes values in a strict subset of  $\mathbb{R}$ . The most common cases are

- Binary:  $y \in \{0, 1\}$
- Multinomial:  $y \in \{0, 1, 2, \dots, k\}$
- Integer:  $y \in \{0, 1, 2, \dots\}$
- Censored:  $y \in \mathbb{R}^+$

The traditional approach to the estimation of limited dependent variable (LDV) models is parametric maximum likelihood. A parametric model is constructed, allowing the construction of the likelihood function. A more modern approach is semi-parametric, eliminating the dependence on a parametric distributional assumption. We will discuss only the first (parametric) approach, due to time constraints. They still constitute the majority of LDV applications. If, however, you were to write a thesis involving LDV estimation, you would be advised to consider employing a semi-parametric estimation approach.

For the parametric approach, estimation is by MLE. A major practical issue is construction of the likelihood function.

### 21.1 Binary Choice

The dependent variable  $y_i \in \{0, 1\}$ . This represents a Yes/No outcome. Given some regressors  $\mathbf{x}_i$ , the goal is to describe  $\Pr(y_i = 1 \mid \mathbf{x}_i)$ , as this is the full conditional distribution.

The linear probability model specifies that

$$\Pr(y_i = 1 \mid \mathbf{x}_i) = \mathbf{x}_i' \boldsymbol{\beta}.$$

As  $\Pr(y_i = 1 \mid \mathbf{x}_i) = \mathbb{E}(y_i \mid \mathbf{x}_i)$ , this yields the regression:  $y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$  which can be estimated by OLS. However, the linear probability model does not impose the restriction that  $0 \leq \Pr(y_i \mid \mathbf{x}_i) \leq 1$ . Even so estimation of a linear probability model is a useful starting point for subsequent analysis.

The standard alternative is to use a function of the form

$$\Pr(y_i = 1 \mid \mathbf{x}_i) = F(\mathbf{x}_i' \boldsymbol{\beta})$$

where  $F(\cdot)$  is a known CDF, typically assumed to be symmetric about zero, so that  $F(u) = 1 - F(-u)$ . The two standard choices for  $F$  are

- Logistic:  $F(u) = (1 + e^{-u})^{-1}$ .

- Normal:  $F(u) = \Phi(u)$ .

If  $F$  is logistic, we call this the **logit** model, and if  $F$  is normal, we call this the **probit** model. This model is identical to the latent variable model

$$\begin{aligned} y_i^* &= \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ e_i &\sim F(\cdot) \\ y_i &= \begin{cases} 1 & \text{if } y_i^* > 0 \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

For then

$$\begin{aligned} \Pr(y_i = 1 \mid \mathbf{x}_i) &= \Pr(y_i^* > 0 \mid \mathbf{x}_i) \\ &= \Pr(\mathbf{x}_i' \boldsymbol{\beta} + e_i > 0 \mid \mathbf{x}_i) \\ &= \Pr(e_i > -\mathbf{x}_i' \boldsymbol{\beta} \mid \mathbf{x}_i) \\ &= 1 - F(-\mathbf{x}_i' \boldsymbol{\beta}) \\ &= F(\mathbf{x}_i' \boldsymbol{\beta}). \end{aligned}$$

Estimation is by maximum likelihood. To construct the likelihood, we need the conditional distribution of an individual observation. Recall that if  $y$  is Bernoulli, such that  $\Pr(y = 1) = p$  and  $\Pr(y = 0) = 1 - p$ , then we can write the density of  $y$  as

$$f(y) = p^y (1 - p)^{1-y}, \quad y = 0, 1.$$

In the Binary choice model,  $y_i$  is conditionally Bernoulli with  $\Pr(y_i = 1 \mid \mathbf{x}_i) = p_i = F(\mathbf{x}_i' \boldsymbol{\beta})$ . Thus the conditional density is

$$\begin{aligned} f(y_i \mid \mathbf{x}_i) &= p_i^{y_i} (1 - p_i)^{1-y_i} \\ &= F(\mathbf{x}_i' \boldsymbol{\beta})^{y_i} (1 - F(\mathbf{x}_i' \boldsymbol{\beta}))^{1-y_i}. \end{aligned}$$

Hence the log-likelihood function is

$$\begin{aligned} \log L(\boldsymbol{\beta}) &= \sum_{i=1}^n \log f(y_i \mid \mathbf{x}_i) \\ &= \sum_{i=1}^n \log (F(\mathbf{x}_i' \boldsymbol{\beta})^{y_i} (1 - F(\mathbf{x}_i' \boldsymbol{\beta}))^{1-y_i}) \\ &= \sum_{i=1}^n [y_i \log F(\mathbf{x}_i' \boldsymbol{\beta}) + (1 - y_i) \log(1 - F(\mathbf{x}_i' \boldsymbol{\beta}))] \\ &= \sum_{y_i=1} \log F(\mathbf{x}_i' \boldsymbol{\beta}) + \sum_{y_i=0} \log(1 - F(\mathbf{x}_i' \boldsymbol{\beta})). \end{aligned}$$

The MLE  $\hat{\boldsymbol{\beta}}$  is the value of  $\boldsymbol{\beta}$  which maximizes  $\log L(\boldsymbol{\beta})$ . Standard errors and test statistics are computed by asymptotic approximations. Details of such calculations are left to more advanced courses.

## 21.2 Count Data

If  $y \in \{0, 1, 2, \dots\}$ , a typical approach is to employ **Poisson regression**. This model specifies that

$$\begin{aligned} \Pr(y_i = k \mid \mathbf{x}_i) &= \frac{\exp(-\lambda_i) \lambda_i^k}{k!}, \quad k = 0, 1, 2, \dots \\ \lambda_i &= \exp(\mathbf{x}_i' \boldsymbol{\beta}). \end{aligned}$$



The conditional density is the Poisson with parameter  $\lambda_i$ . The functional form for  $\lambda_i$  has been picked to ensure that  $\lambda_i > 0$ .

The log-likelihood function is

$$\log L(\boldsymbol{\beta}) = \sum_{i=1}^n \log f(y_i | \mathbf{x}_i) = \sum_{i=1}^n \left( -\exp(\mathbf{x}_i' \boldsymbol{\beta}) + y_i \mathbf{x}_i' \boldsymbol{\beta} - \log(y_i!) \right).$$

The MLE is the value  $\hat{\boldsymbol{\beta}}$  which maximizes  $\log L(\boldsymbol{\beta})$ .

Since

$$\mathbb{E}(y_i | \mathbf{x}_i) = \lambda_i = \exp(\mathbf{x}_i' \boldsymbol{\beta})$$

is the conditional mean, this motivates the label Poisson “regression.”

Also observe that the model implies that

$$\text{var}(y_i | \mathbf{x}_i) = \lambda_i = \exp(\mathbf{x}_i' \boldsymbol{\beta}),$$

so the model imposes the restriction that the conditional mean and variance of  $y_i$  are the same. This may be considered restrictive. A generalization is the negative binomial.

## 21.3 Censored Data

The idea of **censoring** is that some data above or below a threshold are mis-reported at the threshold. Thus the model is that there is some latent process  $y_i^*$  with unbounded support, but we observe only

$$y_i = \begin{cases} y_i^* & \text{if } y_i^* \geq 0 \\ 0 & \text{if } y_i^* < 0 \end{cases} \quad (21.1)$$

(This is written for the case of the threshold being zero, any known value can substitute.) The observed data  $y_i$  therefore come from a mixed continuous/discrete distribution.

Censored models are typically applied when the data set has a meaningful proportion (say 5% or higher) of data at the boundary of the sample support. The censoring process may be explicit in data collection, or it may be a by-product of economic constraints.

An example of a data collection censoring is top-coding of income. In surveys, incomes above a threshold are typically reported at the threshold.

The first censored regression model was developed by Tobin (1958) to explain consumption of durable goods. Tobin observed that for many households, the consumption level (purchases) in a particular period was zero. He proposed the latent variable model

$$y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + e_i \\ e_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

with the observed variable  $y_i$  generated by the censoring equation (21.1). This model (now called the Tobit) specifies that the latent (or ideal) value of consumption may be negative (the household would prefer to sell than buy). All that is reported is that the household purchased zero units of the good.

The naive approach to estimate  $\boldsymbol{\beta}$  is to regress  $y_i$  on  $\mathbf{x}_i$ . This does not work because regression estimates  $\mathbb{E}(y_i | \mathbf{x}_i)$ , not  $\mathbb{E}(y_i^* | \mathbf{x}_i) = \mathbf{x}_i' \boldsymbol{\beta}$ , and the latter is of interest. Thus OLS will be biased for the parameter of interest  $\boldsymbol{\beta}$ .

[Note: it is still possible to estimate  $\mathbb{E}(y_i | \mathbf{x}_i)$  by LS techniques. The Tobit framework postulates that this is not inherently interesting, that the parameter of  $\boldsymbol{\beta}$  is defined by an alternative statistical structure.]

Consistent estimation will be achieved by the MLE. To construct the likelihood, observe that the probability of being censored is

$$\begin{aligned}\Pr(y_i = 0 \mid \mathbf{x}_i) &= \Pr(y_i^* < 0 \mid \mathbf{x}_i) \\ &= \Pr(\mathbf{x}_i' \boldsymbol{\beta} + e_i < 0 \mid \mathbf{x}_i) \\ &= \Pr\left(\frac{e_i}{\sigma} < -\frac{\mathbf{x}_i' \boldsymbol{\beta}}{\sigma} \mid \mathbf{x}_i\right) \\ &= \Phi\left(-\frac{\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right).\end{aligned}$$

The conditional density function above zero is normal:

$$\sigma^{-1} \phi\left(\frac{y - \mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right), \quad y > 0.$$

Therefore, the density function for  $y \geq 0$  can be written as

$$f(y \mid \mathbf{x}_i) = \Phi\left(-\frac{\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right)^{1(y=0)} \left[\sigma^{-1} \phi\left(\frac{y - \mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right)\right]^{1(y>0)},$$

where  $1(\cdot)$  is the indicator function.

Hence the log-likelihood is a mixture of the probit and the normal:

$$\begin{aligned}\log L(\boldsymbol{\beta}) &= \sum_{i=1}^n \log f(y_i \mid \mathbf{x}_i) \\ &= \sum_{y_i=0} \log \Phi\left(-\frac{\mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right) + \sum_{y_i>0} \log \left[\sigma^{-1} \phi\left(\frac{y_i - \mathbf{x}_i' \boldsymbol{\beta}}{\sigma}\right)\right].\end{aligned}$$

The MLE is the value  $\hat{\boldsymbol{\beta}}$  which maximizes  $\log L(\boldsymbol{\beta})$ .

## 21.4 Sample Selection

The problem of sample selection arises when the sample is a non-random selection of potential observations. This occurs when the observed data is systematically different from the population of interest. For example, if you ask for volunteers for an experiment, and they wish to extrapolate the effects of the experiment on a general population, you should worry that the people who volunteer may be systematically different from the general population. This has great relevance for the evaluation of anti-poverty and job-training programs, where the goal is to assess the effect of “training” on the general population, not just on the volunteers.

A simple sample selection model can be written as the latent model

$$\begin{aligned}y_i &= \mathbf{x}_i' \boldsymbol{\beta} + e_{1i} \\ T_i &= 1(z_i' \boldsymbol{\gamma} + e_{0i} > 0)\end{aligned}$$

where  $1(\cdot)$  is the indicator function. The dependent variable  $y_i$  is observed if (and only if)  $T_i = 1$ . Else it is unobserved.

For example,  $y_i$  could be a wage, which can be observed only if a person is employed. The equation for  $T_i$  is an equation specifying the probability that the person is employed.

The model is often completed by specifying that the errors are jointly normal

$$\begin{pmatrix} e_{0i} \\ e_{1i} \end{pmatrix} \sim N\left(0, \begin{pmatrix} 1 & \rho \\ \rho & \sigma^2 \end{pmatrix}\right).$$

It is presumed that we observe  $\{\mathbf{x}_i, \mathbf{z}_i, T_i\}$  for all observations.

Under the normality assumption,

$$e_{1i} = \rho e_{0i} + v_i,$$

where  $v_i$  is independent of  $e_{0i} \sim N(0, 1)$ . A useful fact about the standard normal distribution is that

$$\mathbb{E}(e_{0i} \mid e_{0i} > -x) = \lambda(x) = \frac{\phi(x)}{\Phi(x)},$$

and the function  $\lambda(x)$  is called the inverse Mills ratio.

The naive estimator of  $\beta$  is OLS regression of  $y_i$  on  $\mathbf{x}_i$  for those observations for which  $y_i$  is available. The problem is that this is equivalent to conditioning on the event  $\{T_i = 1\}$ . However,

$$\begin{aligned} \mathbb{E}(e_{1i} \mid T_i = 1, \mathbf{z}_i) &= \mathbb{E}(e_{1i} \mid \{e_{0i} > -\mathbf{z}_i' \gamma\}, \mathbf{z}_i) \\ &= \rho \mathbb{E}(e_{0i} \mid \{e_{0i} > -\mathbf{z}_i' \gamma\}, \mathbf{z}_i) + \mathbb{E}(v_i \mid \{e_{0i} > -\mathbf{z}_i' \gamma\}, \mathbf{z}_i) \\ &= \rho \lambda(\mathbf{z}_i' \gamma), \end{aligned}$$

which is non-zero. Thus

$$e_{1i} = \rho \lambda(\mathbf{z}_i' \gamma) + u_i,$$

where

$$\mathbb{E}(u_i \mid T_i = 1, \mathbf{z}_i) = 0.$$

Hence

$$y_i = \mathbf{x}_i' \beta + \rho \lambda(\mathbf{z}_i' \gamma) + u_i \quad (21.2)$$

is a valid regression equation for the observations for which  $T_i = 1$ .

Heckman (1979) observed that we could consistently estimate  $\beta$  and  $\rho$  from this equation, if  $\gamma$  were known. It is unknown, but also can be consistently estimated by a Probit model for selection. The “Heckit” estimator is thus calculated as follows

- Estimate  $\hat{\gamma}$  from a Probit, using regressors  $\mathbf{z}_i$ . The binary dependent variable is  $T_i$ .
- Estimate  $(\hat{\beta}, \hat{\rho})$  from OLS of  $y_i$  on  $\mathbf{x}_i$  and  $\lambda(\mathbf{z}_i' \hat{\gamma})$ .
- The OLS standard errors will be incorrect, as this is a two-step estimator. They can be corrected using a more complicated formula. Or, alternatively, by viewing the Probit/OLS estimation equations as a large joint GMM problem.

The Heckit estimator is frequently used to deal with problems of sample selection. However, the estimator is built on the assumption of normality, and the estimator can be quite sensitive to this assumption. Some modern econometric research is exploring how to relax the normality assumption.

The estimator can also work quite poorly if  $\lambda(\mathbf{z}_i' \hat{\gamma})$  does not have much in-sample variation. This can happen if the Probit equation does not “explain” much about the selection choice. Another potential problem is that if  $\mathbf{z}_i = \mathbf{x}_i$ , then  $\lambda(\mathbf{z}_i' \hat{\gamma})$  can be highly collinear with  $\mathbf{x}_i$ , so the second step OLS estimator will not be able to precisely estimate  $\beta$ . Based this observation, it is typically recommended to find a valid exclusion restriction: a variable should be in  $\mathbf{z}_i$  which is not in  $\mathbf{x}_i$ . If this is valid, it will ensure that  $\lambda(\mathbf{z}_i' \hat{\gamma})$  is not collinear with  $\mathbf{x}_i$ , and hence improve the second stage estimator’s precision.

## Exercises

**Exercise 21.1** Your model is

$$y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(e_i \mid \mathbf{x}_i) = 0$$

However,  $y_i^*$  is not observed. Instead only a capped version is reported. That is, the dataset contains the variable

$$y_i = \begin{cases} y_i^* & \text{if } y_i^* \leq \tau \\ \tau & \text{if } y_i^* > \tau \end{cases}$$

Suppose you regress  $y_i$  on  $x_i$  using OLS. Is OLS consistent for  $\boldsymbol{\beta}$ ? Describe the nature of the effect of the mis-measured observation on the OLS estimate.

**Exercise 21.2** Take the model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$\mathbb{E}(e_i \mid \mathbf{x}_i) = 0$$

Let  $\hat{\boldsymbol{\beta}}$  denote the OLS estimator for  $\boldsymbol{\beta}$  based on an available sample.

- (a) Suppose that the  $i^{th}$  observation is in the sample only if  $x_{1i} > 0$ , where  $x_{1i}$  is an element of  $\mathbf{x}_i$ . Assume  $\Pr(x_{1i} < 0) > 0$ .
  - i Is  $\hat{\boldsymbol{\beta}}$  consistent for  $\boldsymbol{\beta}$ ?
  - ii If not, can you obtain an expression for its probability limit?  
(For this, you may assume that  $e_i$  is independent of  $\mathbf{x}_i$  and  $N(0, \sigma^2)$ .)
- (b) Suppose that the  $i^{th}$  observation is in the sample only if  $y_i > 0$ .
  - i Is  $\hat{\boldsymbol{\beta}}$  consistent for  $\boldsymbol{\beta}$ ?
  - ii If not, can you obtain an expression for its probability limit?  
(For this, you may assume that  $e_i$  is independent of  $\mathbf{x}_i$  and  $N(0, \sigma^2)$ .)

**Exercise 21.3** The Tobit model is

$$y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + e_i$$

$$e_i \sim N(0, \sigma^2)$$

$$y_i = y_i^* 1(y_i^* \geq 0)$$

where  $1(\cdot)$  is the indicator function.

- (a) Find  $\mathbb{E}(y_i \mid \mathbf{x}_i)$ .

Note: You may use the fact that since  $e_i \sim N(0, \sigma^2)$ ,

$$\mathbb{E}(e_i 1(e_i \geq -u)) = \sigma \lambda(u/\sigma) = \sigma \phi(u/\sigma) / \Phi(u/\sigma).$$

- (b) Use the result from part (a) to suggest a NLLS estimator for the parameter  $\boldsymbol{\beta}$  given a sample  $\{y_i, \mathbf{x}_i\}$ .

**Exercise 21.4** A latent variable  $y_i^*$  is generated by

$$y_i^* = x_i \beta + e_i$$

The distribution of  $e_i$ , conditional on  $x_i$ , is  $N(0, \sigma_i^2)$ , where  $\sigma_i^2 = \gamma_0 + x_i^2 \gamma_1$  with  $\gamma_0 > 0$  and  $\gamma_1 > 0$ . The binary variable  $y_i$  equals 1 if  $y_i^* \geq 0$ , else  $y_i = 0$ . Find the log-likelihood function for the conditional distribution of  $y_i$  given  $x_i$  (the parameters are  $\beta, \gamma_0, \gamma_1$ ).

## Chapter 22

# Nonparametric Density Estimation

### 22.1 Kernel Density Estimation

Let  $X$  be a random variable with continuous distribution  $F(x)$  and density  $f(x) = \frac{d}{dx}F(x)$ . The goal is to estimate  $f(x)$  from a random sample  $(X_1, \dots, X_n)$ . While  $F(x)$  can be estimated by the EDF  $\hat{F}(x) = n^{-1} \sum_{i=1}^n 1(X_i \leq x)$ , we cannot define  $\frac{d}{dx}\hat{F}(x)$  since  $\hat{F}(x)$  is a step function. The standard **nonparametric** method to estimate  $f(x)$  is based on **smoothing** using a kernel.

While we are typically interested in estimating the entire function  $f(x)$ , we can simply focus on the problem where  $x$  is a specific fixed number, and then see how the method generalizes to estimating the entire function.

The most common methods to estimate the density  $f(x)$  is by kernel methods, which are similar to the nonparametric methods introduced in Section 17. As for kernel regression, density estimation uses kernel functions  $k(u)$ , which are density functions symmetric about zero. See Section 17 for a discussion of kernel functions.

The kernel functions are used to smooth the data. The amount of smoothing is controlled by the **bandwidth**  $h > 0$ . Define the rescaled kernel function

$$k_h(u) = \frac{1}{h} k\left(\frac{u}{h}\right).$$

The kernel density estimator of  $f(x)$  is

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n k_h(X_i - x).$$

This estimator is the average of a set of weights. If a large number of the observations  $X_i$  are near  $x$ , then the weights are relatively large and  $\hat{f}(x)$  is larger. Conversely, if only a few  $X_i$  are near  $x$ , then the weights are small and  $\hat{f}(x)$  is small. The bandwidth  $h$  controls the meaning of “near”.

Interestingly, if  $k(u)$  is a second-order kernel then  $\hat{f}(x)$  is a valid density. That is,  $\hat{f}(x) \geq 0$  for all  $x$ , and

$$\begin{aligned} \int_{-\infty}^{\infty} \hat{f}(x) dx &= \int_{-\infty}^{\infty} \frac{1}{n} \sum_{i=1}^n k_h(X_i - x) dx \\ &= \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} k_h(X_i - x) dx \\ &= \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} k(u) du = 1 \end{aligned}$$

where the second-to-last equality makes the change-of-variables  $u = (X_i - x)/h$ .

We can also calculate the moments of the density  $\hat{f}(x)$ . The mean is

$$\begin{aligned}\int_{-\infty}^{\infty} x \hat{f}(x) dx &= \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} x k_h(X_i - x) dx \\ &= \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} (X_i - uh) k(u) du \\ &= \frac{1}{n} \sum_{i=1}^n X_i \int_{-\infty}^{\infty} k(u) du + \frac{1}{n} \sum_{i=1}^n h \int_{-\infty}^{\infty} uk(u) du \\ &= \frac{1}{n} \sum_{i=1}^n X_i\end{aligned}$$

the sample mean of the  $X_i$ , where the second-to-last equality used the change-of-variables  $u = (X_i - x)/h$  which has Jacobian  $h$ .

The second moment of the estimated density is

$$\begin{aligned}\int_{-\infty}^{\infty} x^2 \hat{f}(x) dx &= \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} x^2 k_h(X_i - x) dx \\ &= \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} (X_i - uh)^2 k(u) du \\ &= \frac{1}{n} \sum_{i=1}^n X_i^2 - \frac{2}{n} \sum_{i=1}^n X_i h \int_{-\infty}^{\infty} uk(u) du + \frac{1}{n} \sum_{i=1}^n h^2 \int_{-\infty}^{\infty} u^2 k(u) du \\ &= \frac{1}{n} \sum_{i=1}^n X_i^2 + h^2 \sigma_K^2\end{aligned}$$

where

$$\sigma_K^2 = \int_{-\infty}^{\infty} u^2 k(u) du$$

is the variance of the kernel (see Section 17). It follows that the variance of the density  $\hat{f}(x)$  is

$$\begin{aligned}\int_{-\infty}^{\infty} x^2 \hat{f}(x) dx - \left( \int_{-\infty}^{\infty} x \hat{f}(x) dx \right)^2 &= \frac{1}{n} \sum_{i=1}^n X_i^2 + h^2 \sigma_K^2 - \left( \frac{1}{n} \sum_{i=1}^n X_i \right)^2 \\ &= \hat{\sigma}^2 + h^2 \sigma_K^2\end{aligned}$$

Thus the variance of the estimated density is inflated by the factor  $h^2 \sigma_K^2$  relative to the sample moment.

## 22.2 Asymptotic MSE for Kernel Estimates

For fixed  $x$  and bandwidth  $h$  observe that

$$\begin{aligned}\mathbb{E}k_h(X - x) &= \int_{-\infty}^{\infty} k_h(z - x) f(z) dz \\ &= \int_{-\infty}^{\infty} k_h(uh) f(x + hu) h du \\ &= \int_{-\infty}^{\infty} k(u) f(x + hu) du\end{aligned}$$

The second equality uses the change-of variables  $u = (z - x)/h$ . The last expression shows that the expected value is an average of  $f(z)$  locally about  $x$ .

This integral (typically) is not analytically solvable, so we approximate it using a second order Taylor expansion of  $f(x + hu)$  in the argument  $hu$  about  $hu = 0$ , which is valid as  $h \rightarrow 0$ . Thus

$$f(x + hu) \simeq f(x) + f'(x)hu + \frac{1}{2}f''(x)h^2u^2$$

and therefore

$$\begin{aligned} \mathbb{E}k_h(X - x) &\simeq \int_{-\infty}^{\infty} k(u) \left( f(x) + f'(x)hu + \frac{1}{2}f''(x)h^2u^2 \right) du \\ &= f(x) \int_{-\infty}^{\infty} k(u) du + f'(x)h \int_{-\infty}^{\infty} k(u) u du \\ &\quad + \frac{1}{2}f''(x)h^2 \int_{-\infty}^{\infty} k(u) u^2 du \\ &= f(x) + \frac{1}{2}f''(x)h^2\sigma_k^2. \end{aligned}$$

The bias of  $\hat{f}(x)$  is then

$$\text{Bias}(x) = \mathbb{E}(\hat{f}(x)) - f(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(k_h(X_i - x)) - f(x) = \frac{1}{2}f''(x)h^2\sigma_k^2.$$

We see that the bias of  $\hat{f}(x)$  at  $x$  depends on the second derivative  $f''(x)$ . The sharper the derivative, the greater the bias. Intuitively, the estimator  $\hat{f}(x)$  smooths data local to  $X_i = x$ , so is estimating a smoothed version of  $f(x)$ . The bias results from this smoothing, and is larger the greater the curvature in  $f(x)$ .

We now examine the variance of  $\hat{f}(x)$ . Since it is an average of iid random variables, using first-order Taylor approximations and the fact that  $n^{-1}$  is of smaller order than  $(nh)^{-1}$

$$\begin{aligned} \text{var}(\hat{f}(x)) &= \frac{1}{n} \text{var}(k_h(X_i - x)) \\ &= \frac{1}{n} \mathbb{E}(k_h(X_i - x)^2) - \frac{1}{n} (\mathbb{E}(k_h(X_i - x)))^2 \\ &\simeq \frac{1}{nh^2} \int_{-\infty}^{\infty} k\left(\frac{z-x}{h}\right)^2 f(z) dz - \frac{1}{n} f(x)^2 \\ &= \frac{1}{nh} \int_{-\infty}^{\infty} k(u)^2 f(x + hu) du \\ &\simeq \frac{f(x)}{nh} \int_{-\infty}^{\infty} k(u)^2 du \\ &= \frac{f(x) R_k}{nh}. \end{aligned}$$

where  $R_k = \int_{-\infty}^{\infty} k(u)^2 du$  is called the **roughness** of  $k$  (see Section 17).

Together, the asymptotic mean-squared error (AMSE) for fixed  $x$  is the sum of the approximate squared bias and approximate variance

$$\text{AMSE}_h(x) = \frac{1}{4}f''(x)^2h^4\sigma_k^4 + \frac{f(x) R_k}{nh}.$$

A global measure of precision is the asymptotic mean integrated squared error (AMISE)

$$\text{AMISE}_h = \int \text{AMSE}_h(x) dx = \frac{h^4\sigma_k^4 R(f'')}{4} + \frac{R_k}{nh}. \quad (22.1)$$

where  $R(f'') = \int (f''(x))^2 dx$  is the roughness of  $f''$ . Notice that the first term (the squared bias) is increasing in  $h$  and the second term (the variance) is decreasing in  $nh$ . Thus for the AMISE to decline with  $n$ , we need  $h \rightarrow 0$  but  $nh \rightarrow \infty$ . That is,  $h$  must tend to zero, but at a slower rate than  $n^{-1}$ .

Equation (22.1) is an asymptotic approximation to the MSE. We define the asymptotically optimal bandwidth  $h_0$  as the value which minimizes this approximate MSE. That is,

$$h_0 = \underset{h}{\operatorname{argmin}} AMISE_h$$

It can be found by solving the first order condition

$$\frac{d}{dh} AMISE_h = h^3 \sigma_K^4 R(f'') - \frac{R_k}{nh^2} = 0$$

yielding

$$h_0 = \left( \frac{R_k}{\sigma_K^4 R(f'')} \right)^{1/5} n^{-1/5}. \quad (22.2)$$

This solution takes the form  $h_0 = cn^{-1/5}$  where  $c$  is a function of  $k$  and  $f$ , but not of  $n$ . We thus say that the optimal bandwidth is of order  $O(n^{-1/5})$ . Note that this  $h$  declines to zero, but at a very slow rate.

In practice, how should the bandwidth be selected? This is a difficult problem, and there is a large literature on the subject. The asymptotically optimal choice given in (22.2) depends on  $R_k$ ,  $\sigma_k^2$ , and  $R(f'')$ . The first two are determined by the kernel function and are given in Section 17.

An obvious difficulty is that  $R(f'')$  is unknown. A classic simple solution proposed by Silverman (1986) has come to be known as the **reference bandwidth** or **Silverman's Rule-of-Thumb**. It uses formula (22.2) but replaces  $R(f'')$  with  $\hat{\sigma}^{-5} R(\phi'')$ , where  $\phi$  is the  $N(0, 1)$  distribution and  $\hat{\sigma}^2$  is an estimate of  $\sigma^2 = \operatorname{var}(X)$ . This choice for  $h$  gives an optimal rule when  $f(x)$  is normal, and gives a nearly optimal rule when  $f(x)$  is close to normal. The downside is that if the density is very far from normal, the rule-of-thumb  $h$  can be quite inefficient. We can calculate that  $R(\phi'') = 3/(8\sqrt{\pi})$ . Together with the above table, we find the reference rules for the three kernel functions introduced earlier.

Gaussian Kernel:  $h_{rule} = 1.06\hat{\sigma}n^{-1/5}$

Epanechnikov Kernel:  $h_{rule} = 2.34\hat{\sigma}n^{-1/5}$

Biweight (Quartic) Kernel:  $h_{rule} = 2.78\hat{\sigma}n^{-1/5}$

Unless you delve more deeply into kernel estimation methods the rule-of-thumb bandwidth is a good practical bandwidth choice, perhaps adjusted by visual inspection of the resulting estimate  $\hat{f}(x)$ . There are other approaches, but implementation can be delicate. I now discuss some of these choices. The **plug-in** approach is to estimate  $R(f'')$  in a first step, and then plug this estimate into the formula (22.2). This is more treacherous than may first appear, as the optimal  $h$  for estimation of the roughness  $R(f'')$  is quite different than the optimal  $h$  for estimation of  $f(x)$ . However, there are modern versions of this estimator work well, in particular the iterative method of Sheather and Jones (1991). Another popular choice for selection of  $h$  is **cross-validation**. This works by constructing an estimate of the MISE using leave-one-out estimators. There are some desirable properties of cross-validation bandwidths, but they are also known to converge very slowly to the optimal values. They are also quite ill-behaved when the data has some discretization (as is common in economics), in which case the cross-validation rule can sometimes select very small bandwidths leading to dramatically undersmoothed estimates.



# Appendix A

## Matrix Algebra

### A.1 Notation

A **scalar**  $a$  is a single number.

A **vector**  $\mathbf{a}$  is a  $k \times 1$  list of numbers, typically arranged in a column. We write this as

$$\mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_k \end{pmatrix}$$

Equivalently, a vector  $\mathbf{a}$  is an element of Euclidean  $k$  space, written as  $\mathbf{a} \in \mathbb{R}^k$ . If  $k = 1$  then  $\mathbf{a}$  is a scalar.

A **matrix**  $\mathbf{A}$  is a  $k \times r$  rectangular array of numbers, written as

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1r} \\ a_{21} & a_{22} & \cdots & a_{2r} \\ \vdots & \vdots & & \vdots \\ a_{k1} & a_{k2} & \cdots & a_{kr} \end{bmatrix}$$

By convention  $a_{ij}$  refers to the element in the  $i^{th}$  row and  $j^{th}$  column of  $\mathbf{A}$ . If  $r = 1$  then  $\mathbf{A}$  is a column vector. If  $k = 1$  then  $\mathbf{A}$  is a row vector. If  $r = k = 1$ , then  $\mathbf{A}$  is a scalar.

A standard convention (which we will follow in this text whenever possible) is to denote scalars by lower-case italics ( $a$ ), vectors by lower-case bold italics ( $\mathbf{a}$ ), and matrices by upper-case bold italics ( $\mathbf{A}$ ). Sometimes a matrix  $\mathbf{A}$  is denoted by the symbol  $(a_{ij})$ .

A matrix can be written as a set of column vectors or as a set of row vectors. That is,

$$\mathbf{A} = \begin{bmatrix} \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_r \end{bmatrix} = \begin{bmatrix} \boldsymbol{\alpha}_1 \\ \boldsymbol{\alpha}_2 \\ \vdots \\ \boldsymbol{\alpha}_k \end{bmatrix}$$

where

$$\mathbf{a}_i = \begin{bmatrix} a_{1i} \\ a_{2i} \\ \vdots \\ a_{ki} \end{bmatrix}$$

are column vectors and

$$\boldsymbol{\alpha}_j = \begin{bmatrix} a_{j1} & a_{j2} & \cdots & a_{jr} \end{bmatrix}$$

are row vectors.

The **transpose** of a matrix  $\mathbf{A}$ , denoted  $\mathbf{A}'$ ,  $\mathbf{A}^\top$ , or  $\mathbf{A}^t$ , is obtained by flipping the matrix on its diagonal. (In most of the econometrics literature, and this textbook, we use  $\mathbf{A}'$ , but in the mathematics literature  $\mathbf{A}^\top$  is the convention.) Thus

$$\mathbf{A}' = \begin{bmatrix} a_{11} & a_{21} & \cdots & a_{k1} \\ a_{12} & a_{22} & \cdots & a_{k2} \\ \vdots & \vdots & & \vdots \\ a_{1r} & a_{2r} & \cdots & a_{kr} \end{bmatrix}$$

Alternatively, letting  $\mathbf{B} = \mathbf{A}'$ , then  $b_{ij} = a_{ji}$ . Note that if  $\mathbf{A}$  is  $k \times r$ , then  $\mathbf{A}'$  is  $r \times k$ . If  $\mathbf{a}$  is a  $k \times 1$  vector, then  $\mathbf{a}'$  is a  $1 \times k$  row vector.

A matrix is **square** if  $k = r$ . A square matrix is **symmetric** if  $\mathbf{A} = \mathbf{A}'$ , which requires  $a_{ij} = a_{ji}$ . A square matrix is **diagonal** if the off-diagonal elements are all zero, so that  $a_{ij} = 0$  if  $i \neq j$ . A square matrix is **upper (lower) diagonal** if all elements below (above) the diagonal equal zero.

An important diagonal matrix is the **identity matrix**, which has ones on the diagonal. The  $k \times k$  identity matrix is denoted as

$$\mathbf{I}_k = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}.$$

A **partitioned matrix** takes the form

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} & \cdots & \mathbf{A}_{1r} \\ \mathbf{A}_{21} & \mathbf{A}_{22} & \cdots & \mathbf{A}_{2r} \\ \vdots & \vdots & & \vdots \\ \mathbf{A}_{k1} & \mathbf{A}_{k2} & \cdots & \mathbf{A}_{kr} \end{bmatrix}$$

where the  $\mathbf{A}_{ij}$  denote matrices, vectors and/or scalars.

## A.2 Complex Matrices\*

Scalars, vectors and matrices may contain real or complex numbers as entries. (However, most econometric applications exclusively use real matrices.) If all elements of a vector  $\mathbf{x}$  are real we say that  $\mathbf{x}$  is a real vector, and similarly for matrices.

Recall that a complex number can be written as  $x = a + bi$  where  $i = \sqrt{-1}$  and  $a$  and  $b$  are real numbers. Similarly a vector with complex elements can be written as  $\mathbf{x} = \mathbf{a} + \mathbf{b}i$  where  $\mathbf{a}$  and  $\mathbf{b}$  are real vectors, and a matrix with complex elements can be written as  $\mathbf{X} = \mathbf{A} + \mathbf{B}i$  where  $\mathbf{A}$  and  $\mathbf{B}$  are real matrices.

Recall that the complex conjugate of  $x = a + bi$  is  $x^* = a - bi$ . For matrices, the analogous concept is the conjugate transpose. The conjugate transpose of  $\mathbf{X} = \mathbf{A} + \mathbf{B}i$  is  $\mathbf{X}^* = \mathbf{A}' - \mathbf{B}'i$ . It is obtained by taking the transpose and taking the complex conjugate of each element.

## A.3 Matrix Addition

If the matrices  $\mathbf{A} = (a_{ij})$  and  $\mathbf{B} = (b_{ij})$  are of the same order, we define the sum

$$\mathbf{A} + \mathbf{B} = (a_{ij} + b_{ij}).$$

Matrix addition follows the commutative and associative laws:

$$\begin{aligned} \mathbf{A} + \mathbf{B} &= \mathbf{B} + \mathbf{A} \\ \mathbf{A} + (\mathbf{B} + \mathbf{C}) &= (\mathbf{A} + \mathbf{B}) + \mathbf{C}. \end{aligned}$$

## A.4 Matrix Multiplication

If  $\mathbf{A}$  is  $k \times r$  and  $c$  is real, we define their product as

$$\mathbf{A}c = c\mathbf{A} = (a_{ij}c).$$

If  $\mathbf{a}$  and  $\mathbf{b}$  are both  $k \times 1$ , then their inner product is

$$\mathbf{a}'\mathbf{b} = a_1b_1 + a_2b_2 + \cdots + a_kb_k = \sum_{j=1}^k a_jb_j.$$

Note that  $\mathbf{a}'\mathbf{b} = \mathbf{b}'\mathbf{a}$ . We say that two vectors  $\mathbf{a}$  and  $\mathbf{b}$  are **orthogonal** if  $\mathbf{a}'\mathbf{b} = 0$ .

If  $\mathbf{A}$  is  $k \times r$  and  $\mathbf{B}$  is  $r \times s$ , so that the number of columns of  $\mathbf{A}$  equals the number of rows of  $\mathbf{B}$ , we say that  $\mathbf{A}$  and  $\mathbf{B}$  are **conformable**. In this event the matrix product  $\mathbf{AB}$  is defined. Writing  $\mathbf{A}$  as a set of row vectors and  $\mathbf{B}$  as a set of column vectors (each of length  $r$ ), then the matrix product is defined as

$$\begin{aligned} \mathbf{AB} &= \begin{bmatrix} \mathbf{a}'_1 \\ \mathbf{a}'_2 \\ \vdots \\ \mathbf{a}'_k \end{bmatrix} \begin{bmatrix} \mathbf{b}_1 & \mathbf{b}_2 & \cdots & \mathbf{b}_s \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{a}'_1\mathbf{b}_1 & \mathbf{a}'_1\mathbf{b}_2 & \cdots & \mathbf{a}'_1\mathbf{b}_s \\ \mathbf{a}'_2\mathbf{b}_1 & \mathbf{a}'_2\mathbf{b}_2 & \cdots & \mathbf{a}'_2\mathbf{b}_s \\ \vdots & \vdots & & \vdots \\ \mathbf{a}'_k\mathbf{b}_1 & \mathbf{a}'_k\mathbf{b}_2 & \cdots & \mathbf{a}'_k\mathbf{b}_s \end{bmatrix}. \end{aligned}$$

Matrix multiplication is not commutative: in general  $\mathbf{AB} \neq \mathbf{BA}$ . However, it is associative and distributive:

$$\begin{aligned} \mathbf{A}(\mathbf{BC}) &= (\mathbf{AB})\mathbf{C} \\ \mathbf{A}(\mathbf{B} + \mathbf{C}) &= \mathbf{AB} + \mathbf{AC}. \end{aligned}$$

An alternative way to write the matrix product is to use matrix partitions. For example,

$$\begin{aligned} \mathbf{AB} &= \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{A}_{11}\mathbf{B}_{11} + \mathbf{A}_{12}\mathbf{B}_{21} & \mathbf{A}_{11}\mathbf{B}_{12} + \mathbf{A}_{12}\mathbf{B}_{22} \\ \mathbf{A}_{21}\mathbf{B}_{11} + \mathbf{A}_{22}\mathbf{B}_{21} & \mathbf{A}_{21}\mathbf{B}_{12} + \mathbf{A}_{22}\mathbf{B}_{22} \end{bmatrix}. \end{aligned}$$

As another example,

$$\begin{aligned} \mathbf{AB} &= \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 & \cdots & \mathbf{A}_r \end{bmatrix} \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \\ \vdots \\ \mathbf{B}_r \end{bmatrix} \\ &= \mathbf{A}_1\mathbf{B}_1 + \mathbf{A}_2\mathbf{B}_2 + \cdots + \mathbf{A}_r\mathbf{B}_r \\ &= \sum_{j=1}^r \mathbf{A}_j\mathbf{B}_j. \end{aligned}$$

An important property of the identity matrix is that if  $\mathbf{A}$  is  $k \times r$ , then  $\mathbf{A}\mathbf{I}_r = \mathbf{A}$  and  $\mathbf{I}_k\mathbf{A} = \mathbf{A}$ .

We say two matrices  $\mathbf{A}$  and  $\mathbf{B}$  are **orthogonal** if  $\mathbf{A}'\mathbf{B} = \mathbf{0}$ . This means that all columns of  $\mathbf{A}$  are orthogonal with all columns of  $\mathbf{B}$ .

The  $k \times r$  matrix  $\mathbf{H}$ ,  $r \leq k$ , is called **orthonormal** if  $\mathbf{H}'\mathbf{H} = \mathbf{I}_r$ . This means that the columns of  $\mathbf{H}$  are mutually orthogonal, and each column is normalized to have unit length.

## A.5 Trace

The **trace** of a  $k \times k$  square matrix  $\mathbf{A}$  is the sum of its diagonal elements

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^k a_{ii}.$$

Some straightforward properties for square matrices  $\mathbf{A}$  and  $\mathbf{B}$  and real  $c$  are

$$\begin{aligned}\text{tr}(c\mathbf{A}) &= c \text{tr}(\mathbf{A}) \\ \text{tr}(\mathbf{A}') &= \text{tr}(\mathbf{A}) \\ \text{tr}(\mathbf{A} + \mathbf{B}) &= \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B}) \\ \text{tr}(\mathbf{I}_k) &= k.\end{aligned}$$

Also, for  $k \times r$   $\mathbf{A}$  and  $r \times k$   $\mathbf{B}$  we have

$$\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA}). \quad (\text{A.1})$$

Indeed,

$$\begin{aligned}\text{tr}(\mathbf{AB}) &= \text{tr} \begin{bmatrix} a'_1 b_1 & a'_1 b_2 & \cdots & a'_1 b_k \\ a'_2 b_1 & a'_2 b_2 & \cdots & a'_2 b_k \\ \vdots & \vdots & & \vdots \\ a'_k b_1 & a'_k b_2 & \cdots & a'_k b_k \end{bmatrix} \\ &= \sum_{i=1}^k a'_i b_i \\ &= \sum_{i=1}^k b'_i a_i \\ &= \text{tr}(\mathbf{BA}).\end{aligned}$$

## A.6 Rank and Inverse

The rank of the  $k \times r$  matrix ( $r \leq k$ )

$$\mathbf{A} = [ \mathbf{a}_1 \quad \mathbf{a}_2 \quad \cdots \quad \mathbf{a}_r ]$$

is the number of linearly independent columns  $\mathbf{a}_j$ , and is written as  $\text{rank}(\mathbf{A})$ . We say that  $\mathbf{A}$  has full rank if  $\text{rank}(\mathbf{A}) = r$ .

A square  $k \times k$  matrix  $\mathbf{A}$  is said to be **nonsingular** if it has full rank, e.g.  $\text{rank}(\mathbf{A}) = k$ . This means that there is no  $k \times 1$   $\mathbf{c} \neq \mathbf{0}$  such that  $\mathbf{Ac} = \mathbf{0}$ .

If a square  $k \times k$  matrix  $\mathbf{A}$  is nonsingular then there exists a unique matrix  $k \times k$  matrix  $\mathbf{A}^{-1}$  called the **inverse** of  $\mathbf{A}$  which satisfies

$$\mathbf{AA}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}_k.$$

For non-singular  $\mathbf{A}$  and  $\mathbf{C}$ , some important properties include

$$\begin{aligned}\mathbf{A}\mathbf{A}^{-1} &= \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}_k \\ (\mathbf{A}^{-1})' &= (\mathbf{A}')^{-1} \\ (\mathbf{A}\mathbf{C})^{-1} &= \mathbf{C}^{-1}\mathbf{A}^{-1} \\ (\mathbf{A} + \mathbf{C})^{-1} &= \mathbf{A}^{-1}(\mathbf{A}^{-1} + \mathbf{C}^{-1})^{-1}\mathbf{C}^{-1} \\ \mathbf{A}^{-1} - (\mathbf{A} + \mathbf{C})^{-1} &= \mathbf{A}^{-1}(\mathbf{A}^{-1} + \mathbf{C}^{-1})^{-1}\mathbf{A}^{-1}.\end{aligned}$$

If a  $k \times k$  matrix  $\mathbf{H}$  is orthonormal (so that  $\mathbf{H}'\mathbf{H} = \mathbf{I}_k$ ), then  $\mathbf{H}$  is nonsingular and  $\mathbf{H}^{-1} = \mathbf{H}'$ . Furthermore,  $\mathbf{H}\mathbf{H}' = \mathbf{I}_k$  and  $\mathbf{H}'^{-1} = \mathbf{H}$ .

Another useful result for non-singular  $\mathbf{A}$  is known as the **Woodbury matrix identity**

$$(\mathbf{A} + \mathbf{BCD})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{BC}(\mathbf{C} + \mathbf{CDA}^{-1}\mathbf{BC})^{-1}\mathbf{CDA}^{-1}. \quad (\text{A.2})$$

In particular, for  $\mathbf{C} = -1$ ,  $\mathbf{B} = \mathbf{b}$  and  $\mathbf{D} = \mathbf{b}'$  for vector  $\mathbf{b}$  we find what is known as the **Sherman–Morrison formula**

$$(\mathbf{A} - \mathbf{bb}')^{-1} = \mathbf{A}^{-1} + (1 - \mathbf{b}'\mathbf{A}^{-1}\mathbf{b})^{-1}\mathbf{A}^{-1}\mathbf{bb}'\mathbf{A}^{-1}. \quad (\text{A.3})$$

The following fact about inverting partitioned matrices is quite useful.

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}^{-1} \stackrel{\text{def}}{=} \begin{bmatrix} \mathbf{A}^{11} & \mathbf{A}^{12} \\ \mathbf{A}^{21} & \mathbf{A}^{22} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{11.2}^{-1} & -\mathbf{A}_{11.2}^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \\ -\mathbf{A}_{22.1}^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1} & \mathbf{A}_{22.1}^{-1} \end{bmatrix} \quad (\text{A.4})$$

where  $\mathbf{A}_{11.2} = \mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}$  and  $\mathbf{A}_{22.1} = \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}$ . There are alternative algebraic representations for the components. For example, using the Woodbury matrix identity you can show the following alternative expressions

$$\begin{aligned}\mathbf{A}^{11} &= \mathbf{A}_{11}^{-1} + \mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{A}_{22.1}^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1} \\ \mathbf{A}^{22} &= \mathbf{A}_{22}^{-1} + \mathbf{A}_{22}^{-1}\mathbf{A}_{21}\mathbf{A}_{11.2}^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \\ \mathbf{A}^{12} &= -\mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{A}_{22.1}^{-1} \\ \mathbf{A}^{21} &= -\mathbf{A}_{22}^{-1}\mathbf{A}_{21}\mathbf{A}_{11.2}^{-1}\end{aligned}$$

Even if a matrix  $\mathbf{A}$  does not possess an inverse, we can still define the **Moore–Penrose generalized inverse**  $\mathbf{A}^{-}$  as the matrix which satisfies

$$\begin{aligned}\mathbf{A}\mathbf{A}^{-}\mathbf{A} &= \mathbf{A} \\ \mathbf{A}^{-}\mathbf{A}\mathbf{A}^{-} &= \mathbf{A}^{-} \\ \mathbf{A}\mathbf{A}^{-} &\text{ is symmetric} \\ \mathbf{A}^{-}\mathbf{A} &\text{ is symmetric}\end{aligned}$$

For any matrix  $\mathbf{A}$ , the Moore–Penrose generalized inverse  $\mathbf{A}^{-}$  exists and is unique.

For example, if

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$$

and  $\mathbf{A}_{11}^{-1}$  exists then

$$\mathbf{A}^{-} = \begin{bmatrix} \mathbf{A}_{11}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}.$$

## A.7 Determinant

The **determinant** is a measure of the volume of a square matrix. It is written as  $\det \mathbf{A}$  or  $|\mathbf{A}|$ .

While the determinant is widely used, its precise definition is rarely needed. However, we present the definition here for completeness. Let  $\mathbf{A} = (a_{ij})$  be a  $k \times k$  matrix. Let  $\pi = (j_1, \dots, j_k)$  denote a permutation of  $(1, \dots, k)$ . There are  $k!$  such permutations. There is a unique count of the number of inversions of the indices of such permutations (relative to the natural order  $(1, \dots, k)$ ), and let  $\varepsilon_\pi = +1$  if this count is even and  $\varepsilon_\pi = -1$  if the count is odd. Then the determinant of  $\mathbf{A}$  is defined as

$$\det \mathbf{A} = \sum_{\pi} \varepsilon_{\pi} a_{1j_1} a_{2j_2} \cdots a_{kj_k}.$$

For example, if  $\mathbf{A}$  is  $2 \times 2$ , then the two permutations of  $(1, 2)$  are  $(1, 2)$  and  $(2, 1)$ , for which  $\varepsilon_{(1,2)} = 1$  and  $\varepsilon_{(2,1)} = -1$ . Thus

$$\begin{aligned} \det \mathbf{A} &= \varepsilon_{(1,2)} a_{11} a_{22} + \varepsilon_{(2,1)} a_{21} a_{12} \\ &= a_{11} a_{22} - a_{12} a_{21}. \end{aligned}$$

For a square matrix  $\mathbf{A}$ , the **minor**  $M_{ij}$  of the  $ij^{th}$  element  $a_{ij}$  is the determinant of the matrix obtained by removing the  $i^{th}$  row and  $j^{th}$  column of  $\mathbf{A}$ . The **cofactor** of the  $ij^{th}$  element is  $C_{ij} = (-1)^{i+j} M_{ij}$ . An important representation known as Laplace's expansion relates the determinant of  $\mathbf{A}$  to its cofactors:

$$\det \mathbf{A} = \sum_{j=1}^k a_{ij} C_{ij}.$$

This holds for all  $i = 1, \dots, k$ . This is often presented as a method for computation of a determinant.

**Theorem A.7.1** *Properties of the determinant*

1.  $\det(\mathbf{A}) = \det(\mathbf{A}')$
2.  $\det(c\mathbf{A}) = c^k \det \mathbf{A}$
3.  $\det(\mathbf{AB}) = \det(\mathbf{BA}) = (\det \mathbf{A})(\det \mathbf{B})$
4.  $\det(\mathbf{A}^{-1}) = (\det \mathbf{A})^{-1}$
5.  $\det \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = (\det \mathbf{D}) \det(\mathbf{A} - \mathbf{BD}^{-1}\mathbf{C})$  if  $\det \mathbf{D} \neq 0$
6.  $\det \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{0} & \mathbf{D} \end{bmatrix} = \det(\mathbf{A})(\det \mathbf{D})$  and  $\det \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \det(\mathbf{A})(\det \mathbf{D})$
7. If  $\mathbf{A}$  is  $p \times q$  and  $\mathbf{B}$  is  $q \times p$  then  $\det(\mathbf{I}_p + \mathbf{AB}) = \det(\mathbf{I}_q + \mathbf{BA})$
8. If  $\mathbf{A}$  and  $\mathbf{D}$  are invertible then  $\det(\mathbf{A} - \mathbf{BD}^{-1}\mathbf{C}) = \frac{\det(\mathbf{A})}{\det(\mathbf{D})} \det(\mathbf{D} - \mathbf{CA}^{-1}\mathbf{B})$
9.  $\det \mathbf{A} \neq 0$  if and only if  $\mathbf{A}$  is nonsingular
10. If  $\mathbf{A}$  is triangular (upper or lower), then  $\det \mathbf{A} = \prod_{i=1}^k a_{ii}$
11. If  $\mathbf{A}$  is orthonormal, then  $\det \mathbf{A} = \pm 1$
12.  $\mathbf{A}^{-1} = (\det \mathbf{A})^{-1} \mathbf{C}$  where  $\mathbf{C} = (C_{ij})$  is the matrix of cofactors

## A.8 Eigenvalues

The **characteristic equation** of a  $k \times k$  square matrix  $\mathbf{A}$  is

$$\det(\lambda \mathbf{I}_k - \mathbf{A}) = 0.$$

The left side is a polynomial of degree  $k$  in  $\lambda$  so it has exactly  $k$  roots, which are not necessarily distinct and may be real or complex. They are called the **latent roots**, **characteristic roots**, or **eigenvalues** of  $\mathbf{A}$ . If  $\lambda$  is an eigenvalue of  $\mathbf{A}$ , then  $\lambda \mathbf{I}_k - \mathbf{A}$  is singular so there exists a non-zero vector  $\mathbf{h}$  such that  $(\lambda \mathbf{I}_k - \mathbf{A}) \mathbf{h} = \mathbf{0}$  or

$$\mathbf{A}\mathbf{h} = \mathbf{h}\lambda.$$

The vector  $\mathbf{h}$  is called a **latent vector**, **characteristic vector**, or **eigenvector** of  $\mathbf{A}$  corresponding to  $\lambda$ . They are typically normalized so that  $\mathbf{h}'\mathbf{h} = 1$  and thus  $\lambda = \mathbf{h}'\mathbf{A}\mathbf{h}$ .

Set  $\mathbf{H} = [\mathbf{h}_1 \cdots \mathbf{h}_k]$  and  $\mathbf{\Lambda} = \text{diag}\{\lambda_1, \dots, \lambda_k\}$ . A matrix expression is

$$\mathbf{A}\mathbf{H} = \mathbf{H}\mathbf{\Lambda}$$

We now state some useful properties.

**Theorem A.8.1** *Properties of eigenvalues. Let  $\lambda_i$  and  $\mathbf{h}_i$ ,  $i = 1, \dots, k$ , denote the  $k$  eigenvalues and eigenvectors of a square matrix  $\mathbf{A}$ .*

1.  $\det(\mathbf{A}) = \prod_{i=1}^k \lambda_i$
2.  $\text{tr}(\mathbf{A}) = \sum_{i=1}^k \lambda_i$
3.  $\mathbf{A}$  is non-singular if and only if all its eigenvalues are non-zero.
4. If  $\mathbf{A}$  has distinct eigenvalues, there exists a nonsingular matrix  $\mathbf{P}$  such that  $\mathbf{A} = \mathbf{P}^{-1}\mathbf{\Lambda}\mathbf{P}$  and  $\mathbf{P}\mathbf{A}\mathbf{P}^{-1} = \mathbf{\Lambda}$ .
5. The non-zero eigenvalues of  $\mathbf{A}\mathbf{B}$  and  $\mathbf{B}\mathbf{A}$  are identical.
6. If  $\mathbf{B}$  is non-singular then  $\mathbf{A}$  and  $\mathbf{B}^{-1}\mathbf{A}\mathbf{B}$  have the same eigenvalues.
7. If  $\mathbf{A}\mathbf{h} = \mathbf{h}\lambda$  then  $(\mathbf{I} - \mathbf{A})\mathbf{h} = \mathbf{h}(1 - \lambda)$ . So  $\mathbf{I} - \mathbf{A}$  has the eigenvalue  $1 - \lambda$  and associated eigenvector  $\mathbf{h}$ .

Most eigenvalue applications in econometrics concern the case where the matrix  $\mathbf{A}$  is real and symmetric. In this case all eigenvalues of  $\mathbf{A}$  are real and its eigenvectors are mutually orthogonal. Thus  $\mathbf{H}$  is orthonormal so  $\mathbf{H}'\mathbf{H} = \mathbf{I}_k$  and  $\mathbf{H}\mathbf{H}' = \mathbf{I}_k$ . When the eigenvalues are all real it is conventional to write them in descending order  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$ .

The following is a very important property of real symmetric matrices, which follows directly from the equations  $\mathbf{A}\mathbf{H} = \mathbf{H}\mathbf{\Lambda}$  and  $\mathbf{H}'\mathbf{H} = \mathbf{I}_k$ .

**Spectral Decomposition.** If  $\mathbf{A}$  is a  $k \times k$  real symmetric matrix, then  $\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$  where  $\mathbf{H}$  contains the eigenvectors and  $\mathbf{\Lambda}$  is a diagonal matrix with the eigenvalues on the diagonal. The eigenvalues are all real and the eigenvector matrix satisfies  $\mathbf{H}'\mathbf{H} = \mathbf{I}_k$ . The decomposition can be alternatively written as  $\mathbf{H}'\mathbf{A}\mathbf{H} = \mathbf{\Lambda}$ .

If  $\mathbf{A}$  is real, symmetric, and invertible, then by the spectral decomposition and the properties of orthonormal matrices,  $\mathbf{A}^{-1} = \mathbf{H}'^{-1}\mathbf{\Lambda}^{-1}\mathbf{H}^{-1} = \mathbf{H}\mathbf{\Lambda}^{-1}\mathbf{H}'$ . Thus the columns of  $\mathbf{H}$  are also the eigenvectors of  $\mathbf{A}^{-1}$ , and its eigenvalues are  $\lambda_1^{-1}, \lambda_2^{-1}, \dots, \lambda_k^{-1}$ .

## A.9 Positive Definite Matrices

We say that a  $k \times k$  real symmetric square matrix  $\mathbf{A}$  is **positive semi-definite** if for all  $\mathbf{c} \neq \mathbf{0}$ ,  $\mathbf{c}'\mathbf{A}\mathbf{c} \geq 0$ . This is written as  $\mathbf{A} \geq 0$ . We say that  $\mathbf{A}$  is **positive definite** if for all  $\mathbf{c} \neq \mathbf{0}$ ,  $\mathbf{c}'\mathbf{A}\mathbf{c} > 0$ . This is written as  $\mathbf{A} > 0$ .

Some properties include:

**Theorem A.9.1** *Properties of positive semi-definite matrices*

1. If  $\mathbf{A} = \mathbf{G}'\mathbf{B}\mathbf{G}$  with  $\mathbf{B} \geq 0$  and some matrix  $\mathbf{G}$ , then  $\mathbf{A}$  is positive semi-definite. (For any  $\mathbf{c} \neq \mathbf{0}$ ,  $\mathbf{c}'\mathbf{A}\mathbf{c} = \alpha'\mathbf{B}\alpha \geq 0$  where  $\alpha = \mathbf{G}\mathbf{c}$ .) If  $\mathbf{G}$  has full column rank and  $\mathbf{B} > 0$ , then  $\mathbf{A}$  is positive definite.
2. If  $\mathbf{A}$  is positive definite, then  $\mathbf{A}$  is non-singular and  $\mathbf{A}^{-1}$  exists. Furthermore,  $\mathbf{A}^{-1} > 0$ .
3.  $\mathbf{A} > 0$  if and only if it is symmetric and all its eigenvalues are positive.
4. By the spectral decomposition,  $\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$  where  $\mathbf{H}'\mathbf{H} = \mathbf{I}_k$  and  $\mathbf{\Lambda}$  is diagonal with non-negative diagonal elements. All diagonal elements of  $\mathbf{\Lambda}$  are strictly positive if (and only if)  $\mathbf{A} > 0$ .
5. The rank of  $\mathbf{A}$  equals the number of strictly positive eigenvalues.
6. If  $\mathbf{A} > 0$  then  $\mathbf{A}^{-1} = \mathbf{H}\mathbf{\Lambda}^{-1}\mathbf{H}'$ .
7. If  $\mathbf{A} \geq 0$  and  $\text{rank}(\mathbf{A}) = r \leq k$  then the Moore-Penrose generalized inverse of  $\mathbf{A}$  is  $\mathbf{A}^- = \mathbf{H}\mathbf{\Lambda}^-\mathbf{H}'$  where  $\mathbf{\Lambda}^- = \text{diag}(\lambda_1^{-1}, \lambda_2^{-1}, \dots, \lambda_r^{-1}, 0, \dots, 0)$ .
8. If  $\mathbf{A} \geq 0$  we can find a matrix  $\mathbf{B}$  such that  $\mathbf{A} = \mathbf{B}\mathbf{B}'$ . We call  $\mathbf{B}$  a **matrix square root** of  $\mathbf{A}$  and is typically written as  $\mathbf{B} = \mathbf{A}^{1/2}$ . The matrix  $\mathbf{B}$  need not be unique. One matrix square root is obtained using the spectral decomposition  $\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$ . Then  $\mathbf{B} = \mathbf{H}\mathbf{\Lambda}^{1/2}\mathbf{H}'$  is itself symmetric and positive definite and satisfies  $\mathbf{A} = \mathbf{B}\mathbf{B}$ . Another matrix square root is the Cholesky decomposition, described in Section A.14.

## A.10 Generalized Eigenvalues

Let  $\mathbf{A}$  and  $\mathbf{B}$  be  $k \times k$  matrices. The generalized characteristic equation is

$$\det(\mu\mathbf{B} - \mathbf{A}) = 0.$$

The solutions  $\mu$  are known as **generalized eigenvalues** of  $\mathbf{A}$  with respect to  $\mathbf{B}$ . Associated with each generalized eigenvalue  $\mu$  is a **generalized eigenvector**  $\mathbf{v}$  which satisfies

$$\mathbf{A}\mathbf{v} = \mathbf{B}\mathbf{v}\mu.$$

They are typically normalized so that  $\mathbf{v}'\mathbf{B}\mathbf{v} = 1$  and thus  $\mu = \mathbf{v}'\mathbf{A}\mathbf{v}$ .

A matrix expression is

$$\mathbf{A}\mathbf{V} = \mathbf{B}\mathbf{V}\mathbf{M}$$

where  $\mathbf{M} = \text{diag}\{\mu_1, \dots, \mu_k\}$ .

If  $\mathbf{A}$  and  $\mathbf{B}$  are real and symmetric then the generalized eigenvalues are real.

Suppose in addition that  $\mathbf{B}$  is invertible. Then the generalized eigenvalues of  $\mathbf{A}$  with respect to  $\mathbf{B}$  are equal to the eigenvalues of  $\mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2}$ . The generalized eigenvectors  $\mathbf{V}$  of  $\mathbf{A}$  with respect to  $\mathbf{B}$  are related to the eigenvectors  $\mathbf{H}$  of  $\mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2}$  by the relationship  $\mathbf{V} = \mathbf{B}^{-1/2}\mathbf{H}$ . This implies  $\mathbf{V}'\mathbf{B}\mathbf{V} = \mathbf{I}_k$ . Thus the generalized eigenvectors are orthogonalized with respect to the matrix  $\mathbf{B}$ .



If  $\mathbf{A}\mathbf{v} = \mathbf{B}\mathbf{v}\mu$  then  $(\mathbf{B} - \mathbf{A})\mathbf{v} = \mathbf{B}\mathbf{v}(1 - \mu)$ . So a generalized eigenvalue of  $\mathbf{B} - \mathbf{A}$  with respect to  $\mathbf{B}$  is  $1 - \mu$  with associated eigenvector  $\mathbf{v}$ .

Generalized eigenvalue equations have an interesting dual property. The following is based on Lemma A.9 of Johansen (1995).

**Theorem A.10.1** *Suppose that  $\mathbf{B}$  and  $\mathbf{C}$  are invertible  $p \times p$  and  $r \times r$  matrices, respectively, and  $\mathbf{A}$  is  $p \times r$ . Then the generalized eigenvalue problems*

$$\det(\mu\mathbf{B} - \mathbf{A}\mathbf{C}^{-1}\mathbf{A}') = 0 \quad (\text{A.5})$$

and

$$\det(\mu\mathbf{C} - \mathbf{A}'\mathbf{B}^{-1}\mathbf{A}) = 0 \quad (\text{A.6})$$

have the same non-zero generalized eigenvalues. Furthermore, for any such generalized eigenvalue  $\mu$ , if  $\mathbf{v}$  and  $\mathbf{w}$  are the associated generalized eigenvectors of (A.5) and (A.6), then

$$\mathbf{w} = \mu^{-1/2}\mathbf{C}^{-1}\mathbf{A}'\mathbf{v}. \quad (\text{A.7})$$

**Proof:.** Let  $\mu \neq 0$  be an eigenvalue of (A.5). Then using Theorem A.7.1.8

$$\begin{aligned} 0 &= \det(\mu\mathbf{B} - \mathbf{A}\mathbf{C}^{-1}\mathbf{A}') \\ &= \frac{\det(\mu\mathbf{B})}{\det(\mathbf{C})} \det(\mathbf{C} - \mathbf{A}'(\mu\mathbf{B})^{-1}\mathbf{A}) \\ &= \frac{\det(\mathbf{B})}{\det(\mathbf{C})} \det(\mu\mathbf{C} - \mathbf{A}'\mathbf{B}^{-1}\mathbf{A}). \end{aligned}$$

Since  $\det(\mathbf{B})/\det(\mathbf{C}) \neq 0$  this implies (A.7) holds. Hence  $\mu$  is an eigenvalue of (A.6), as claimed.

We next show that (A.7) is an eigenvector of (A.6). Note that the solutions to (A.5) and (A.6) satisfy

$$\mathbf{B}\mathbf{v}\mu = \mathbf{A}\mathbf{C}^{-1}\mathbf{A}'\mathbf{v} \quad (\text{A.8})$$

and

$$\mathbf{C}\mathbf{w}\mu = \mathbf{A}'\mathbf{B}^{-1}\mathbf{A}\mathbf{w} \quad (\text{A.9})$$

and are normalized so that  $\mathbf{v}'\mathbf{B}\mathbf{v} = 1$  and  $\mathbf{w}'\mathbf{C}\mathbf{w} = 1$ . We show that (A.7) satisfies (A.9). Using (A.7), we find that the left-side of (A.9) equals

$$\mathbf{C}(\mu^{-1/2}\mathbf{C}^{-1}\mathbf{A}')\mu = \mathbf{A}'\mu^{1/2} = \mathbf{A}'\mathbf{B}^{-1}\mathbf{B}\mathbf{v}\mu^{1/2} = \mathbf{A}'\mathbf{B}^{-1}\mathbf{A}\mathbf{C}^{-1}\mathbf{A}'\mathbf{v}\mu^{-1/2} = \mathbf{A}'\mathbf{B}^{-1}\mathbf{A}\mathbf{w}$$

The third equality is (A.8) and the final is (A.7). This shows that (A.9) holds and thus (A.7) is an eigenvector of (A.6) as stated. ■

## A.11 Extrema of Quadratic Forms

The extrema of quadratic forms in real symmetric matrices can be conveniently be written in terms of eigenvalues and eigenvectors.

Let  $\mathbf{A}$  denote a  $k \times k$  real symmetric matrix. Let  $\lambda_1 \geq \dots \geq \lambda_k$  be the ordered eigenvalues of  $\mathbf{A}$  and  $\mathbf{h}_1, \dots, \mathbf{h}_k$  the associated ordered eigenvectors.

We start with results for the extrema of  $\mathbf{x}'\mathbf{A}\mathbf{x}$ . Throughout this Section, when we refer to the “solution” of an extremum problem, it is the solution to the normalized expression.

- $\max_{\mathbf{x}'\mathbf{x}=1} \mathbf{x}'\mathbf{A}\mathbf{x} = \max_{\mathbf{x}} \frac{\mathbf{x}'\mathbf{A}\mathbf{x}}{\mathbf{x}'\mathbf{x}} = \lambda_1$ . The solution is  $\mathbf{x} = \mathbf{h}_1$ . (That is, the maximizer of  $\mathbf{x}'\mathbf{A}\mathbf{x}$  over  $\mathbf{x}'\mathbf{x} = 1$ .)

- $\min_{\mathbf{x}'\mathbf{x}=1} \mathbf{x}'\mathbf{A}\mathbf{x} = \min_{\mathbf{x}} \frac{\mathbf{x}'\mathbf{A}\mathbf{x}}{\mathbf{x}'\mathbf{x}} = \lambda_k$ . The solution is  $\mathbf{x} = \mathbf{h}_k$ .

Multivariate generalizations can involve either the trace or the determinant.

- $\max_{\mathbf{X}'\mathbf{X}=\mathbf{I}_\ell} \text{tr}(\mathbf{X}'\mathbf{A}\mathbf{X}) = \max_{\mathbf{X}} \text{tr}\left((\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{A}\mathbf{X})\right) = \sum_{i=1}^{\ell} \lambda_i$ .

The solution is  $\mathbf{X} = [\mathbf{h}_1, \dots, \mathbf{h}_\ell]$ .

- $\min_{\mathbf{X}'\mathbf{X}=\mathbf{I}_\ell} \text{tr}(\mathbf{X}'\mathbf{A}\mathbf{X}) = \min_{\mathbf{X}} \text{tr}\left((\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{A}\mathbf{X})\right) = \sum_{i=1}^{\ell} \lambda_{k-i+1}$ .

The solution is  $\mathbf{X} = [\mathbf{h}_{k-\ell+1}, \dots, \mathbf{h}_k]$ .

For a proof, see Theorem 11.13 of Magnus and Neudecker (1988).

Suppose as well that  $\mathbf{A} > 0$  with ordered eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$  and eigenvectors  $[\mathbf{h}_1, \dots, \mathbf{h}_k]$

- $\max_{\mathbf{X}'\mathbf{X}=\mathbf{I}_\ell} \det(\mathbf{X}'\mathbf{A}\mathbf{X}) = \max_{\mathbf{X}} \frac{\det(\mathbf{X}'\mathbf{A}\mathbf{X})}{\det(\mathbf{X}'\mathbf{X})} = \prod_{i=1}^{\ell} \lambda_i$ . The solution is  $\mathbf{X} = [\mathbf{h}_1, \dots, \mathbf{h}_\ell]$ .

- $\min_{\mathbf{X}'\mathbf{X}=\mathbf{I}_\ell} \det(\mathbf{X}'\mathbf{A}\mathbf{X}) = \min_{\mathbf{X}} \frac{\det(\mathbf{X}'\mathbf{A}\mathbf{X})}{\det(\mathbf{X}'\mathbf{X})} = \prod_{i=1}^{\ell} \lambda_{k-i+1}$ . The solution is  $\mathbf{X} = [\mathbf{h}_{k-\ell+1}, \dots, \mathbf{h}_k]$ .

- $\max_{\mathbf{X}'\mathbf{X}=\mathbf{I}_\ell} \det(\mathbf{X}'(\mathbf{I} - \mathbf{A})\mathbf{X}) = \max_{\mathbf{X}} \frac{\det(\mathbf{X}'(\mathbf{I} - \mathbf{A})\mathbf{X})}{\det(\mathbf{X}'\mathbf{X})} = \prod_{i=1}^{\ell} (1 - \lambda_{k-i+1})$ . The solution is  $\mathbf{X} = [\mathbf{h}_{k-\ell+1}, \dots, \mathbf{h}_k]$ .

- $\min_{\mathbf{X}'\mathbf{X}=\mathbf{I}_\ell} \det(\mathbf{X}'(\mathbf{I} - \mathbf{A})\mathbf{X}) = \min_{\mathbf{X}} \frac{\det(\mathbf{X}'(\mathbf{I} - \mathbf{A})\mathbf{X})}{\det(\mathbf{X}'\mathbf{X})} = \prod_{i=1}^{\ell} (1 - \lambda_i)$ . The solution is  $\mathbf{X} = [\mathbf{h}_1, \dots, \mathbf{h}_\ell]$ .

For a proof, see Theorem 11.15 of Magnus and Neudecker (1988).

We can extend the above results to incorporate generalized eigenvalue equations.

Let  $\mathbf{A}$  and  $\mathbf{B}$  be  $k \times k$  real symmetric matrices with  $\mathbf{B} > 0$ . Let  $\mu_1 \geq \dots \geq \mu_k$  be the ordered generalized eigenvalues of  $\mathbf{A}$  with respect to  $\mathbf{B}$  and  $\mathbf{v}_1, \dots, \mathbf{v}_k$  the associated ordered eigenvectors.

- $\max_{\mathbf{x}'\mathbf{B}\mathbf{x}=1} \mathbf{x}'\mathbf{A}\mathbf{x} = \max_{\mathbf{x}} \frac{\mathbf{x}'\mathbf{A}\mathbf{x}}{\mathbf{x}'\mathbf{B}\mathbf{x}} = \mu_1$ . The solution is  $\mathbf{x} = \mathbf{v}_1$ .

- $\min_{\mathbf{x}'\mathbf{B}\mathbf{x}=1} \mathbf{x}'\mathbf{A}\mathbf{x} = \min_{\mathbf{x}} \frac{\mathbf{x}'\mathbf{A}\mathbf{x}}{\mathbf{x}'\mathbf{B}\mathbf{x}} = \mu_k$ . The solution is  $\mathbf{x} = \mathbf{v}_k$ .

- $\max_{\mathbf{X}'\mathbf{B}\mathbf{X}=\mathbf{I}_\ell} \text{tr}(\mathbf{X}'\mathbf{A}\mathbf{X}) = \max_{\mathbf{X}} \text{tr}\left((\mathbf{X}'\mathbf{B}\mathbf{X})^{-1}(\mathbf{X}'\mathbf{A}\mathbf{X})\right) = \sum_{i=1}^{\ell} \mu_i$ .

The solution is  $\mathbf{X} = [\mathbf{v}_1, \dots, \mathbf{v}_\ell]$ .

- $\min_{\mathbf{X}'\mathbf{B}\mathbf{X}=\mathbf{I}_\ell} \text{tr}(\mathbf{X}'\mathbf{A}\mathbf{X}) = \min_{\mathbf{X}} \text{tr}\left((\mathbf{X}'\mathbf{B}\mathbf{X})^{-1}(\mathbf{X}'\mathbf{A}\mathbf{X})\right) = \sum_{i=1}^{\ell} \mu_{k-i+1}$ .

The solution is  $\mathbf{X} = [\mathbf{v}_{k-\ell+1}, \dots, \mathbf{v}_k]$ .

Suppose as well that  $\mathbf{A} > 0$ .

- $\max_{\mathbf{X}'\mathbf{B}\mathbf{X}=\mathbf{I}_\ell} \det(\mathbf{X}'\mathbf{A}\mathbf{X}) = \max_{\mathbf{X}} \frac{\det(\mathbf{X}'\mathbf{A}\mathbf{X})}{\det(\mathbf{X}'\mathbf{B}\mathbf{X})} = \prod_{i=1}^{\ell} \mu_i.$

The solution is  $\mathbf{X} = [\mathbf{v}_1, \dots, \mathbf{v}_\ell].$

- $\min_{\mathbf{X}'\mathbf{B}\mathbf{X}=\mathbf{I}_\ell} \det(\mathbf{X}'\mathbf{A}\mathbf{X}) = \min_{\mathbf{X}} \frac{\det(\mathbf{X}'\mathbf{A}\mathbf{X})}{\det(\mathbf{X}'\mathbf{B}\mathbf{X})} = \prod_{i=1}^{\ell} \mu_{k-i+1}.$

The solution is  $\mathbf{X} = [\mathbf{v}_{k-\ell+1}, \dots, \mathbf{v}_k].$

- $\max_{\mathbf{X}'\mathbf{B}\mathbf{X}=\mathbf{I}_\ell} \det(\mathbf{X}'(\mathbf{I} - \mathbf{A})\mathbf{X}) = \max_{\mathbf{X}} \frac{\det(\mathbf{X}'(\mathbf{I} - \mathbf{A})\mathbf{X})}{\det(\mathbf{X}'\mathbf{B}\mathbf{X})} = \prod_{i=1}^{\ell} (1 - \mu_{k-i+1}).$

The solution is  $\mathbf{X} = [\mathbf{v}_{k-\ell+1}, \dots, \mathbf{v}_k].$

- $\min_{\mathbf{X}'\mathbf{B}\mathbf{X}=\mathbf{I}_\ell} \det(\mathbf{X}'(\mathbf{I} - \mathbf{A})\mathbf{X}) = \min_{\mathbf{X}} \frac{\det(\mathbf{X}'(\mathbf{I} - \mathbf{A})\mathbf{X})}{\det(\mathbf{X}'\mathbf{B}\mathbf{X})} = \prod_{i=1}^{\ell} (1 - \mu_i).$

The solution is  $\mathbf{X} = [\mathbf{v}_1, \dots, \mathbf{v}_\ell].$

By change-of-variables, we can re-express one eigenvalue problem in terms of another. For example, let  $\mathbf{A} > 0$ ,  $\mathbf{B} > 0$ , and  $\mathbf{C} > 0$ . Then

$$\max_{\mathbf{X}} \frac{\det(\mathbf{X}'\mathbf{C}\mathbf{A}\mathbf{C}\mathbf{X})}{\det(\mathbf{X}'\mathbf{C}\mathbf{B}\mathbf{C}\mathbf{X})} = \max_{\mathbf{X}} \frac{\det(\mathbf{X}'\mathbf{A}\mathbf{X})}{\det(\mathbf{X}'\mathbf{B}\mathbf{X})}$$

and

$$\min_{\mathbf{X}} \frac{\det(\mathbf{X}'\mathbf{C}\mathbf{A}\mathbf{C}\mathbf{X})}{\det(\mathbf{X}'\mathbf{C}\mathbf{B}\mathbf{C}\mathbf{X})} = \min_{\mathbf{X}} \frac{\det(\mathbf{X}'\mathbf{A}\mathbf{X})}{\det(\mathbf{X}'\mathbf{B}\mathbf{X})}.$$

## A.12 Idempotent Matrices

A  $k \times k$  square matrix  $\mathbf{A}$  is **idempotent** if  $\mathbf{A}\mathbf{A} = \mathbf{A}$ . When  $k = 1$  the only idempotent numbers are 1 and 0. For  $k > 1$  there are many possibilities. For example, the following matrix is idempotent

$$\mathbf{A} = \begin{bmatrix} 1/2 & -1/2 \\ -1/2 & 1/2 \end{bmatrix}.$$

If  $\mathbf{A}$  is idempotent and symmetric with rank  $r$ , then it has  $r$  eigenvalues which equal 1 and  $k-r$  eigenvalues which equal 0. To see this, by the spectral decomposition we can write  $\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$  where  $\mathbf{H}$  is orthonormal and  $\mathbf{\Lambda}$  contains the eigenvalues. Then

$$\mathbf{A} = \mathbf{A}\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'\mathbf{H}\mathbf{\Lambda}\mathbf{H}' = \mathbf{H}\mathbf{\Lambda}^2\mathbf{H}'.$$

We deduce that  $\mathbf{\Lambda}^2 = \mathbf{\Lambda}$  and  $\lambda_i^2 = \lambda_i$  for  $i = 1, \dots, k$ . Hence each  $\lambda_i$  must equal either 0 or 1. Since the rank of  $\mathbf{A}$  is  $r$ , and the rank equals the number of positive eigenvalues, it follows that

$$\mathbf{\Lambda} = \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{k-r} \end{bmatrix}.$$

Thus the spectral decomposition of an idempotent matrix  $\mathbf{A}$  takes the form

$$\mathbf{A} = \mathbf{H} \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{k-r} \end{bmatrix} \mathbf{H}' \tag{A.10}$$

with  $\mathbf{H}'\mathbf{H} = \mathbf{I}_k$ . Additionally,  $\text{tr}(\mathbf{A}) = \text{rank}(\mathbf{A})$  and  $\mathbf{A}$  is positive semi-definite.

If  $\mathbf{A}$  is idempotent and symmetric with rank  $r < k$  then it does not possess an inverse, but its Moore-Penrose generalized inverse takes the simple form  $\mathbf{A}^- = \mathbf{A}$ . This can be verified by checking the conditions for the Moore-Penrose generalized inverse, for example  $\mathbf{A}\mathbf{A}^-\mathbf{A} = \mathbf{A}\mathbf{A}\mathbf{A} = \mathbf{A}$ .

If  $\mathbf{A}$  is idempotent then  $\mathbf{I} - \mathbf{A}$  is also idempotent.

One useful fact is that if  $\mathbf{A}$  is idempotent then for any conformable vector  $\mathbf{c}$ ,

$$\mathbf{c}'\mathbf{A}\mathbf{c} \leq \mathbf{c}'\mathbf{c} \quad (\text{A.11})$$

$$\mathbf{c}'(\mathbf{I} - \mathbf{A})\mathbf{c} \leq \mathbf{c}'\mathbf{c} \quad (\text{A.12})$$

To see this, note that

$$\mathbf{c}'\mathbf{c} = \mathbf{c}'\mathbf{A}\mathbf{c} + \mathbf{c}'(\mathbf{I} - \mathbf{A})\mathbf{c}.$$

Since  $\mathbf{A}$  and  $\mathbf{I} - \mathbf{A}$  are idempotent, they are both positive semi-definite, so both  $\mathbf{c}'\mathbf{A}\mathbf{c}$  and  $\mathbf{c}'(\mathbf{I} - \mathbf{A})\mathbf{c}$  are non-negative. Thus they must satisfy (A.11)-(A.12).

### A.13 Singular Values

The singular values of a  $k \times r$  real matrix  $\mathbf{A}$  are the positive square roots of the eigenvalues of  $\mathbf{A}'\mathbf{A}$ . Thus for  $j = 1, \dots, r$

$$s_j = \sqrt{\lambda_j(\mathbf{A}'\mathbf{A})}$$

Since  $\mathbf{A}'\mathbf{A}$  is positive semi-definite, its eigenvalues are non-negative. Thus singular values are always real and non-negative.

The non-zero singular values of  $\mathbf{A}$  and  $\mathbf{A}'$  are the same.

When  $\mathbf{A}$  is positive semi-definite then the singular values of  $\mathbf{A}$  correspond to its eigenvalues.

The singular value decomposition of a  $k \times r$  real matrix  $\mathbf{A}$  takes the form  $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}'$  where  $\mathbf{U}$  is  $k \times k$ ,  $\mathbf{\Lambda}$  is  $k \times r$  and  $\mathbf{V}$  is  $r \times r$ , with  $\mathbf{U}$  and  $\mathbf{V}$  orthonormal ( $\mathbf{U}'\mathbf{U} = \mathbf{I}_k$  and  $\mathbf{V}'\mathbf{V} = \mathbf{I}_r$ ) and  $\mathbf{\Lambda}$  is a diagonal matrix with the singular values of  $\mathbf{A}$  on the diagonal.

It is convention to write the singular values in descending order  $s_1 \geq s_2 \geq \dots \geq s_r$ .

### A.14 Cholesky Decomposition

For a  $k \times k$  positive definite matrix  $\mathbf{A}$ , its **Cholesky decomposition** takes the form

$$\mathbf{A} = \mathbf{L}\mathbf{L}'$$

where  $\mathbf{L}$  is **lower triangular**, and thus takes the form

$$\mathbf{L} = \begin{bmatrix} L_{11} & 0 & \cdots & 0 \\ L_{21} & L_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ L_{k1} & L_{k2} & \cdots & L_{kk} \end{bmatrix}.$$

The diagonal elements of  $\mathbf{L}$  are all strictly positive.

The Cholesky decomposition is unique (for positive definite  $\mathbf{A}$ ). One intuition is that the matrices  $\mathbf{A}$  and  $\mathbf{L}$  each have  $k(k+1)/2$  free elements.

The decomposition is very useful for a range of computations, especially when a matrix square root is required. Algorithms for computation are available in standard packages (for example, `chol` in either MATLAB or R).

Lower triangular matrices such as  $\mathbf{L}$  have special properties. One is that its determinant equals the product of the diagonal elements.

Proofs of uniqueness are algorithmic. Here is one such argument for the case  $k = 3$ . Write out

$$\begin{bmatrix} A_{11} & A_{21} & A_{31} \\ A_{21} & A_{22} & A_{32} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} = \mathbf{A} = \mathbf{L}\mathbf{L}' = \begin{bmatrix} L_{11} & 0 & 0 \\ L_{21} & L_{22} & 0 \\ L_{31} & L_{32} & L_{33} \end{bmatrix} \begin{bmatrix} L_{11} & L_{21} & L_{31} \\ 0 & L_{22} & L_{32} \\ 0 & 0 & L_{33} \end{bmatrix} \\ = \begin{bmatrix} L_{11}^2 & L_{11}L_{21} & L_{11}L_{31} \\ L_{11}L_{21} & L_{21}^2 + L_{22}^2 & L_{31}L_{21} + L_{32}L_{22} \\ L_{11}L_{31} & L_{31}L_{21} + L_{32}L_{22} & L_{31}^2 + L_{32}^2 + L_{33}^2 \end{bmatrix}$$

There are six equations, six knowns (the elements of  $\mathbf{A}$ ) and six unknowns (the elements of  $\mathbf{L}$ ). We can solve for the latter by starting with the first column, moving from top to bottom. The first element has the simple solution

$$L_{11} = \sqrt{A_{11}}.$$

This has a real solution since  $A_{11} > 0$ . Moving down, since  $L_{11}$  is known, for the entries beneath  $L_{11}$  we solve and find

$$L_{21} = \frac{A_{21}}{L_{11}} = \frac{A_{21}}{\sqrt{A_{11}}} \\ L_{31} = \frac{A_{31}}{L_{11}} = \frac{A_{31}}{\sqrt{A_{11}}}$$

Next we move to the second column. We observe that  $L_{21}$  is known. Then we solve for  $L_{22}$

$$L_{22} = \sqrt{A_{22} - L_{21}^2} = \sqrt{A_{22} - \frac{A_{21}^2}{A_{11}}}.$$

This has a real solution since  $\mathbf{A} > 0$ . Then since  $L_{22}$  is known we can move down the column to find

$$L_{32} = \frac{A_{32} - L_{31}L_{21}}{L_{22}} = \frac{A_{32} - \frac{A_{31}A_{21}}{A_{11}}}{\sqrt{A_{22} - \frac{A_{21}^2}{A_{11}}}}.$$

Finally we take the third column. All elements except  $L_{33}$  are known. So we solve to find

$$L_{33} = \sqrt{A_{33} - L_{31}^2 - L_{32}^2} = \sqrt{A_{33} - \frac{A_{31}^2}{A_{11}} - \frac{\left(A_{32} - \frac{A_{31}A_{21}}{A_{11}}\right)^2}{A_{22} - \frac{A_{21}^2}{A_{11}}}}.$$

## A.15 Matrix Calculus

Let  $\mathbf{x} = (x_1, \dots, x_k)'$  be  $k \times 1$  and  $g(\mathbf{x}) = g(x_1, \dots, x_k) : \mathbb{R}^k \rightarrow \mathbb{R}$ . The vector derivative is

$$\frac{\partial}{\partial \mathbf{x}} g(\mathbf{x}) = \begin{pmatrix} \frac{\partial}{\partial x_1} g(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial x_k} g(\mathbf{x}) \end{pmatrix}$$

and

$$\frac{\partial}{\partial \mathbf{x}'} g(\mathbf{x}) = \left( \frac{\partial}{\partial x_1} g(\mathbf{x}) \quad \cdots \quad \frac{\partial}{\partial x_k} g(\mathbf{x}) \right).$$

Some properties are now summarized.

**Theorem A.15.1** *Properties of matrix derivatives*

$$1. \quad \frac{\partial}{\partial \mathbf{x}} (\mathbf{a}'\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}'\mathbf{a}) = \mathbf{a}$$

2.  $\frac{\partial}{\partial \mathbf{x}'} (\mathbf{A}\mathbf{x}) = \mathbf{A}$
3.  $\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}'\mathbf{A}\mathbf{x}) = (\mathbf{A} + \mathbf{A}')\mathbf{x}$
4.  $\frac{\partial^2}{\partial \mathbf{x}\partial \mathbf{x}'} (\mathbf{x}'\mathbf{A}\mathbf{x}) = \mathbf{A} + \mathbf{A}'$
5.  $\frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{B}\mathbf{A}) = \mathbf{B}'$
6.  $\frac{\partial}{\partial \mathbf{A}} \log \det(\mathbf{A}) = (\mathbf{A}^{-1})'$

The final two results require some justification. Recall from Section A.5 that we can write out explicitly

$$\text{tr}(\mathbf{B}\mathbf{A}) = \sum_i \sum_j a_{ij} b_{ji}.$$

Thus if we take the derivative with respect to  $a_{ij}$  we find

$$\frac{\partial}{\partial a_{ij}} \text{tr}(\mathbf{B}\mathbf{A}) = b_{ji}.$$

which is the  $ij^{\text{th}}$  element of  $\mathbf{B}'$ , establishing part 5.

For part 6, recall Laplace's expansion

$$\det \mathbf{A} = \sum_{j=1}^k a_{ij} C_{ij}.$$

where  $C_{ij}$  is the  $ij^{\text{th}}$  cofactor of  $\mathbf{A}$ . Set  $\mathbf{C} = (C_{ij})$ . Observe that  $C_{ij}$  for  $j = 1, \dots, k$  are not functions of  $a_{ij}$ . Thus the derivative with respect to  $a_{ij}$  is

$$\frac{\partial}{\partial a_{ij}} \log \det(\mathbf{A}) = (\det \mathbf{A})^{-1} \frac{\partial}{\partial a_{ij}} \det \mathbf{A} = (\det \mathbf{A})^{-1} C_{ij}$$

Together this implies

$$\frac{\partial}{\partial \mathbf{A}} \log \det(\mathbf{A}) = (\det \mathbf{A})^{-1} \mathbf{C} = \mathbf{A}^{-1}$$

where the second equality is Theorem A.7.1.12.

## A.16 Kronecker Products and the Vec Operator

Let  $\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_n]$  be  $m \times n$ . The **vec** of  $\mathbf{A}$ , denoted by  $\text{vec}(\mathbf{A})$ , is the  $mn \times 1$  vector

$$\text{vec}(\mathbf{A}) = \begin{pmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \vdots \\ \mathbf{a}_n \end{pmatrix}.$$

Let  $\mathbf{A} = (a_{ij})$  be an  $m \times n$  matrix and let  $\mathbf{B}$  be any matrix. The **Kronecker product** of  $\mathbf{A}$  and  $\mathbf{B}$ , denoted  $\mathbf{A} \otimes \mathbf{B}$ , is the matrix

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \cdots & a_{2n}\mathbf{B} \\ \vdots & \vdots & & \vdots \\ a_{m1}\mathbf{B} & a_{m2}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{bmatrix}.$$

Some important properties are now summarized. These results hold for matrices for which all matrix multiplications are conformable.

**Theorem A.16.1** *Properties of the Kronecker product*

1.  $(\mathbf{A} + \mathbf{B}) \otimes \mathbf{C} = \mathbf{A} \otimes \mathbf{C} + \mathbf{B} \otimes \mathbf{C}$
2.  $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{AC} \otimes \mathbf{BD}$
3.  $\mathbf{A} \otimes (\mathbf{B} \otimes \mathbf{C}) = (\mathbf{A} \otimes \mathbf{B}) \otimes \mathbf{C}$
4.  $(\mathbf{A} \otimes \mathbf{B})' = \mathbf{A}' \otimes \mathbf{B}'$
5.  $\text{tr}(\mathbf{A} \otimes \mathbf{B}) = \text{tr}(\mathbf{A}) \text{tr}(\mathbf{B})$
6. If  $\mathbf{A}$  is  $m \times m$  and  $\mathbf{B}$  is  $n \times n$ ,  $\det(\mathbf{A} \otimes \mathbf{B}) = (\det(\mathbf{A}))^n (\det(\mathbf{B}))^m$
7.  $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$
8. If  $\mathbf{A} > 0$  and  $\mathbf{B} > 0$  then  $\mathbf{A} \otimes \mathbf{B} > 0$
9.  $\text{vec}(\mathbf{ABC}) = (\mathbf{C}' \otimes \mathbf{A}) \text{vec}(\mathbf{B})$
10.  $\text{tr}(\mathbf{ABCD}) = \text{vec}(\mathbf{D}')' (\mathbf{C}' \otimes \mathbf{A}) \text{vec}(\mathbf{B})$

## A.17 Vector Norms

Given any vector space  $V$  (such as Euclidean space  $\mathbb{R}^m$ ) a **norm** on  $V$  is a function  $\rho : V \rightarrow \mathbb{R}$  with the properties

1.  $\rho(c\mathbf{a}) = |c| \rho(\mathbf{a})$  for any complex number  $c$  and  $\mathbf{a} \in V$
2.  $\rho(\mathbf{a} + \mathbf{b}) \leq \rho(\mathbf{a}) + \rho(\mathbf{b})$
3. If  $\rho(\mathbf{a}) = 0$  then  $\mathbf{a} = \mathbf{0}$

A seminorm on  $V$  is a function which satisfies the first two properties. The second property is known as the triangle inequality, and it is the one property which typically needs a careful demonstration (as the other two properties typically hold by inspection).

The typical norm used for Euclidean space  $\mathbb{R}^m$  is the **Euclidean norm**

$$\begin{aligned} \|\mathbf{a}\| &= (\mathbf{a}'\mathbf{a})^{1/2} \\ &= \left( \sum_{i=1}^m a_i^2 \right)^{1/2}. \end{aligned}$$

An alternative norm is the  $p$ -norm (for  $p \geq 1$ )

$$\|\mathbf{a}\|_p = \left( \sum_{i=1}^m |a_i|^p \right)^{1/p}.$$

Special cases include the Euclidean norm ( $p = 2$ ), the 1-norm

$$\|\mathbf{a}\|_1 = \sum_{i=1}^m |a_i|$$

and the sup-norm

$$\|\mathbf{a}\|_\infty = \max(|a_1|, \dots, |a_m|).$$

For real numbers ( $m = 1$ ) these norms coincide.

Some standard inequalities for Euclidean space are now given. The Minkowski inequality given below establishes that any  $p$ -norm with  $p \geq 1$  (including the Euclidean norm) satisfies the triangle inequality and is thus a valid norm.

**Jensen's Inequality.** If  $g(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$  is convex, then for any non-negative weights  $a_j$  such that  $\sum_{j=1}^m a_j = 1$ , and any real numbers  $x_j$

$$g\left(\sum_{j=1}^m a_j x_j\right) \leq \sum_{j=1}^m a_j g(x_j). \quad (\text{A.13})$$

In particular, setting  $a_j = 1/m$ , then

$$g\left(\frac{1}{m} \sum_{j=1}^m x_j\right) \leq \frac{1}{m} \sum_{j=1}^m g(x_j). \quad (\text{A.14})$$

If  $g(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$  is concave then the inequalities in (A.13) and (A.14) are reversed.

**Weighted Geometric Mean Inequality.** For any non-negative real weights  $a_j$  such that  $\sum_{j=1}^m a_j = 1$ , and any non-negative real numbers  $x_j$

$$x_1^{a_1} x_2^{a_2} \cdots x_m^{a_m} \leq \sum_{j=1}^m a_j x_j \quad (\text{A.15})$$

**Loève's  $c_r$  Inequality.** For  $r > 0$ ,

$$\left| \sum_{j=1}^m a_j \right|^r \leq c_r \sum_{j=1}^m |a_j|^r \quad (\text{A.16})$$

where  $c_r = 1$  when  $r \leq 1$  and  $c_r = m^{r-1}$  when  $r \geq 1$ .

**$c_2$  Inequality.** For any  $m \times 1$  vectors  $\mathbf{a}$  and  $\mathbf{b}$ ,

$$(\mathbf{a} + \mathbf{b})'(\mathbf{a} + \mathbf{b}) \leq 2\mathbf{a}'\mathbf{a} + 2\mathbf{b}'\mathbf{b} \quad (\text{A.17})$$

**Hölder's Inequality.** If  $p > 1$ ,  $q > 1$ , and  $1/p + 1/q = 1$ , then for any  $m \times 1$  vectors  $\mathbf{a}$  and  $\mathbf{b}$ ,

$$\sum_{j=1}^m |a_j b_j| \leq \|\mathbf{a}\|_p \|\mathbf{b}\|_q \quad (\text{A.18})$$

**Minkowski's Inequality.** For any  $m \times 1$  vectors  $\mathbf{a}$  and  $\mathbf{b}$ , if  $p \geq 1$ , then

$$\|\mathbf{a} + \mathbf{b}\|_p \leq \|\mathbf{a}\|_p + \|\mathbf{b}\|_p \quad (\text{A.19})$$

**Schwarz Inequality.** For any  $m \times 1$  vectors  $\mathbf{a}$  and  $\mathbf{b}$ ,

$$|\mathbf{a}'\mathbf{b}| \leq \|\mathbf{a}\| \|\mathbf{b}\|. \quad (\text{A.20})$$

**Proof of Jensen's Inequality (A.13).** By the definition of convexity, for any  $\lambda \in [0, 1]$

$$g(\lambda x_1 + (1 - \lambda) x_2) \leq \lambda g(x_1) + (1 - \lambda) g(x_2). \quad (\text{A.21})$$



This implies

$$\begin{aligned} g\left(\sum_{j=1}^m a_j x_j\right) &= g\left(a_1 x_1 + (1 - a_1) \sum_{j=2}^m \frac{a_j}{1 - a_1} x_j\right) \\ &\leq a_1 g(x_1) + (1 - a_1) g\left(\sum_{j=2}^m b_j x_j\right) \end{aligned}$$

where  $b_j = a_j/(1 - a_1)$  and  $\sum_{j=2}^m b_j = 1$ . By another application of (A.21) this is bounded by

$$\begin{aligned} &a_1 g(x_1) + (1 - a_1) \left( b_2 g(x_2) + (1 - b_2) g\left(\sum_{j=2}^m c_j x_j\right) \right) \\ &= a_1 g(x_1) + a_2 g(x_2) + (1 - a_1)(1 - b_2) g\left(\sum_{j=2}^m c_j x_j\right) \end{aligned}$$

where  $c_j = b_j/(1 - b_2)$ . By repeated application of (A.21) we obtain (A.13).  $\blacksquare$

**Proof of Weighted Geometric Mean Inequality.** Since the logarithm is strictly concave, by Jensen's inequality

$$\log(x_1^{a_1} x_2^{a_2} \cdots x_m^{a_m}) = \sum_{j=1}^m a_j \log x_j \leq \log\left(\sum_{j=1}^m a_j x_j\right).$$

Applying the exponential yields (A.15).  $\blacksquare$

**Proof of Loève's  $c_r$  Inequality.** For  $r \geq 1$  this is simply a rewriting of the finite form Jensen's inequality (A.14) with  $g(u) = u^r$ . For  $r < 1$ , define  $b_j = |a_j| / \left(\sum_{j=1}^m |a_j|\right)$ . The facts that  $0 \leq b_j \leq 1$  and  $r < 1$  imply  $b_j \leq b_j^r$  and thus

$$1 = \sum_{j=1}^m b_j \leq \sum_{j=1}^m b_j^r$$

which implies

$$\left(\sum_{j=1}^m |a_j|\right)^r \leq \sum_{j=1}^m |a_j|^r.$$

The proof is completed by observing that

$$\left(\sum_{j=1}^m a_j\right)^r \leq \left(\sum_{j=1}^m |a_j|\right)^r.$$

$\blacksquare$

**Proof of  $c_2$  Inequality.** By the  $c_r$  inequality,  $(a_j + b_j)^2 \leq 2a_j^2 + 2b_j^2$ . Thus

$$\begin{aligned} (\mathbf{a} + \mathbf{b})'(\mathbf{a} + \mathbf{b}) &= \sum_{j=1}^m (a_j + b_j)^2 \\ &\leq 2 \sum_{j=1}^m a_j^2 + 2 \sum_{j=1}^m b_j^2 \\ &= 2\mathbf{a}'\mathbf{a} + 2\mathbf{b}'\mathbf{b} \end{aligned}$$

■

**Proof of Hölder's Inequality.** Set  $u_j = |a_j|^p / \|\mathbf{a}\|_p^p$  and  $v_j = |b_j|^q / \|\mathbf{b}\|_q^q$  and observe that  $\sum_{j=1}^m u_j = 1$  and  $\sum_{j=1}^m v_j = 1$ . By the weighted geometric mean inequality,

$$u_j^{1/p} v_j^{1/q} \leq \frac{u_j}{p} + \frac{v_j}{q}.$$

Then since  $\sum_{j=1}^m u_j = 1$ ,  $\sum_{j=1}^m v_j = 1$  and  $1/p + 1/q = 1$

$$\frac{\sum_{j=1}^m |a_j b_j|}{\|\mathbf{a}\|_p \|\mathbf{b}\|_q} = \sum_{j=1}^m u_j^{1/p} v_j^{1/q} \leq \sum_{j=1}^m \left( \frac{u_j}{p} + \frac{v_j}{q} \right) = 1$$

which is (A.18). ■

**Proof of Minkowski's Inequality.** Set  $q = p/(p-1)$  so that  $1/p + 1/q = 1$ . Using the triangle inequality for real numbers and two applications of Hölder's inequality

$$\begin{aligned} \|\mathbf{a} + \mathbf{b}\|_p^p &= \sum_{j=1}^m |a_j + b_j|^p \\ &= \sum_{j=1}^m |a_j + b_j| |a_j + b_j|^{p-1} \\ &\leq \sum_{j=1}^m |a_j| |a_j + b_j|^{p-1} + \sum_{j=1}^m |b_j| |a_j + b_j|^{p-1} \\ &\leq \|\mathbf{a}\|_p \left( \sum_{j=1}^m |a_j + b_j|^{(p-1)q} \right)^{1/q} + \|\mathbf{b}\|_p \left( \sum_{j=1}^m |a_j + b_j|^{(p-1)q} \right)^{1/q} \\ &= (\|\mathbf{a}\|_p + \|\mathbf{b}\|_p) \|\mathbf{a} + \mathbf{b}\|_p^{p-1} \end{aligned}$$

Solving, we find (A.19). ■

**Proof of Schwarz Inequality.** Using Hölder's inequality with  $p = q = 2$

$$|\mathbf{a}'\mathbf{b}| \leq \sum_{j=1}^m |a_j b_j| \leq \|\mathbf{a}\| \|\mathbf{b}\|$$

■

## A.18 Matrix Norms

Two common norms used for matrix spaces are the **Frobenius norm** and the **spectral norm**. We can write either as  $\|\mathbf{A}\|$ , but may write  $\|\mathbf{A}\|_F$  or  $\|\mathbf{A}\|_2$  when we want to be specific.

The **Frobenius norm** of an  $m \times k$  matrix  $\mathbf{A}$  is the Euclidean norm applied to its elements

$$\begin{aligned} \|\mathbf{A}\|_F &= \|\text{vec}(\mathbf{A})\| \\ &= (\text{tr}(\mathbf{A}'\mathbf{A}))^{1/2} \\ &= \left( \sum_{i=1}^m \sum_{j=1}^k a_{ij}^2 \right)^{1/2}. \end{aligned}$$

When  $m \times m$   $\mathbf{A}$  is real symmetric then

$$\|\mathbf{A}\|_F = \left( \sum_{\ell=1}^m \lambda_\ell^2 \right)^{1/2}$$

where  $\lambda_\ell$ ,  $\ell = 1, \dots, m$  are the eigenvalues of  $\mathbf{A}$ . To see this, by the spectral decomposition  $\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$  with  $\mathbf{H}'\mathbf{H} = \mathbf{I}$  and  $\mathbf{\Lambda} = \text{diag}\{\lambda_1, \dots, \lambda_m\}$ , so

$$\|\mathbf{A}\|_F = (\text{tr}(\mathbf{H}\mathbf{\Lambda}\mathbf{H}'\mathbf{H}\mathbf{\Lambda}\mathbf{H}'))^{1/2} = (\text{tr}(\mathbf{\Lambda}\mathbf{\Lambda}))^{1/2} = \left( \sum_{\ell=1}^m \lambda_\ell^2 \right)^{1/2}. \quad (\text{A.22})$$

A useful calculation is for any  $m \times 1$  vectors  $\mathbf{a}$  and  $\mathbf{b}$ , using (A.1),

$$\|\mathbf{a}\mathbf{b}'\|_F = \text{tr}(\mathbf{b}\mathbf{a}'\mathbf{a}\mathbf{b}')^{1/2} = (\mathbf{b}'\mathbf{b}\mathbf{a}'\mathbf{a})^{1/2} = \|\mathbf{a}\| \|\mathbf{b}\| \quad (\text{A.23})$$

and in particular

$$\|\mathbf{a}\mathbf{a}'\|_F = \|\mathbf{a}\|^2. \quad (\text{A.24})$$

The **spectral norm** of an  $m \times k$  real matrix  $\mathbf{A}$  is its largest singular value

$$\|\mathbf{A}\|_2 = s_{\max}(\mathbf{A}) = (\lambda_{\max}(\mathbf{A}'\mathbf{A}))^{1/2}$$

where  $\lambda_{\max}(\mathbf{B})$  denotes the largest eigenvalue of the matrix  $\mathbf{B}$ . Notice that

$$\lambda_{\max}(\mathbf{A}'\mathbf{A}) = \|\mathbf{A}'\mathbf{A}\|_2$$

so

$$\|\mathbf{A}\|_2 = \|\mathbf{A}'\mathbf{A}\|_2^{1/2}.$$

If  $\mathbf{A}$  is  $m \times m$  and symmetric with eigenvalues  $\lambda_j$  then

$$\|\mathbf{A}\|_2 = \max_{j \leq m} |\lambda_j|.$$

The Frobenius and spectral norms are closely related. They are equivalent when applied to a matrix of rank 1, since  $\|\mathbf{a}\mathbf{b}'\|_2 = \|\mathbf{a}\| \|\mathbf{b}\| = \|\mathbf{a}\mathbf{b}'\|_F$ . In general, for  $m \times k$  matrix  $\mathbf{A}$  with rank  $r$

$$\|\mathbf{A}\|_2 = (\lambda_{\max}(\mathbf{A}'\mathbf{A}))^{1/2} \leq \left( \sum_{j=1}^k \lambda_j(\mathbf{A}'\mathbf{A}) \right)^{1/2} = \|\mathbf{A}\|_F.$$

Since  $\mathbf{A}'\mathbf{A}$  also has rank at most  $r$ , it has at most  $r$  non-zero eigenvalues, and hence

$$\|\mathbf{A}\|_F = \left( \sum_{j=1}^k \lambda_j(\mathbf{A}'\mathbf{A}) \right)^{1/2} = \left( \sum_{j=1}^r \lambda_j(\mathbf{A}'\mathbf{A}) \right)^{1/2} \leq (r \lambda_{\max}(\mathbf{A}'\mathbf{A}))^{1/2} = \sqrt{r} \|\mathbf{A}\|_2.$$

Given any vector norm  $\|\mathbf{a}\|$  the **induced matrix norm** is defined as

$$\|\mathbf{A}\| = \sup_{\mathbf{x}'\mathbf{x}=1} \|\mathbf{A}\mathbf{x}\| = \sup_{\mathbf{x} \neq 0} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|}.$$

To see that this is a norm we need to check that it satisfies the triangle inequality. Indeed

$$\|\mathbf{A} + \mathbf{B}\| = \sup_{\mathbf{x}'\mathbf{x}=1} \|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{x}\| \leq \sup_{\mathbf{x}'\mathbf{x}=1} \|\mathbf{A}\mathbf{x}\| + \sup_{\mathbf{x}'\mathbf{x}=1} \|\mathbf{B}\mathbf{x}\| = \|\mathbf{A}\| + \|\mathbf{B}\|.$$

For any vector  $\mathbf{x}$ , by the definition of the induced norm

$$\|\mathbf{Ax}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$$

a property which is called consistent norms.

Let  $\mathbf{A}$  and  $\mathbf{B}$  be conformable and  $\|\mathbf{A}\|$  an induced matrix norm. Then using the property of consistent norms

$$\|\mathbf{AB}\| = \sup_{\mathbf{x}'\mathbf{x}=1} \|\mathbf{ABx}\| \leq \sup_{\mathbf{x}'\mathbf{x}=1} \|\mathbf{A}\| \|\mathbf{Bx}\| = \|\mathbf{A}\| \|\mathbf{B}\|.$$

A matrix norm which satisfies this property is called a **sub-multiplicative norm**, and is a matrix form of the Schwarz inequality.

Of particular interest, the matrix norm induced by the Euclidean vector norm is the spectral norm. Indeed,

$$\sup_{\mathbf{x}'\mathbf{x}=1} \|\mathbf{Ax}\|^2 = \sup_{\mathbf{x}'\mathbf{x}=1} \mathbf{x}'\mathbf{A}'\mathbf{Ax} = \lambda_{\max}(\mathbf{A}'\mathbf{A}) = \|\mathbf{A}\|_2^2.$$

It follows that the spectral norm is consistent with the Euclidean norm, and is sub-multiplicative.

## A.19 Matrix Inequalities

**Schwarz Matrix Inequality:** For any  $m \times k$  and  $k \times m$  matrices  $\mathbf{A}$  and  $\mathbf{B}$ , and either the Frobenius or spectral norm,

$$\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|. \quad (\text{A.25})$$

**Triangle Inequality:** For any  $m \times k$  matrices  $\mathbf{A}$  and  $\mathbf{B}$ , and either the Frobenius or spectral norm,

$$\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|. \quad (\text{A.26})$$

**Trace Inequality.** For any  $m \times m$  matrices  $\mathbf{A}$  and  $\mathbf{B}$  such that  $\mathbf{A}$  is symmetric and  $\mathbf{B} \geq 0$

$$\text{tr}(\mathbf{AB}) \leq \|\mathbf{A}\|_2 \text{tr}(\mathbf{B}). \quad (\text{A.27})$$

**Quadratic Inequality.** For any  $m \times 1$   $\mathbf{b}$  and  $m \times m$  symmetric matrix  $\mathbf{A}$

$$\mathbf{b}'\mathbf{Ab} \leq \|\mathbf{A}\|_2 \mathbf{b}'\mathbf{b} \quad (\text{A.28})$$

**Strong Schwarz Matrix Inequality.** For any conformable matrices  $\mathbf{A}$  and  $\mathbf{B}$

$$\|\mathbf{AB}\|_F \leq \|\mathbf{A}\|_2 \|\mathbf{B}\|_F. \quad (\text{A.29})$$

**Norm Equivalence.** For any  $m \times k$  matrix  $\mathbf{A}$  of rank  $r$

$$\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F \leq \sqrt{r} \|\mathbf{A}\|_2. \quad (\text{A.30})$$

**Eigenvalue Product Inequality.** For any  $m \times m$  real symmetric matrices  $\mathbf{A} \geq 0$  and  $\mathbf{B} \geq 0$ , the eigenvalues  $\lambda_\ell(\mathbf{AB})$  are real and satisfy

$$\lambda_{\min}(\mathbf{A}) \lambda_{\min}(\mathbf{B}) \leq \lambda_\ell(\mathbf{AB}) \leq \lambda_{\max}(\mathbf{A}) \lambda_{\max}(\mathbf{B}) \quad (\text{A.31})$$

(Zhang and Zhang, 2006, Corollary 11)

---

**Proof of Schwarz Matrix Inequality:** The inequality holds for the spectral norm since it is an induced norm. Now consider the Frobenius norm. Partition  $\mathbf{A}' = [\mathbf{a}_1, \dots, \mathbf{a}_n]$  and  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_n]$ .

Then by partitioned matrix multiplication, the definition of the Frobenius norm and the Schwarz inequality for vectors

$$\begin{aligned}
\|\mathbf{AB}\|_F &= \left\| \begin{bmatrix} \mathbf{a}'_1 \mathbf{b}_1 & \mathbf{a}'_1 \mathbf{b}_2 & \cdots \\ \mathbf{a}'_2 \mathbf{b}_1 & \mathbf{a}'_2 \mathbf{b}_2 & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix} \right\|_F \\
&\leq \left\| \begin{bmatrix} \|\mathbf{a}_1\| \|\mathbf{b}_1\| & \|\mathbf{a}_1\| \|\mathbf{b}_2\| & \cdots \\ \|\mathbf{a}_2\| \|\mathbf{b}_1\| & \|\mathbf{a}_2\| \|\mathbf{b}_2\| & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix} \right\|_F \\
&= \left( \sum_{i=1}^m \sum_{j=1}^m \|\mathbf{a}_i\|^2 \|\mathbf{b}_j\|^2 \right)^{1/2} \\
&= \left( \sum_{i=1}^m \|\mathbf{a}_i\|^2 \right)^{1/2} \left( \sum_{j=1}^m \|\mathbf{b}_j\|^2 \right)^{1/2} \\
&= \left( \sum_{i=1}^k \sum_{j=1}^m \mathbf{a}_{ji}^2 \right)^{1/2} \left( \sum_{i=1}^m \sum_{j=1}^k \|\mathbf{b}_{ji}\|^2 \right)^{1/2} \\
&= \|\mathbf{A}\|_F \|\mathbf{B}\|_F
\end{aligned}$$

■

**Proof of Triangle Inequality:** The inequality holds for the spectral norm since it is an induced norm. Now consider the Frobenius norm. Let  $\mathbf{a} = \text{vec}(\mathbf{A})$  and  $\mathbf{b} = \text{vec}(\mathbf{B})$ . Then by the definition of the Frobenius norm and the Schwarz Inequality for vectors

$$\begin{aligned}
\|\mathbf{A} + \mathbf{B}\|_F &= \|\text{vec}(\mathbf{A} + \mathbf{B})\|_F \\
&= \|\mathbf{a} + \mathbf{b}\| \\
&\leq \|\mathbf{a}\| + \|\mathbf{b}\| \\
&= \|\mathbf{A}\|_F + \|\mathbf{B}\|_F
\end{aligned}$$

■

**Proof of Trace Inequality.** By the spectral decomposition for symmetric matrices,  $\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$  where  $\mathbf{\Lambda}$  has the eigenvalues  $\lambda_j$  of  $\mathbf{A}$  on the diagonal and  $\mathbf{H}$  is orthonormal. Define  $\mathbf{C} = \mathbf{H}'\mathbf{B}\mathbf{H}$  which has non-negative diagonal elements  $C_{jj}$  since  $\mathbf{B}$  is positive semi-definite. Then

$$\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{\Lambda C}) = \sum_{j=1}^m \lambda_j C_{jj} \leq \max_j |\lambda_j| \sum_{j=1}^m C_{jj} = \|\mathbf{A}\|_2 \text{tr}(\mathbf{C})$$

where the inequality uses the fact that  $C_{jj} \geq 0$ . But note that

$$\text{tr}(\mathbf{C}) = \text{tr}(\mathbf{H}'\mathbf{B}\mathbf{H}) = \text{tr}(\mathbf{H}\mathbf{H}'\mathbf{B}) = \text{tr}(\mathbf{B})$$

since  $\mathbf{H}$  is orthonormal. Thus  $\text{tr}(\mathbf{AB}) \leq \|\mathbf{A}\|_2 \text{tr}(\mathbf{B})$  as stated. ■

**Proof of Quadratic Inequality:** In the Trace Inequality set  $\mathbf{B} = \mathbf{b}\mathbf{b}'$  and note  $\text{tr}(\mathbf{AB}) = \mathbf{b}'\mathbf{A}\mathbf{b}$  and  $\text{tr}(\mathbf{B}) = \mathbf{b}'\mathbf{b}$ . ■

**Proof of Strong Schwarz Matrix Inequality.** By the definition of the Frobenius norm, the property of the trace, the Trace Inequality (noting that both  $\mathbf{A}'\mathbf{A}$  and  $\mathbf{B}\mathbf{B}'$  are symmetric and

positive semi-definite), and the Schwarz matrix inequality

$$\begin{aligned}
 \|\mathbf{A}\mathbf{B}\|_F &= (\text{tr}(\mathbf{B}'\mathbf{A}'\mathbf{A}\mathbf{B}))^{1/2} \\
 &= (\text{tr}(\mathbf{A}'\mathbf{A}\mathbf{B}\mathbf{B}'))^{1/2} \\
 &\leq (\|\mathbf{A}'\mathbf{A}\|_2 \text{tr}(\mathbf{B}\mathbf{B}'))^{1/2} \\
 &= \|\mathbf{A}\|_2 \|\mathbf{B}\|_F.
 \end{aligned}$$

■

## Appendix B

# Probability Inequalities

The following bounds are used frequently in econometric theory, predominantly in asymptotic analysis.

**Monotone Probability Inequality.** For any events  $A$  and  $B$  such that  $A \subset B$ ,

$$\Pr(A) \leq \Pr(B). \quad (\text{B.1})$$

**Union Equality.** For any events  $A$  and  $B$ ,

$$\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B). \quad (\text{B.2})$$

**Boole's Inequality (Union Bound).** For any events  $A$  and  $B$ ,

$$\Pr(A \cup B) \leq \Pr(A) + \Pr(B). \quad (\text{B.3})$$

**Bonferroni's Inequality.** For any events  $A$  and  $B$ ,

$$\Pr(A \cap B) \geq \Pr(A) + \Pr(B) - 1. \quad (\text{B.4})$$

**Jensen's Inequality.** If  $g(\cdot) : \mathbb{R}^m \rightarrow \mathbb{R}$  is convex, then for any random vector  $\mathbf{x}$  for which

$$\mathbb{E} \|\mathbf{x}\| < \infty \text{ and } \mathbb{E} |g(\mathbf{x})| < \infty,$$

$$g(\mathbb{E}(\mathbf{x})) \leq \mathbb{E}(g(\mathbf{x})). \quad (\text{B.5})$$

If  $g(\cdot)$  concave, then the inequality is reversed.

**Conditional Jensen's Inequality.** If  $g(\cdot) : \mathbb{R}^m \rightarrow \mathbb{R}$  is convex, then for any random vectors  $(\mathbf{y}, \mathbf{x})$  for which  $\mathbb{E} \|\mathbf{y}\| < \infty$  and  $\mathbb{E} \|g(\mathbf{y})\| < \infty$ ,

$$g(\mathbb{E}(\mathbf{y} \mid \mathbf{x})) \leq \mathbb{E}(g(\mathbf{y}) \mid \mathbf{x}). \quad (\text{B.6})$$

If  $g(\cdot)$  concave, then the inequality is reversed.

**Conditional Expectation Inequality.** For any  $r \geq 1$  such that  $\mathbb{E} |y|^r < \infty$ , then

$$\mathbb{E} (|\mathbb{E}(y \mid \mathbf{x})|^r) \leq \mathbb{E} (|y|^r) < \infty. \quad (\text{B.7})$$

**Expectation Inequality.** For any random matrix  $\mathbf{Y}$  for which  $\mathbb{E} \|\mathbf{Y}\| < \infty$ ,

$$\|\mathbb{E}(\mathbf{Y})\| \leq \mathbb{E} \|\mathbf{Y}\|. \quad (\text{B.8})$$

**Hölder's Inequality.** If  $p > 1$  and  $q > 1$  and  $\frac{1}{p} + \frac{1}{q} = 1$ , then for any random  $m \times n$  matrices  $\mathbf{X}$  and  $\mathbf{Y}$ ,

$$\mathbb{E} \|\mathbf{X}'\mathbf{Y}\| \leq (\mathbb{E}(\|\mathbf{X}\|^p))^{1/p} (\mathbb{E}(\|\mathbf{Y}\|^q))^{1/q}. \quad (\text{B.9})$$

**Cauchy-Schwarz Inequality.** For any random  $m \times n$  matrices  $\mathbf{X}$  and  $\mathbf{Y}$ ,

$$\mathbb{E} \|\mathbf{X}'\mathbf{Y}\| \leq \left( \mathbb{E}(\|\mathbf{X}\|^2) \right)^{1/2} \left( \mathbb{E}(\|\mathbf{Y}\|^2) \right)^{1/2}. \quad (\text{B.10})$$

**Matrix Cauchy-Schwarz Inequality.** Tripathi (1999). For any random  $\mathbf{x} \in \mathbb{R}^m$  and  $\mathbf{y} \in \mathbb{R}^\ell$ ,

$$\mathbb{E}(\mathbf{y}\mathbf{x}') (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}\mathbf{y}') \leq \mathbb{E}(\mathbf{y}\mathbf{y}') \quad (\text{B.11})$$

**Minkowski's Inequality.** For any random  $m \times n$  matrices  $\mathbf{X}$  and  $\mathbf{Y}$ ,

$$(\mathbb{E}(\|\mathbf{X} + \mathbf{Y}\|^p))^{1/p} \leq (\mathbb{E}(\|\mathbf{X}\|^p))^{1/p} + (\mathbb{E}(\|\mathbf{Y}\|^p))^{1/p} \quad (\text{B.12})$$

**Liapunov's Inequality.** For any random  $m \times n$  matrix  $\mathbf{X}$  and  $1 \leq r \leq p$ ,

$$(\mathbb{E}(\|\mathbf{X}\|^r))^{1/r} \leq (\mathbb{E}(\|\mathbf{X}\|^p))^{1/p} \quad (\text{B.13})$$

**Markov's Inequality (standard form).** For any random vector  $\mathbf{x}$  and non-negative function  $g(\mathbf{x}) \geq 0$ ,

$$\Pr(g(\mathbf{x}) > \alpha) \leq \alpha^{-1} \mathbb{E}(g(\mathbf{x})). \quad (\text{B.14})$$

**Markov's Inequality (strong form).** For any random vector  $\mathbf{x}$  and non-negative function  $g(\mathbf{x}) \geq 0$ ,

$$\Pr(g(\mathbf{x}) > \alpha) \leq \alpha^{-1} \mathbb{E}(g(\mathbf{x}) \mathbf{1}(g(\mathbf{x}) > \alpha)). \quad (\text{B.15})$$

**Chebyshev's Inequality.** For any random variable  $x$ ,

$$\Pr(|x - \mathbb{E}x| > \alpha) \leq \frac{\text{var}(x)}{\alpha^2}. \quad (\text{B.16})$$

**Proof of Monotone Probability Inequality.** Since  $A \subset B$  then  $B = A \cup \{B \cap A^c\}$  where  $A^c$  is the complement of  $A$ . The sets  $A$  and  $\{B \cap A^c\}$  are disjoint. Thus

$$\Pr(B) = \Pr(A \cup \{B \cap A^c\}) = \Pr(A) + \Pr(B \cap A^c) \geq \Pr(A)$$

since probabilities are non-negative. Thus  $\Pr(A) \leq \Pr(B)$  as claimed.  $\blacksquare$

**Proof of Union Equality.**  $\{A \cup B\} = A \cup \{B \cap A^c\}$  where  $A$  and  $\{B \cap A^c\}$  are disjoint. Also  $B = \{B \cap A\} \cup \{B \cap A^c\}$  where  $\{B \cap A\}$  and  $\{B \cap A^c\}$  are disjoint. These two relationships imply

$$\begin{aligned} \Pr(A \cup B) &= \Pr(A) + \Pr(B \cap A^c) \\ \Pr(B) &= \Pr(B \cap A) + \Pr(B \cap A^c). \end{aligned}$$

Subtracting,

$$\Pr(A \cup B) - \Pr(B) = \Pr(A) - \Pr(B \cap A)$$



which is (B.2) upon rearrangement. ■

**Proof of Boole's Inequality.** From the Union Equality and  $\Pr(A \cap B) \geq 0$ ,

$$\begin{aligned}\Pr(A \cup B) &= \Pr(A) + \Pr(B) - \Pr(A \cap B) \\ &\leq \Pr(A) + \Pr(B)\end{aligned}$$

as claimed. ■

**Proof of Bonferroni's Inequality.** Rearranging the Union Equality and using  $\Pr(A \cup B) \leq 1$

$$\begin{aligned}\Pr(A \cap B) &= \Pr(A) + \Pr(B) - \Pr(A \cup B) \\ &\geq \Pr(A) + \Pr(B) - 1\end{aligned}$$

which is (B.4). ■

**Proof of Jensen's Inequality.** Since  $g(\mathbf{u})$  is convex, at any point  $\mathbf{u}$  there is a nonempty set of subderivatives (linear surfaces touching  $g(\mathbf{u})$  at  $\mathbf{u}$  but lying below  $g(\mathbf{u})$  for all  $\mathbf{u}$ ). Let  $a + \mathbf{b}'\mathbf{u}$  be a subderivative of  $g(\mathbf{u})$  at  $\mathbf{u} = \mathbb{E}(\mathbf{x})$ . Then for all  $\mathbf{u}$ ,  $g(\mathbf{u}) \geq a + \mathbf{b}'\mathbf{u}$  yet  $g(\mathbb{E}(\mathbf{x})) = a + \mathbf{b}'\mathbb{E}(\mathbf{x})$ . Applying expectations,  $\mathbb{E}(g(\mathbf{x})) \geq a + \mathbf{b}'\mathbb{E}(\mathbf{x}) = g(\mathbb{E}(\mathbf{x}))$ , as stated. ■

**Proof of Conditional Jensen's Inequality.** The same as the proof of Jensen's Inequality, but using conditional expectations. The conditional expectations exist since  $\mathbb{E}\|\mathbf{y}\| < \infty$  and  $\mathbb{E}\|g(\mathbf{y})\| < \infty$ . ■

**Proof of Conditional Expectation Inequality.** As the function  $|u|^r$  is convex for  $r \geq 1$ , the Conditional Jensen's inequality implies

$$|\mathbb{E}(y \mid \mathbf{x})|^r \leq \mathbb{E}(|y|^r \mid \mathbf{x}).$$

Taking unconditional expectations and the law of iterated expectations, we obtain

$$\mathbb{E}(|\mathbb{E}(y \mid \mathbf{x})|^r) \leq \mathbb{E}(\mathbb{E}(|y|^r \mid \mathbf{x})) = \mathbb{E}(|y|^r) < \infty$$

as required. ■

**Proof of Expectation Inequality.** By the Triangle inequality, for  $\lambda \in [0, 1]$ ,

$$\|\lambda \mathbf{U}_1 + (1 - \lambda) \mathbf{U}_2\| \leq \lambda \|\mathbf{U}_1\| + (1 - \lambda) \|\mathbf{U}_2\|$$

which shows that the matrix norm  $g(\mathbf{U}) = \|\mathbf{U}\|$  is convex. Applying Jensen's Inequality (B.5) we find (B.8). ■

**Proof of Hölder's Inequality.** Since  $\frac{1}{p} + \frac{1}{q} = 1$  an application of the discrete Jensen's Inequality (A.13) shows that for any real  $a$  and  $b$

$$\exp \left[ \frac{1}{p}a + \frac{1}{q}b \right] \leq \frac{1}{p} \exp(a) + \frac{1}{q} \exp(b).$$

Setting  $u = \exp(a)$  and  $v = \exp(b)$  this implies

$$u^{1/p} v^{1/q} \leq \frac{u}{p} + \frac{v}{q}$$

and this inequality holds for any  $u > 0$  and  $v > 0$ .

Set  $u = \|\mathbf{X}\|^p / \mathbb{E}(\|\mathbf{X}\|^p)$  and  $v = \|\mathbf{Y}\|^q / \mathbb{E}(\|\mathbf{Y}\|^q)$ . Note that  $\mathbb{E}(u) = \mathbb{E}(v) = 1$ . By the matrix Schwarz Inequality (A.25),  $\|\mathbf{X}'\mathbf{Y}\| \leq \|\mathbf{X}\| \|\mathbf{Y}\|$ . Thus

$$\begin{aligned} \frac{\mathbb{E} \|\mathbf{X}'\mathbf{Y}\|}{(\mathbb{E}(\|\mathbf{X}\|^p))^{1/p} (\mathbb{E}(\|\mathbf{Y}\|^q))^{1/q}} &\leq \frac{\mathbb{E}(\|\mathbf{X}\| \|\mathbf{Y}\|)}{(\mathbb{E}(\|\mathbf{X}\|^p))^{1/p} (\mathbb{E}(\|\mathbf{Y}\|^q))^{1/q}} \\ &= \mathbb{E}\left(u^{1/p} v^{1/q}\right) \\ &\leq \mathbb{E}\left(\frac{u}{p} + \frac{v}{q}\right) \\ &= \frac{1}{p} + \frac{1}{q} \\ &= 1, \end{aligned}$$

which is (B.9).  $\blacksquare$

**Proof of Cauchy-Schwarz Inequality.** Special case of Hölder's with  $p = q = 2$ .

**Proof of Matrix Cauchy-Schwarz Inequality.** Define  $e = \mathbf{y} - (\mathbb{E}(\mathbf{y}\mathbf{x}'))(\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1}\mathbf{x}$ . Note that  $\mathbb{E}(ee') \geq 0$  is positive semi-definite. We can calculate that

$$\mathbb{E}(ee') = \mathbb{E}(\mathbf{y}\mathbf{y}') - (\mathbb{E}(\mathbf{y}\mathbf{x}'))(\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1}\mathbb{E}(\mathbf{x}\mathbf{y}').$$

Since the left-hand-side is positive semi-definite, so is the right-hand-side, which means  $\mathbb{E}(\mathbf{y}\mathbf{y}') \geq (\mathbb{E}(\mathbf{y}\mathbf{x}'))(\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1}\mathbb{E}(\mathbf{x}\mathbf{y}')$  as stated.  $\blacksquare$

**Proof of Liapunov's Inequality.** The function  $g(u) = u^{p/r}$  is convex for  $u > 0$  since  $p \geq r$ . Set  $u = \|\mathbf{X}\|^r$ . By Jensen's inequality,  $g(\mathbb{E}(u)) \leq \mathbb{E}(g(u))$  or

$$(\mathbb{E}(\|\mathbf{X}\|^r))^{p/r} \leq \mathbb{E}\left((\|\mathbf{X}\|^r)^{p/r}\right) = \mathbb{E}(\|\mathbf{X}\|^p).$$

Raising both sides to the power  $1/p$  yields  $(\mathbb{E}(\|\mathbf{X}\|^r))^{1/r} \leq (\mathbb{E}(\|\mathbf{X}\|^p))^{1/p}$  as claimed.  $\blacksquare$

**Proof of Minkowski's Inequality.** Note that by rewriting, using the triangle inequality (A.26), and then Hölder's Inequality to the two expectations

$$\begin{aligned} \mathbb{E}(\|\mathbf{X} + \mathbf{Y}\|^p) &= \mathbb{E}\left(\|\mathbf{X} + \mathbf{Y}\| \|\mathbf{X} + \mathbf{Y}\|^{p-1}\right) \\ &\leq \mathbb{E}\left(\|\mathbf{X}\| \|\mathbf{X} + \mathbf{Y}\|^{p-1}\right) + \mathbb{E}\left(\|\mathbf{Y}\| \|\mathbf{X} + \mathbf{Y}\|^{p-1}\right) \\ &\leq (\mathbb{E}(\|\mathbf{X}\|^p))^{1/p} \mathbb{E}\left(\left(\|\mathbf{X} + \mathbf{Y}\|^{q(p-1)}\right)^{1/q}\right) \\ &\quad + (\mathbb{E}(\|\mathbf{Y}\|^p))^{1/p} \mathbb{E}\left(\left(\|\mathbf{X} + \mathbf{Y}\|^{q(p-1)}\right)^{1/q}\right) \\ &= \left((\mathbb{E}(\|\mathbf{X}\|^p))^{1/p} + (\mathbb{E}(\|\mathbf{Y}\|^p))^{1/p}\right) \mathbb{E}\left((\|\mathbf{X} + \mathbf{Y}\|^p)^{(p-1)/p}\right) \end{aligned}$$

where the second equality picks  $q$  to satisfy  $1/p + 1/q = 1$ , and the final equality uses this fact to make the substitution  $q = p/(p-1)$  and then collects terms. Dividing both sides by  $\mathbb{E}\left((\|\mathbf{X} + \mathbf{Y}\|^p)^{(p-1)/p}\right)$ , we obtain (B.12).  $\blacksquare$

**Proof of Markov's Inequality.** Let  $F$  denote the distribution function of  $\mathbf{x}$ . Then

$$\begin{aligned} \Pr(g(\mathbf{x}) \geq \alpha) &= \int_{\{g(\mathbf{u}) \geq \alpha\}} dF(\mathbf{u}) \\ &\leq \int_{\{g(\mathbf{u}) \geq \alpha\}} \frac{g(\mathbf{u})}{\alpha} dF(\mathbf{u}) \\ &= \alpha^{-1} \int 1(g(\mathbf{u}) > \alpha) g(\mathbf{u}) dF(\mathbf{u}) \\ &= \alpha^{-1} \mathbb{E}(g(\mathbf{x}) 1(g(\mathbf{x}) > \alpha)) \end{aligned}$$

the inequality using the region of integration  $\{g(\mathbf{u}) > \alpha\}$ . This establishes the strong form (B.15). Since  $1(g(\mathbf{x}) > \alpha) \leq 1$ , the final expression is less than  $\alpha^{-1} \mathbb{E}(g(\mathbf{x}))$ , establishing the standard form (B.14). ■

**Proof of Chebyshev's Inequality.** Define  $y = (x - \mathbb{E}x)^2$  and note that  $\mathbb{E}(y) = \text{var}(x)$ . The events  $\{|x - \mathbb{E}x| > \alpha\}$  and  $\{y > \alpha^2\}$  are equal, so by an application Markov's inequality we find

$$\Pr(|x - \mathbb{E}x| > \alpha) = \Pr(y > \alpha^2) \leq \alpha^{-2} \mathbb{E}(y) = \alpha^{-2} \text{var}(x)$$

as stated. ■

# Bibliography

- [1] Abadir, Karim M. and Jan R. Magnus (2005): *Matrix Algebra*, Cambridge University Press.
- [2] Acemoglu, Daron, Simon Johnson, James A. Robinson (2001): “The Colonial Origins of Comparative Development: An Empirical Investigation,” *American Economic Review*, 91, 1369-1401.
- [3] Acemoglu, Daron, Simon Johnson, James A. Robinson (2012): “The Colonial Origins of Comparative Development: An Empirical Investigation: Reply,” *American Economic Review*, 102, 3077–3110.
- [4] Aitken, A.C. (1935): “On least squares and linear combinations of observations,” *Proceedings of the Royal Statistical Society*, 55, 42-48.
- [5] Akaike, H. (1973): “Information theory and an extension of the maximum likelihood principle.” In B. Petroc and F. Csake, eds., *Second International Symposium on Information Theory*.
- [6] Anderson, T.W. and H. Rubin (1949): “Estimation of the parameters of a single equation in a complete system of stochastic equations,” *The Annals of Mathematical Statistics*, 20, 46-63.
- [7] Andrews, Donald W. K. (1988): “Laws of large numbers for dependent non-identically distributed random variables,” *Econometric Theory*, 4, 458-467.
- [8] Andrews, Donald W. K. (1991), “Asymptotic normality of series estimators for nonparametric and semiparametric regression models,” *Econometrica*, 59, 307-345.
- [9] Andrews, Donald W. K. (1993), “Tests for parameter instability and structural change with unknown change point,” *Econometrica*, 61, 821-8516.
- [10] Andrews, Donald W. K. and Moshe Buchinsky: (2000): “A three-step method for choosing the number of bootstrap replications,” *Econometrica*, 68, 23-51.
- [11] Andrews, Donald W. K. and Werner Ploberger (1994): “Optimal tests when a nuisance parameter is present only under the alternative,” *Econometrica*, 62, 1383-1414.
- [12] Angrist, Joshua D., Guido W. Imbens, and Donald B. Rubin (1996): “Identification of causal effects using instrumental variables,” *Journal of the American Statistical Association*, 55, 650-659.
- [13] Angrist, Joshua D. and Alan B. Krueger (1991): “Does compulsory school attendance affect schooling and earnings?” *Quarterly Journal of Economics*, 91, 444-455.
- [14] Ash, Robert B. (1972): *Real Analysis and Probability*, Academic Press.
- [15] Barro, Robert J. (1977): “Unanticipated money growth and unemployment in the United States,” *American Economic Review*, 67, 101–115

- [16] Basmann, R. L. (1957): "A generalized classical method of linear estimation of coefficients in a structural equation," *Econometrica*, 25, 77-83.
- [17] Basmann, R. L. (1960): "On finite sample distributions of generalized classical linear identifiability test statistics," *Journal of the American Statistical Association*, 55, 650-659.
- [18] Baum, Christopher F, Mark E. Schaffer, and Steven Stillman (2003): "Instrumental variables and GMM: Estimation and testing," *The Stata Journal*, 3, 1-31.
- [19] Bekker, P.A. (1994): "Alternative approximations to the distributions of instrumental variable estimators," *Econometrica*, 62, 657-681.
- [20] Billingsley, Patrick (1968): *Convergence of Probability Measures*. New York: Wiley.
- [21] Billingsley, Patrick (1995): *Probability and Measure*, 3rd Edition, New York: Wiley.
- [22] Bose, A. (1988): "Edgeworth correction by bootstrap in autoregressions," *Annals of Statistics*, 16, 1709-1722.
- [23] Box, George E. P. and Dennis R. Cox, (1964). "An analysis of transformations," *Journal of the Royal Statistical Society, Series B*, 26, 211-252.
- [24] Breusch, T.S. and A.R. Pagan (1979): "The Lagrange multiplier test and its application to model specification in econometrics," *Review of Economic Studies*, 47, 239-253.
- [25] Brown, B. W. and Whitney K. Newey (2002): "GMM, efficient bootstrapping, and improved inference," *Journal of Business and Economic Statistics*.
- [26] Card, David (1995): "Using geographic variation in college proximity to estimate the return to schooling," in *Aspects of Labor Market Behavior: Essays in Honour of John Vanderkamp*, L.N. Christofides, E.K. Grant, and R. Swidinsky, editors. Toronto: University of Toronto Press.
- [27] Carlstein, E. (1986): "The use of subseries methods for estimating the variance of a general statistic from a stationary time series," *Annals of Statistics*, 14, 1171-1179.
- [28] Casella, George and Roger L. Berger (2002): *Statistical Inference*, 2nd Edition, Duxbury Press.
- [29] Chamberlain, Gary (1987): "Asymptotic efficiency in estimation with conditional moment restrictions," *Journal of Econometrics*, 34, 305-334.
- [30] Chow, G.C. (1960): "Tests of equality between sets of coefficients in two linear regressions," *Econometrica*, 28, 591-603.
- [31] Cragg, John G. (1992): "Quasi-Aitken Estimation for Heterskedasticity of Unknown Form" *Journal of Econometrics*, 54, 179-201.
- [32] Cragg, John G. and Stephen G. Donald (1993): "Testing identifiability and specification in instrumental variable models," *Econometric Theory*, 9, 222-240.
- [33] Davidson, James (1994): *Stochastic Limit Theory: An Introduction for Econometricians*. Oxford: Oxford University Press.
- [34] Davison, A.C. and D.V. Hinkley (1997): *Bootstrap Methods and their Application*. Cambridge University Press.
- [35] De Luca, Giuseppe, Jan R. Magnus, and Franco Peracchi (2017): "Balanced variable addition in linear models" *Journal of Economic Surveys*, 31.

- [36] Dickey, D.A. and W.A. Fuller (1979): "Distribution of the estimators for autoregressive time series with a unit root," *Journal of the American Statistical Association*, 74, 427-431.
- [37] Donald Stephen G. and Whitney K. Newey (2001): "Choosing the number of instruments," *Econometrica*, 69, 1161-1191.
- [38] Duflo, Esther, Pascaline Dupas, and Michael Kremer (2011): "Peer effects, teacher incentives, and the impact of tracking: Evidence from a randomized evaluation in Kenya," *American Economic Review*, 101, 1739-1774.
- [39] Dufour, Jean-Marie (1997): "Some impossibility theorems in econometrics with applications to structural and dynamic models," *Econometrica*, 65, 1365-1387.
- [40] Durbin, James (1954): "Errors in variables," *Review of the International Statistical Institute*, 22, 23-32.
- [41] Efron, Bradley (1979): "Bootstrap methods: Another look at the jackknife," *Annals of Statistics*, 7, 1-26.
- [42] Efron, Bradley (1982): *The Jackknife, the Bootstrap, and Other Resampling Plans*. Society for Industrial and Applied Mathematics.
- [43] Efron, Bradley and R.J. Tibshirani (1993): *An Introduction to the Bootstrap*, New York: Chapman-Hall.
- [44] Eichenbaum, Martin S., Lars Peter Hansen, and Kenneth J. Singleton (1988): "A time series analysis of representative agent models of consumption and leisure choice," *The Quarterly Journal of Economics*, 103, 51-78.
- [45] Eicker, F. (1963): "Asymptotic normality and consistency of the least squares estimators for families of linear regressions," *Annals of Mathematical Statistics*, 34, 447-456.
- [46] Engle, Robert F. and Clive W. J. Granger (1987): "Co-integration and error correction: Representation, estimation and testing," *Econometrica*, 55, 251-276.
- [47] Frisch, Ragnar (1933): "Editorial," *Econometrica*, 1, 1-4.
- [48] Frisch, Ragnar and F. Waugh (1933): "Partial time regressions as compared with individual trends," *Econometrica*, 1, 387-401.
- [49] Gallant, A. Ronald and D. W. Nychka (1987): "Seminonparametric maximum likelihood estimation," *Econometrica*, 55, 363-390.
- [50] Gallant, A. Ronald and Halbert White (1988): *A Unified Theory of Estimation and Inference for Nonlinear Dynamic Models*. New York: Basil Blackwell.
- [51] Galton, Francis (1886): "Regression Towards Mediocrity in Hereditary Stature," *The Journal of the Anthropological Institute of Great Britain and Ireland*, 15, 246-263.
- [52] Goldberger, Arthur S. (1964): *Econometric Theory*, Wiley.
- [53] Goldberger, Arthur S. (1968): *Topics in Regression Analysis*, Macmillan
- [54] Goldberger, Arthur S. (1991): *A Course in Econometrics*. Cambridge: Harvard University Press.
- [55] Goffe, W.L., G.D. Ferrier and J. Rogers (1994): "Global optimization of statistical functions with simulated annealing," *Journal of Econometrics*, 60, 65-99.

- [56] Gosset, William S. (a.k.a. "Student") (1908): "The probable error of a mean," *Biometrika*, 6, 1-25.
- [57] Gauss, K. F. (1809): "Theoria motus corporum coelestium," in *Werke*, Vol. VII, 240-254.
- [58] Granger, Clive W. J. (1969): "Investigating causal relations by econometric models and cross-spectral methods," *Econometrica*, 37, 424-438.
- [59] Granger, Clive W. J. (1981): "Some properties of time series data and their use in econometric specification," *Journal of Econometrics*, 16, 121-130.
- [60] Granger, Clive W. J. and Timo Teräsvirta (1993): *Modelling Nonlinear Economic Relationships*, Oxford University Press, Oxford.
- [61] Gregory, A. and M. Veall (1985): "On formulating Wald tests of nonlinear restrictions," *Econometrica*, 53, 1465-1468,
- [62] Haavelmo, T. (1944): "The probability approach in econometrics," *Econometrica*, supplement, 12.
- [63] Hall, A. R. (2000): "Covariance matrix estimation and the power of the overidentifying restrictions test," *Econometrica*, 68, 1517-1527.
- [64] Hall, B. H. and R. E. Hall (1993): "The Value and Performance of U.S. Corporations" (1993) *Brookings Papers on Economic Activity*, 1-49.
- [65] Hall, Peter (1992): *The Bootstrap and Edgeworth Expansion*, New York: Springer-Verlag.
- [66] Hall, Peter (1994): "Methodology and theory for the bootstrap," *Handbook of Econometrics*, Vol. IV, eds. R.F. Engle and D.L. McFadden. New York: Elsevier Science.
- [67] Hall, Peter and Joel L. Horowitz (1996): "Bootstrap critical values for tests based on Generalized-Method-of-Moments estimation," *Econometrica*, 64, 891-916.
- [68] Hahn, Jinyong (1996): "A note on bootstrapping generalized method of moments estimators," *Econometric Theory*, 12, 187-197.
- [69] Hamilton, James D. (1994) *Time Series Analysis*.
- [70] Hansen, Bruce E. (1992): "Efficient estimation and testing of cointegrating vectors in the presence of deterministic trends," *Journal of Econometrics*, 53, 87-121.
- [71] Hansen, Bruce E. (1996): "Inference when a nuisance parameter is not identified under the null hypothesis," *Econometrica*, 64, 413-430.
- [72] Hansen, Bruce E. (1999): "Threshold effects in non-dynamic panels: Estimation, testing and inference," *Journal of Econometrics*, 93, 345-368.
- [73] Hansen, Bruce E. (2006): "Edgeworth expansions for the Wald and GMM statistics for nonlinear restrictions," *Econometric Theory and Practice: Frontiers of Analysis and Applied Research*, edited by Dean Corbae, Steven N. Durlauf and Bruce E. Hansen. Cambridge University Press.
- [74] Hansen, Bruce E. and Seojeong Lee (2018): "Inference for iterated GMM under misspecification and clustering", working paper.
- [75] Hansen, Christopher B. (2007): "Asymptotic properties of a robust variance matrix estimator for panel data when  $T$  is large," *Journal of Econometrics*, 141, 595-620.

- [76] Hansen, Lars Peter (1982): "Large sample properties of generalized method of moments estimators," *Econometrica*, 50, 1029-1054.
- [77] Hansen, Lars Peter, John Heaton, and A. Yaron (1996): "Finite sample properties of some alternative GMM estimators," *Journal of Business and Economic Statistics*, 14, 262-280.
- [78] Hausman, J.A. (1978): "Specification tests in econometrics," *Econometrica*, 46, 1251-1271.
- [79] Heckman, J. (1979): "Sample selection bias as a specification error," *Econometrica*, 47, 153-161.
- [80] Hinkley, D. V. (1977): "Jackknifing in unbalanced situations," *Technometrics*, 19, 285-292.
- [81] Horn, S.D., R.A. Horn, and D.B. Duncan. (1975) "Estimating heteroscedastic variances in linear model," *Journal of the American Statistical Association*, 70, 380-385.
- [82] Horowitz, Joel (2001): "The Bootstrap," *Handbook of Econometrics*, Vol. 5, J.J. Heckman and E.E. Leamer, eds., Elsevier Science, 3159-3228.
- [83] Imbens, Guido W. (1997): "One step estimators for over-identified generalized method of moments models," *Review of Economic Studies*, 64, 359-383.
- [84] Imbens, Guido W., and Joshua D. Angrist (1994): "Identification and estimation of local average treatment effects," *Econometrica*, 62, 467-476.
- [85] Imbens, G.W., R.H. Spady and P. Johnson (1998): "Information theoretic approaches to inference in moment condition models," *Econometrica*, 66, 333-357.
- [86] Jarque, C.M. and A.K. Bera (1980): "Efficient tests for normality, homoskedasticity and serial independence of regression residuals," *Economic Letters*, 6, 255-259.
- [87] Johansen, S. (1988): "Statistical analysis of cointegrating vectors," *Journal of Economic Dynamics and Control*, 12, 231-254.
- [88] Johansen, S. (1991): "Estimation and hypothesis testing of cointegration vectors in the presence of linear trend," *Econometrica*, 59, 1551-1580.
- [89] Johansen, S. (1995): *Likelihood-Based Inference in Cointegrated Vector Auto-Regressive Models*, Oxford University Press.
- [90] Johansen, S. and K. Juselius (1992): "Testing structural hypotheses in a multivariate cointegration analysis of the PPP and the UIP for the UK," *Journal of Econometrics*, 53, 211-244.
- [91] Kilian, Lutz and Helmut Lütkepohl: (2017): *Structural Vector Autoregressive Analysis*, Cambridge University Press, forthcoming.
- [92] Kitamura, Y. (2001): "Asymptotic optimality and empirical likelihood for testing moment restrictions," *Econometrica*, 69, 1661-1672.
- [93] Kitamura, Y. and M. Stutzer (1997): "An information-theoretic alternative to generalized method of moments," *Econometrica*, 65, 861-874..
- [94] Koenker, Roger (2005): *Quantile Regression*. Cambridge University Press.
- [95] Kunsch, H.R. (1989): "The jackknife and the bootstrap for general stationary observations," *Annals of Statistics*, 17, 1217-1241.



- [96] Kwiatkowski, D., P.C.B. Phillips, P. Schmidt, and Y. Shin (1992): "Testing the null hypothesis of stationarity against the alternative of a unit root: How sure are we that economic time series have a unit root?" *Journal of Econometrics*, 54, 159-178.
- [97] Lafontaine, F. and K.J. White (1986): "Obtaining any Wald statistic you want," *Economics Letters*, 21, 35-40.
- [98] Legendre, Adrien-Marie (1805): *Nouvelles methodes pour la determination des orbites de cometes [New Methods for the Determination of the Orbits of Comets]*, Pris: F. Didot.
- [99] Lehmann, E.L. and George Casella (1998): *Theory of Point Estimation*, 2<sup>nd</sup> Edition, Springer.
- [100] Lehmann, E.L. and Joseph P. Romano (2005): *Testing Statistical Hypotheses*, 3<sup>rd</sup> Edition, Springer.
- [101] Lindeberg, Jarl Waldemar, (1922): "Eine neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeitsrechnung," *Mathematische Zeitschrift*, 15, 211-225.
- [102] Li, Qi and Jeffrey Racine (2007) *Nonparametric Econometrics*.
- [103] Lovell, M.C. (1963): "Seasonal adjustment of economic time series," *Journal of the American Statistical Association*, 58, 993-1010.
- [104] MacKinnon, James G. (1990): "Critical values for cointegration," in Engle, R.F. and C.W. Granger (eds.) *Long-Run Economic Relationships: Readings in Cointegration*, Oxford, Oxford University Press.
- [105] MacKinnon, James G. and Halbert White (1985): "Some heteroskedasticity-consistent covariance matrix estimators with improved finite sample properties," *Journal of Econometrics*, 29, 305-325.
- [106] Magnus, J. R., and H. Neudecker (1988): *Matrix Differential Calculus with Applications in Statistics and Econometrics*, New York: John Wiley and Sons.
- [107] Mankiw, N. Gregory, David Romer, and David N. Weil (1992): "A contribution to the empirics of economic growth," *The Quarterly Journal of Economics*, 107, 407-437.
- [108] Mann, H.B. and A. Wald (1943). "On stochastic limit and order relationships," *The Annals of Mathematical Statistics* 14, 217-226.
- [109] Muirhead, R.J. (1982): *Aspects of Multivariate Statistical Theory*. New York: Wiley.
- [110] Nelder, J. and R. Mead (1965): "A simplex method for function minimization," *Computer Journal*, 7, 308-313.
- [111] Nerlove, Marc (1963): "Returns to Scale in Electricity Supply," Chapter 7 of *Measurement in Economics* (C. Christ, et al, eds.). Stanford: Stanford University Press, 167-198.
- [112] Newey, Whitney K. (1990): "Semiparametric efficiency bounds," *Journal of Applied Econometrics*, 5, 99-135.
- [113] Newey, Whitney K. (1995): "Generalized method of moments specification testing," *Journal of Econometrics*, 29, 229-256.
- [114] Newey, Whitney K. (1997): "Convergence rates and asymptotic normality for series estimators," *Journal of Econometrics*, 79, 147-168.

- [115] Newey, Whitney K. and Daniel L. McFadden (1994): "Large Sample Estimation and Hypothesis Testing," in Robert Engle and Daniel McFadden, (eds.) *Handbook of Econometrics*, vol. IV, 2111-2245, North Holland: Amsterdam.
- [116] Newey, Whitney K. and Kenneth D. West (1987): "Hypothesis testing with efficient method of moments estimation," *International Economic Review*, 28, 777-787.
- [117] Owen, Art B. (1988): "Empirical likelihood ratio confidence intervals for a single functional," *Biometrika*, 75, 237-249.
- [118] Owen, Art B. (2001): *Empirical Likelihood*. New York: Chapman & Hall.
- [119] Pagan, Adrian (1984): "Econometric issues in the analysis of regressions with generated regressors," *International Economic Review*, 25, 221-247.
- [120] Pagan, Adrian (1986): "Two stage and related estimators and their applications," *Review of Economic Studies*, 53, 517-538.
- [121] Park, Joon Y. and Peter C. B. Phillips (1988): "On the formulation of Wald tests of nonlinear restrictions," *Econometrica*, 56, 1065-1083,
- [122] Phillips, Peter C.B. (1983): "Exact small sample theory in the simultaneous equations model," *Handbook of Econometrics, Volume 1*, edited by Z. Griliches and M. D. Intriligator, North-Holland.
- [123] Phillips, Peter C.B. and Sam Ouliaris (1990): "Asymptotic properties of residual based tests for cointegration," *Econometrica*, 58, 165-193.
- [124] Politis, D.N. and J.P. Romano (1996): "The stationary bootstrap," *Journal of the American Statistical Association*, 89, 1303-1313.
- [125] Potscher, B.M. (1991): "Effects of model selection on inference," *Econometric Theory*, 7, 163-185.
- [126] Qin, J. and J. Lawless (1994): "Empirical likelihood and general estimating equations," *The Annals of Statistics*, 22, 300-325.
- [127] Ramsey, J. B. (1969): "Tests for specification errors in classical linear least-squares regression analysis," *Journal of the Royal Statistical Society, Series B*, 31, 350-371.
- [128] Rudin, W. (1987): *Real and Complex Analysis*, 3rd edition. New York: McGraw-Hill.
- [129] Runge, Carl (1901): "Über empirische Funktionen und die Interpolation zwischen äquidistanten Ordinaten," *Zeitschrift für Mathematik und Physik*, 46, 224-243.
- [130] Said, S.E. and D.A. Dickey (1984): "Testing for unit roots in autoregressive-moving average models of unknown order," *Biometrika*, 71, 599-608.
- [131] Sargan, J.D. (1958): "The estimation of economic relationships using instrumental variables," *Econometrica*, 26, 393-415.
- [132] Secrist, Horace (1933): *The Triumph of Mediocrity in Business*. Evanston: Northwestern University.
- [133] Shao, Jun and D. Tu (1995): *The Jackknife and Bootstrap*. NY: Springer.
- [134] Shao, Jun (2003): *Mathematical Statistics*, 2nd edition, Springer.

- [135] Sheather, S.J. and M.C. Jones (1991): "A reliable data-based bandwidth selection method for kernel density estimation, *Journal of the Royal Statistical Society, Series B*, 53, 683-690.
- [136] Shin, Y. (1994): "A residual-based test of the null of cointegration against the alternative of no cointegration," *Econometric Theory*, 10, 91-115.
- [137] Silverman, B.W. (1986): *Density Estimation for Statistics and Data Analysis*. London: Chapman and Hall.
- [138] Sims, C.A. (1972): "Money, income and causality," *American Economic Review*, 62, 540-552.
- [139] Sims, C.A. (1980): "Macroeconomics and reality," *Econometrica*, 48, 1-48.
- [140] Staiger, D. and James H. Stock (1997): "Instrumental variables regression with weak instruments," *Econometrica*, 65, 557-586.
- [141] Stock, James H. (1987): "Asymptotic properties of least squares estimators of cointegrating vectors," *Econometrica*, 55, 1035-1056.
- [142] Stock, James H. (1991): "Confidence intervals for the largest autoregressive root in U.S. macroeconomic time series," *Journal of Monetary Economics*, 28, 435-460.
- [143] Stock, James H. and Jonathan H. Wright (2000): "GMM with weak identification," *Econometrica*, 68, 1055-1096.
- [144] Stock, James H. and Mark W. Watson (2014): *Introduction to Econometrics*, 3<sup>rd</sup> edition, Pearson.
- [145] Stock, James H. and Motohiro Yogo (2005): "Testing for weak instruments in linear IV regression," in *Identification and Inference for Econometric Models: Essays in Honor of Thomas Rothenberg*, eds Donald W.K. Andrews and James H. Stock, Cambridge University Press, 80-108.
- [146] Stone, Marshall H. (1937): "Applications of the Theory of Boolean Rings to General Topology," *Transactions of the American Mathematical Society*, 41, 375-481.
- [147] Stone, Marshall H. (1948): "The Generalized Weierstrass Approximation Theorem," *Mathematics Magazine*, 21, 167-184.
- [148] Theil, Henri. (1953): "Repeated least squares applied to complete equation systems," The Hague, Central Planning Bureau, mimeo.
- [149] Theil, Henri (1961): *Economic Forecasts and Policy*. Amsterdam: North Holland.
- [150] Theil, Henri. (1971): *Principles of Econometrics*, New York: Wiley.
- [151] Tobin, James (1958): "Estimation of relationships for limited dependent variables," *Econometrica*, 26, 24-36.
- [152] Tripathi, Gautam (1999): "A matrix extension of the Cauchy-Schwarz inequality," *Economics Letters*, 63, 1-3.
- [153] van der Vaart, A.W. (1998): *Asymptotic Statistics*, Cambridge University Press.
- [154] Wald, Abraham. (1940): "The fitting of straight lines if both variables are subject to error," *The Annals of Mathematical Statistics*, 11, 283-300

- [155] Wald, Abraham. (1943): “Tests of statistical hypotheses concerning several parameters when the number of observations is large,” *Transactions of the American Mathematical Society*, 54, 426-482.
- [156] Weierstrass, K. (1885): “Über die analytische Darstellbarkeit sogenannter willkürlicher Functionen einer reellen Veränderlichen,” *Sitzungsberichte der Königlich Preussischen Akademie der Wissenschaften zu Berlin*, 1885.
- [157] White, Halbert (1980): “A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity,” *Econometrica*, 48, 817-838.
- [158] White, Halbert (1984): *Asymptotic Theory for Econometricians*, Academic Press.
- [159] Wooldridge, Jeffrey M. (1995): “Score diagnostics for linear models estimated by two stage least squares,” In *Advances in Econometrics and Quantitative Economics: Essays in honor of Professor C. R. Rao*, eds. G. S. Maddala, P.C.B. Phillips, and T.N. Srinivasan, 66-87. Cambridge: Blackwell.
- [160] Wooldridge, Jeffrey M. (2010) *Econometric Analysis of Cross Section and Panel Data*, 2<sup>nd</sup> edition, MIT Press.
- [161] Wooldridge, Jeffrey M. (2015) *Introductory Econometrics: A Modern Approach*, 6<sup>th</sup> edition, Southwestern.
- [162] Wu, De-Min (1973): Alternative tests of independence between stochastic regressors and disturbances,” *Econometrica*, 41, 733-750.
- [163] Zellner, Arnold. (1962): “An efficient method of estimating seemingly unrelated regressions, and tests for aggregation bias,” *Journal of the American Statistical Association*, 57, 348-368.
- [164] Zhang, Fuzhen and Qingling Zhang (2006): “Eigenvalue inequalities for matrix product,” *IEEE Transactions on Automatic Control*, 51, 1506-1509.)