

**Stanislav Anatolyev**

**Intermediate and advanced econometrics: problems and solutions**

**Third edition**

KL/2009/018

Moscow

2009

**Анатольев С.А.** Задачи и решения по эконометрике. #KL/2009/018. М.: Российская экономическая школа, 2009 г. – 178 с. (Англ.)

Данное пособие – сборник задач, которые использовались автором при преподавании эконометрики промежуточного и продвинутого уровней в Российской Экономической Школе в течение последних нескольких лет. Все задачи сопровождаются решениями.

**Ключевые слова:** асимптотическая теория, бутстрап, линейная регрессия, метод наименьших квадратов, нелинейная регрессия, непараметрическая регрессия, экстремальное оценивание, метод наибольшего правдоподобия, инструментальные переменные, обобщенный метод моментов, эмпирическое правдоподобие, анализ панельных данных, условные ограничения на моменты, альтернативная асимптотика, асимптотика высокого порядка.

**Anatolyev, Stanislav A.** Intermediate and advanced econometrics: problems and solutions. #KL 2009/018 – Moscow, New Economic School, 2009 – 178 pp. (Eng.)

This manual is a collection of problems that the author has been using in teaching intermediate and advanced level econometrics courses at the New Economic School during last several years. All problems are accompanied by sample solutions.

**Key words:** asymptotic theory, bootstrap, linear regression, ordinary and generalized least squares, nonlinear regression, nonparametric regression, extremum estimation, maximum likelihood, instrumental variables, generalized method of moments, empirical likelihood, panel data analysis, conditional moment restrictions, alternative asymptotics, higher-order asymptotics

ISBN

© Анатольев С.А., 2009 г.

© Российская экономическая школа, 2009 г.

# CONTENTS

---

<b>I</b>	<b>Problems</b>	<b>15</b>
<b>1</b>	<b>Asymptotic theory: general and independent data</b>	<b>17</b>
1.1	Asymptotics of transformations . . . . .	17
1.2	Asymptotics of rotated logarithms . . . . .	17
1.3	Escaping probability mass . . . . .	17
1.4	Asymptotics of $t$ -ratios . . . . .	18
1.5	Creeping bug on simplex . . . . .	18
1.6	Asymptotics of sample variance . . . . .	19
1.7	Asymptotics of roots . . . . .	19
1.8	Second-order Delta-method . . . . .	19
1.9	Asymptotics with shrinking regressor . . . . .	19
1.10	Power trends . . . . .	20
<b>2</b>	<b>Asymptotic theory: time series</b>	<b>21</b>
2.1	Trended vs. differenced regression . . . . .	21
2.2	Long run variance for AR(1) . . . . .	21
2.3	Asymptotics of averages of AR(1) and MA(1) . . . . .	21
2.4	Asymptotics for impulse response functions . . . . .	22
<b>3</b>	<b>Bootstrap</b>	<b>23</b>
3.1	Brief and exhaustive . . . . .	23
3.2	Bootstrapping $t$ -ratio . . . . .	23
3.3	Bootstrap bias correction . . . . .	23
3.4	Bootstrap in linear model . . . . .	24
3.5	Bootstrap for impulse response functions . . . . .	24
<b>4</b>	<b>Regression and projection</b>	<b>25</b>
4.1	Regressing and projecting dice . . . . .	25
4.2	Mixture of normals . . . . .	25
4.3	Bernoulli regressor . . . . .	25
4.4	Best polynomial approximation . . . . .	26
4.5	Handling conditional expectations . . . . .	26
<b>5</b>	<b>Linear regression and OLS</b>	<b>27</b>
5.1	Fixed and random regressors . . . . .	27
5.2	Consistency of OLS under serially correlated errors . . . . .	27
5.3	Estimation of linear combination . . . . .	28
5.4	Incomplete regression . . . . .	28
5.5	Generated coefficient . . . . .	29
5.6	OLS in nonlinear model . . . . .	29
5.7	Long and short regressions . . . . .	29
5.8	Ridge regression . . . . .	29
5.9	Inconsistency under alternative . . . . .	30
5.10	Returns to schooling . . . . .	30

<b>6</b>	<b>Heteroskedasticity and GLS</b>	<b>31</b>
6.1	Conditional variance estimation . . . . .	31
6.2	Exponential heteroskedasticity . . . . .	31
6.3	OLS and GLS are identical . . . . .	31
6.4	OLS and GLS are equivalent . . . . .	32
6.5	Equicorrelated observations . . . . .	32
6.6	Unbiasedness of certain FGLS estimators . . . . .	32
<b>7</b>	<b>Variance estimation</b>	<b>33</b>
7.1	White estimator . . . . .	33
7.2	HAC estimation under homoskedasticity . . . . .	33
7.3	Expectations of White and Newey–West estimators in IID setting . . . . .	34
<b>8</b>	<b>Nonlinear regression</b>	<b>35</b>
8.1	Local and global identification . . . . .	35
8.2	Identification when regressor is nonrandom . . . . .	35
8.3	Cobb–Douglas production function . . . . .	35
8.4	Exponential regression . . . . .	35
8.5	Power regression . . . . .	36
8.6	Transition regression . . . . .	36
8.7	Nonlinear consumption function . . . . .	36
<b>9</b>	<b>Extremum estimators</b>	<b>37</b>
9.1	Regression on constant . . . . .	37
9.2	Quadratic regression . . . . .	37
9.3	Nonlinearity at left hand side . . . . .	38
9.4	Least fourth powers . . . . .	38
9.5	Asymmetric loss . . . . .	38
<b>10</b>	<b>Maximum likelihood estimation</b>	<b>39</b>
10.1	Normal distribution . . . . .	39
10.2	Pareto distribution . . . . .	39
10.3	Comparison of ML tests . . . . .	39
10.4	Invariance of ML tests to reparameterizations of null . . . . .	40
10.5	Misspecified maximum likelihood . . . . .	40
10.6	Individual effects . . . . .	41
10.7	Irregular confidence interval . . . . .	41
10.8	Trivial parameter space . . . . .	41
10.9	Nuisance parameter in density . . . . .	41
10.10	MLE versus OLS . . . . .	42
10.11	MLE versus GLS . . . . .	42
10.12	MLE in heteroskedastic time series regression . . . . .	42
10.13	Does the link matter? . . . . .	43
10.14	Maximum likelihood and binary variables . . . . .	43
10.15	Maximum likelihood and binary dependent variable . . . . .	43
10.16	Poisson regression . . . . .	44
10.17	Bootstrapping ML tests . . . . .	44

<b>11 Instrumental variables</b>	<b>45</b>
11.1 Invalid 2SLS . . . . .	45
11.2 Consumption function . . . . .	45
11.3 Optimal combination of instruments . . . . .	46
11.4 Trade and growth . . . . .	46
<b>12 Generalized method of moments</b>	<b>47</b>
12.1 Nonlinear simultaneous equations . . . . .	47
12.2 Improved GMM . . . . .	47
12.3 Minimum Distance estimation . . . . .	47
12.4 Formation of moment conditions . . . . .	48
12.5 What CMM estimates . . . . .	48
12.6 Trinity for GMM . . . . .	49
12.7 All about J . . . . .	49
12.8 Interest rates and future inflation . . . . .	49
12.9 Spot and forward exchange rates . . . . .	50
12.10 Returns from financial market . . . . .	50
12.11 Instrumental variables in ARMA models . . . . .	51
12.12 Hausman may not work . . . . .	51
12.13 Testing moment conditions . . . . .	52
12.14 Bootstrapping OLS . . . . .	52
12.15 Bootstrapping DD . . . . .	52
<b>13 Panel data</b>	<b>53</b>
13.1 Alternating individual effects . . . . .	53
13.2 Time invariant regressors . . . . .	53
13.3 Within and Between . . . . .	54
13.4 Panels and instruments . . . . .	54
13.5 Differencing transformations . . . . .	54
13.6 Nonlinear panel data model . . . . .	55
13.7 Durbin–Watson statistic and panel data . . . . .	55
13.8 Higher-order dynamic panel . . . . .	56
<b>14 Nonparametric estimation</b>	<b>57</b>
14.1 Nonparametric regression with discrete regressor . . . . .	57
14.2 Nonparametric density estimation . . . . .	57
14.3 Nadaraya–Watson density estimator . . . . .	57
14.4 First difference transformation and nonparametric regression . . . . .	58
14.5 Unbiasedness of kernel estimates . . . . .	58
14.6 Shape restriction . . . . .	58
14.7 Nonparametric hazard rate . . . . .	59
14.8 Nonparametrics and perfect fit . . . . .	59
14.9 Nonparametrics and extreme observations . . . . .	59
<b>15 Conditional moment restrictions</b>	<b>61</b>
15.1 Usefulness of skedastic function . . . . .	61
15.2 Symmetric regression error . . . . .	61
15.3 Optimal instrumentation of consumption function . . . . .	61
15.4 Optimal instrument in AR-ARCH model . . . . .	62
15.5 Optimal instrument in AR with nonlinear error . . . . .	62
15.6 Optimal IV estimation of a constant . . . . .	62

15.7	Negative binomial distribution and PML . . . . .	63
15.8	Nesting and PML . . . . .	63
15.9	Misspecification in variance . . . . .	63
15.10	Modified Poisson regression and PML estimators . . . . .	64
15.11	Optimal instrument and regression on constant . . . . .	64
<b>16</b>	<b>Empirical Likelihood</b>	<b>65</b>
16.1	Common mean . . . . .	65
16.2	Kullback–Leibler Information Criterion . . . . .	65
16.3	Empirical likelihood as IV estimation . . . . .	66
<b>17</b>	<b>Advanced asymptotic theory</b>	<b>67</b>
17.1	Maximum likelihood and asymptotic bias . . . . .	67
17.2	Empirical likelihood and asymptotic bias . . . . .	67
17.3	Asymptotically irrelevant instruments . . . . .	67
17.4	Weakly endogenous regressors . . . . .	68
17.5	Weakly invalid instruments . . . . .	68
<b>II</b>	<b>Solutions</b>	<b>69</b>
<b>1</b>	<b>Asymptotic theory: general and independent data</b>	<b>71</b>
1.1	Asymptotics of transformations . . . . .	71
1.2	Asymptotics of rotated logarithms . . . . .	71
1.3	Escaping probability mass . . . . .	72
1.4	Asymptotics of $t$ -ratios . . . . .	72
1.5	Creeping bug on simplex . . . . .	74
1.6	Asymptotics of sample variance . . . . .	74
1.7	Asymptotics of roots . . . . .	74
1.8	Second-order Delta-method . . . . .	75
1.9	Asymptotics with shrinking regressor . . . . .	76
1.10	Power trends . . . . .	77
<b>2</b>	<b>Asymptotic theory: time series</b>	<b>79</b>
2.1	Trended vs. differenced regression . . . . .	79
2.2	Long run variance for AR(1) . . . . .	80
2.3	Asymptotics of averages of AR(1) and MA(1) . . . . .	80
2.4	Asymptotics for impulse response functions . . . . .	81
<b>3</b>	<b>Bootstrap</b>	<b>83</b>
3.1	Brief and exhaustive . . . . .	83
3.2	Bootstrapping $t$ -ratio . . . . .	83
3.3	Bootstrap bias correction . . . . .	83
3.4	Bootstrap in linear model . . . . .	84
3.5	Bootstrap for impulse response functions . . . . .	85
<b>4</b>	<b>Regression and projection</b>	<b>87</b>
4.1	Regressing and projecting dice . . . . .	87
4.2	Mixture of normals . . . . .	88
4.3	Bernoulli regressor . . . . .	89
4.4	Best polynomial approximation . . . . .	89
4.5	Handling conditional expectations . . . . .	90

<b>5</b>	<b>Linear regression and OLS</b>	<b>91</b>
5.1	Fixed and random regressors . . . . .	91
5.2	Consistency of OLS under serially correlated errors . . . . .	91
5.3	Estimation of linear combination . . . . .	92
5.4	Incomplete regression . . . . .	92
5.5	Generated coefficient . . . . .	93
5.6	OLS in nonlinear model . . . . .	94
5.7	Long and short regressions . . . . .	94
5.8	Ridge regression . . . . .	94
5.9	Inconsistency under alternative . . . . .	95
5.10	Returns to schooling . . . . .	96
<b>6</b>	<b>Heteroskedasticity and GLS</b>	<b>97</b>
6.1	Conditional variance estimation . . . . .	97
6.2	Exponential heteroskedasticity . . . . .	97
6.3	OLS and GLS are identical . . . . .	97
6.4	OLS and GLS are equivalent . . . . .	98
6.5	Equicorrelated observations . . . . .	99
6.6	Unbiasedness of certain FGLS estimators . . . . .	99
<b>7</b>	<b>Variance estimation</b>	<b>101</b>
7.1	White estimator . . . . .	101
7.2	HAC estimation under homoskedasticity . . . . .	101
7.3	Expectations of White and Newey–West estimators in IID setting . . . . .	102
<b>8</b>	<b>Nonlinear regression</b>	<b>103</b>
8.1	Local and global identification . . . . .	103
8.2	Identification when regressor is nonrandom . . . . .	103
8.3	Cobb–Douglas production function . . . . .	103
8.4	Exponential regression . . . . .	104
8.5	Power regression . . . . .	104
8.6	Simple transition regression . . . . .	104
8.7	Nonlinear consumption function . . . . .	105
<b>9</b>	<b>Extremum estimators</b>	<b>107</b>
9.1	Regression on constant . . . . .	107
9.2	Quadratic regression . . . . .	108
9.3	Nonlinearity at left hand side . . . . .	109
9.4	Least fourth powers . . . . .	110
9.5	Asymmetric loss . . . . .	111
<b>10</b>	<b>Maximum likelihood estimation</b>	<b>113</b>
10.1	Normal distribution . . . . .	113
10.2	Pareto distribution . . . . .	113
10.3	Comparison of ML tests . . . . .	114
10.4	Invariance of ML tests to reparameterizations of null . . . . .	115
10.5	Misspecified maximum likelihood . . . . .	116
10.6	Individual effects . . . . .	117
10.7	Irregular confidence interval . . . . .	117
10.8	Trivial parameter space . . . . .	118
10.9	Nuisance parameter in density . . . . .	118

10.10	MLE versus OLS . . . . .	119
10.11	MLE versus GLS . . . . .	120
10.12	MLE in heteroskedastic time series regression . . . . .	121
10.13	Does the link matter? . . . . .	123
10.14	Maximum likelihood and binary variables . . . . .	123
10.15	Maximum likelihood and binary dependent variable . . . . .	124
10.16	Poisson regression . . . . .	125
10.17	Bootstrapping ML tests . . . . .	127
<b>11</b>	<b>Instrumental variables</b>	<b>129</b>
11.1	Invalid 2SLS . . . . .	129
11.2	Consumption function . . . . .	130
11.3	Optimal combination of instruments . . . . .	131
11.4	Trade and growth . . . . .	132
<b>12</b>	<b>Generalized method of moments</b>	<b>133</b>
12.1	Nonlinear simultaneous equations . . . . .	133
12.2	Improved GMM . . . . .	133
12.3	Minimum Distance estimation . . . . .	134
12.4	Formation of moment conditions . . . . .	135
12.5	What CMM estimates . . . . .	136
12.6	Trinity for GMM . . . . .	136
12.7	All about J . . . . .	137
12.8	Interest rates and future inflation . . . . .	137
12.9	Spot and forward exchange rates . . . . .	138
12.10	Returns from financial market . . . . .	139
12.11	Instrumental variables in ARMA models . . . . .	140
12.12	Hausman may not work . . . . .	141
12.13	Testing moment conditions . . . . .	141
12.14	Bootstrapping OLS . . . . .	142
12.15	Bootstrapping DD . . . . .	142
<b>13</b>	<b>Panel data</b>	<b>143</b>
13.1	Alternating individual effects . . . . .	143
13.2	Time invariant regressors . . . . .	145
13.3	Within and Between . . . . .	146
13.4	Panels and instruments . . . . .	146
13.5	Differencing transformations . . . . .	147
13.6	Nonlinear panel data model . . . . .	148
13.7	Durbin–Watson statistic and panel data . . . . .	148
13.8	Higher-order dynamic panel . . . . .	149
<b>14</b>	<b>Nonparametric estimation</b>	<b>151</b>
14.1	Nonparametric regression with discrete regressor . . . . .	151
14.2	Nonparametric density estimation . . . . .	151
14.3	Nadaraya–Watson density estimator . . . . .	152
14.4	First difference transformation and nonparametric regression . . . . .	153
14.5	Unbiasedness of kernel estimates . . . . .	155
14.6	Shape restriction . . . . .	155
14.7	Nonparametric hazard rate . . . . .	156
14.8	Nonparametrics and perfect fit . . . . .	157



14.9	Nonparametrics and extreme observations . . . . .	158
<b>15</b>	<b>Conditional moment restrictions</b>	<b>159</b>
15.1	Usefulness of skedastic function . . . . .	159
15.2	Symmetric regression error . . . . .	160
15.3	Optimal instrumentation of consumption function . . . . .	161
15.4	Optimal instrument in AR-ARCH model . . . . .	162
15.5	Optimal instrument in AR with nonlinear error . . . . .	162
15.6	Optimal IV estimation of a constant . . . . .	163
15.7	Negative binomial distribution and PML . . . . .	164
15.8	Nesting and PML . . . . .	164
15.9	Misspecification in variance . . . . .	164
15.10	Modified Poisson regression and PML estimators . . . . .	165
15.11	Optimal instrument and regression on constant . . . . .	167
<b>16</b>	<b>Empirical Likelihood</b>	<b>169</b>
16.1	Common mean . . . . .	169
16.2	Kullback–Leibler Information Criterion . . . . .	171
16.3	Empirical likelihood as IV estimation . . . . .	172
<b>17</b>	<b>Advanced asymptotic theory</b>	<b>173</b>
17.1	Maximum likelihood and asymptotic bias . . . . .	173
17.2	Empirical likelihood and asymptotic bias . . . . .	173
17.3	Asymptotically irrelevant instruments . . . . .	174
17.4	Weakly endogenous regressors . . . . .	175
17.5	Weakly invalid instruments . . . . .	176



# PREFACE

---

This manual is a third edition of the collection of problems that I have been using in teaching intermediate and advanced level econometrics courses at the New Economic School (NES), Moscow, for already a decade. All problems are accompanied by sample solutions. Approximately, chapters 1–8 and 14 of the collection belong to a course in intermediate level econometrics (“Econometrics III” in the NES internal course structure); chapters 9–13 – to a course in advanced level econometrics (“Econometrics IV”, respectively). The problems in chapters 15–17 require knowledge of more advanced and special material. They have been used in the NES course “Topics in Econometrics”.

Many of the problems are not new. Some are inspired by my former teachers of econometrics at PhD studies: Hyungtaik Ahn, Mahmoud El-Gamal, Bruce Hansen, Yuichi Kitamura, Charles Manski, Gautam Tripathi, Kenneth West. Some problems are borrowed from their problem sets, as well as problem sets of other leading econometrics scholars or their textbooks. Some originate from the *Problems and Solutions* section of the journal *Econometric Theory*, where the author has published several problems.

The release of this collection would be hard without valuable help of my teaching assistants during various years: Andrey Vasnev, Viktor Subbotin, Semyon Polbennikov, Alexander Vaschilko, Denis Sokolov, Oleg Itskhoki, Andrey Shabalin, Stanislav Kolenikov, Anna Mikusheva, Dmitry Shakin, Oleg Shibanov, Vadim Cherepanov, Pavel Stetsenko, Ivan Lazarev, Yulia Shkurat, Dmitry Muravyev, Artem Shamgunov, Danila Deliya, Viktoria Stepanova, Boris Gershman, Alexander Migita, Ivan Mirgorodsky, Roman Chkoller, Andrey Savochkin, Alexander Knobel, Ekaterina Lavrinenko, Yulia Vakhrutdinova, Elena Pikulina, to whom go my deepest thanks. My thanks also go to my students and assistants who spotted errors and typos that crept into the first and second editions of this manual, especially Dmitry Shakin, Denis Sokolov, Pavel Stetsenko, Georgy Kartashov, and Roman Chkoller. Preparation of this manual was supported in part by the Swedish Professorship (2000–2003) from the Economics Education and Research Consortium, with funds provided by the Government of Sweden through the Eurasia Foundation, and by the Access Industries Professorship (2003–2009) from Access Industries.

I will be grateful to everyone who finds errors, mistakes and typos in this collection and reports them to [sanatoly@nes.ru](mailto:sanatoly@nes.ru).



# NOTATION AND ABBREVIATIONS

---

- ID – identification  
FOC/SOC – first/second order condition(s)  
CDF – cumulative distribution function, typically denoted as  $F$   
PDF – probability density function, typically denoted as  $f$   
LIME – law of iterated (mathematical) expectations  
LLN – law of large numbers  
CLT – central limit theorem  
 $\mathbb{I}_{\{A\}}$  – indicator function equalling unity when  $A$  holds and zero otherwise  
 $\Pr\{A\}$  – probability of  $A$   
 $\mathbb{E}[y|x]$  – mathematical expectation (mean) of  $y$  conditional on  $x$   
 $\mathbb{V}[y|x]$  – variance of  $y$  conditional on  $x$   
 $\mathbb{C}[x, y]$  – covariance between  $x$  and  $y$   
BLP,  $\mathbb{BLP}[y|x]$  – best linear predictor  
 $I_k$  –  $k \times k$  identity matrix  
plim – probability limit  
 $\sim$  typically means “distributed as”  
 $\mathcal{N}$  – normal (Gaussian) distribution  
 $\chi_k^2$  – chi-squared distribution with  $k$  degrees of freedom  
 $\chi_k^2(\delta)$  – non-central chi-squared distribution with  $k$  degrees of freedom and non-centrality parameter  $\delta$   
 $\mathcal{B}(p)$  – Bernoulli distribution with success probability  $p$   
IID – independently and identically distributed  
 $n$  – typically sample size in cross-sections  
 $T$  – typically sample size in time series  
 $k$  – typically number of parameters in parametric models  
 $\ell$  – typically number of instruments or moment conditions  
 $\mathcal{X}, \mathcal{Y}, \mathcal{Z}, \mathcal{E}, \widehat{\mathcal{E}}$  – data matrices of regressors, dependent variables, instruments, errors, residuals  
 $\mathcal{L}(\circ)$  – (conditional) likelihood function  
 $\ell_n(\circ)$  – (conditional) loglikelihood function  
 $s(\circ)$  – (conditional) score function  
 $m(\circ)$  – moment function  
 $Q_f$  – typically  $\mathbb{E}[f]$ , for example,  $Q_{xx} = \mathbb{E}[xx']$ ,  $Q_{gge^2} = \mathbb{E}[g\theta g' e^2]$ ,  $Q_{\partial m} = \mathbb{E}[\partial m / \partial \theta]$ , etc.  
 $\mathcal{I}$  – Information matrix  
 $\mathcal{W}$  – Wald test statistic  
 $\mathcal{LR}$  – likelihood ratio test statistic  
 $\mathcal{LM}$  – Lagrange multiplier (score) test statistic  
 $\mathcal{J}$  – Hansen’s J test statistic



**Part I**  
**Problems**





# 1. ASYMPTOTIC THEORY: GENERAL AND INDEPENDENT DATA

---

## 1.1 Asymptotics of transformations

1. Suppose that  $\sqrt{T}(\hat{\phi} - 2\pi) \xrightarrow{d} \mathcal{N}(0, 1)$ . Find the limiting distribution of  $T(1 - \cos \hat{\phi})$ .
2. Suppose that  $T(\hat{\psi} - 2\pi) \xrightarrow{d} \mathcal{N}(0, 1)$ . Find the limiting distribution of  $T \sin \hat{\psi}$ .
3. Suppose that  $T\hat{\theta} \xrightarrow{d} \chi_1^2$ . Find the limiting distribution of  $T \log \hat{\theta}$ .

## 1.2 Asymptotics of rotated logarithms

Let the positive random vector  $(U_n, V_n)'$  be such that

$$\sqrt{n} \left( \begin{pmatrix} U_n \\ V_n \end{pmatrix} - \begin{pmatrix} \mu_u \\ \mu_v \end{pmatrix} \right) \xrightarrow{d} \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \omega_{uu} & \omega_{uv} \\ \omega_{uv} & \omega_{vv} \end{pmatrix} \right)$$

as  $n \rightarrow \infty$ . Find the joint asymptotic distribution of

$$\begin{pmatrix} \ln U_n - \ln V_n \\ \ln U_n + \ln V_n \end{pmatrix}.$$

What is the condition under which  $\ln U_n - \ln V_n$  and  $\ln U_n + \ln V_n$  are asymptotically independent?

## 1.3 Escaping probability mass

Let  $X = \{x_1, \dots, x_n\}$  be a random sample from some population of  $x$  with  $\mathbb{E}[x] = \mu$  and  $\mathbb{V}[x] = \sigma^2$ . Let  $A_n$  denote an event such that  $\mathbb{P}\{A_n\} = 1 - \frac{1}{n}$ , and let the distribution of  $A_n$  be independent of the distribution of  $x$ . Now construct the following randomized estimator of  $\mu$ :

$$\hat{\mu}_n = \begin{cases} \bar{x}_n & \text{if } A_n \text{ happens,} \\ n & \text{otherwise.} \end{cases}$$

- (i) Find the bias, variance, and  $\text{MSE}$  of  $\hat{\mu}_n$ . Show how they behave as  $n \rightarrow \infty$ .
- (ii) Is  $\hat{\mu}_n$  a consistent estimator of  $\mu$ ? Find the asymptotic distribution of  $\sqrt{n}(\hat{\mu}_n - \mu)$ .
- (iii) Use this distribution to construct an approximately  $(1 - \alpha) \times 100\%$  confidence interval for  $\mu$ . Compare this CI with the one obtained by using  $\bar{x}_n$  as an estimator of  $\mu$ .

## 1.4 Asymptotics of $t$ -ratios

Let  $\{x_i\}_{i=1}^n$  be a random sample of a scalar random variable  $x$  with  $\mathbb{E}[x] = \mu$ ,  $\mathbb{V}[x] = \sigma^2$ ,  $\mathbb{E}[(x - \mu)^3] = 0$ ,  $\mathbb{E}[(x - \mu)^4] = \tau$ , where all parameters are finite.

- (a) Define  $T_n \equiv \frac{\bar{x}}{\hat{\sigma}}$ , where

$$\bar{x} \equiv \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\sigma}^2 \equiv \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2.$$

Derive the limiting distribution of  $\sqrt{n}T_n$  under the assumption  $\mu = 0$ .

- (b) Now suppose it is not assumed that  $\mu = 0$ . Derive the limiting distribution of

$$\sqrt{n} \left( T_n - \operatorname{plim}_{n \rightarrow \infty} T_n \right).$$

Be sure your answer reduces to the result of part (a) when  $\mu = 0$ .

- (c) Define  $R_n \equiv \frac{\bar{x}}{\bar{\sigma}}$ , where

$$\bar{\sigma}^2 \equiv \frac{1}{n} \sum_{i=1}^n x_i^2$$

is the constrained estimator of  $\sigma^2$  under the (possibly incorrect) assumption  $\mu = 0$ . Derive the limiting distribution of

$$\sqrt{n} \left( R_n - \operatorname{plim}_{n \rightarrow \infty} R_n \right)$$

for arbitrary  $\mu$  and  $\sigma^2 > 0$ . Under what conditions on  $\mu$  and  $\sigma^2$  will this asymptotic distribution be the same as in part (b)?

## 1.5 Creeping bug on simplex

Consider a positive  $(x, y)$  orthant  $\mathbb{R}_+^2$  and its unit simplex, i.e. the line segment  $x + y = 1$ ,  $x \geq 0$ ,  $y \geq 0$ . Take an arbitrary natural number  $k \in \mathbb{N}$ . Imagine a bug starting creeping from the origin  $(x, y) = (0, 0)$ . Each second the bug goes either in the positive  $x$  direction with probability  $p$ , or in the positive  $y$  direction with probability  $1 - p$ , each time covering distance  $\frac{1}{k}$ . Evidently, this way the bug reaches the simplex in  $k$  seconds. Suppose it arrives there at point  $(x_k, y_k)$ . Now let  $k \rightarrow \infty$ , i.e. as if the bug shrinks in size and physical abilities per second. Determine

- the probability limit of  $(x_k, y_k)$ ;
- the rate of convergence;
- the asymptotic distribution of  $(x_k, y_k)$ .

## 1.6 Asymptotics of sample variance

Let  $x_1, \dots, x_n$  be a random sample from a population of  $x$  with finite fourth moments. Let  $\bar{x}_n$  and  $\overline{x_n^2}$  be the sample averages of  $x$  and  $x^2$ , respectively. Find constants  $a$  and  $b$  and function  $c(n)$  such that the vector sequence

$$c(n) \begin{pmatrix} \bar{x}_n - a \\ \overline{x_n^2} - b \end{pmatrix}$$

converges to a nontrivial distribution, and determine this limiting distribution. Derive the asymptotic distribution of the sample variance  $\overline{x_n^2} - (\bar{x}_n)^2$ .

## 1.7 Asymptotics of roots

Suppose we are interested in the inference about the root  $\theta$  of the nonlinear system

$$F(a, \theta) = 0,$$

where  $F : \mathbb{R}^p \times \mathbb{R}^k \rightarrow \mathbb{R}^k$ , and  $a$  is a vector of constants. Let available be  $\hat{a}$ , a consistent and asymptotically normal estimator of  $a$ . Assuming that  $\theta$  is the unique solution of the above system, and  $\hat{\theta}$  is the unique solution of the system

$$F(\hat{a}, \hat{\theta}) = 0,$$

derive the asymptotic distribution of  $\hat{\theta}$ . Assume that all needed smoothness conditions are satisfied.

## 1.8 Second-order Delta-method

Let  $S_n = n^{-1} \sum_{i=1}^n X_i$ , where  $X_i, i = 1, \dots, n$ , is a random sample of scalar random variables with  $\mathbb{E}[X_i] = \mu$  and  $\mathbb{V}[X_i] = 1$ . It is easy to show that  $\sqrt{n}(S_n - \mu) \xrightarrow{d} \mathcal{N}(0, 1)$  when  $\mu \neq 0$ .

- Find the asymptotic distribution of  $S_n^2$  when  $\mu = 0$ , by taking a square of the asymptotic distribution of  $S_n$ .
- Find the asymptotic distribution of  $\cos(S_n)$ . Hint: take a higher-order Taylor expansion applied to  $\cos(S_n)$ .
- Using the technique of part (b), formulate and prove an analog of the Delta-method for the case when the function is scalar-valued, has zero first derivative and nonzero second derivative (when the derivatives are evaluated at the probability limit). For simplicity, let all involved random variables be scalars.

## 1.9 Asymptotics with shrinking regressor

Suppose that

$$y_i = \alpha + \beta x_i + u_i,$$

where  $\{u_i\}$  are IID with  $\mathbb{E}[u_i] = 0$ ,  $\mathbb{E}[u_i^2] = \sigma^2$  and  $\mathbb{E}[u_i^3] = \nu$ , while the regressor  $x_i$  is deterministically shrinking:  $x_i = \rho^i$  with  $\rho \in (0, 1)$ . Let the sample size be  $n$ . Discuss as fully as you can the asymptotic behavior of the OLS estimates  $(\hat{\alpha}, \hat{\beta}, \hat{\sigma}^2)$  of  $(\alpha, \beta, \sigma^2)$  as  $n \rightarrow \infty$ .

## 1.10 Power trends

Suppose that

$$y_i = \beta x_i + \sigma_i \varepsilon_i, \quad i = 1, \dots, n,$$

where  $\varepsilon_i \sim IID(0, 1)$  while  $x_i = i^\lambda$  for some known  $\lambda$ , and  $\sigma_i^2 = \delta i^\mu$  for some known  $\mu$ .

1. Under what conditions on  $\lambda$  and  $\mu$  is the OLS estimator of  $\beta$  consistent? Derive its asymptotic distribution when it is consistent.
2. Under what conditions on  $\lambda$  and  $\mu$  is the GLS estimator of  $\beta$  consistent? Derive its asymptotic distribution when it is consistent.

## 2. ASYMPTOTIC THEORY: TIME SERIES

### 2.1 Trended vs. differenced regression

Consider a linear model with a linearly trending regressor:

$$y_t = \alpha + \beta t + \varepsilon_t,$$

where the sequence  $\varepsilon_t$  is independently and identically distributed according to some distribution  $\mathcal{D}$  with mean zero and variance  $\sigma^2$ . The object of interest is  $\beta$ .

1. Write out the OLS estimator  $\hat{\beta}$  of  $\beta$  in deviations form and find its asymptotic distribution.
2. A researcher suggests removing the trending regressor by taking differences to obtain

$$y_t - y_{t-1} = \beta + \varepsilon_t - \varepsilon_{t-1}$$

and then estimating  $\beta$  by OLS. Write out the OLS estimator  $\check{\beta}$  of  $\beta$  and find its asymptotic distribution.

3. Compare the estimators  $\hat{\beta}$  and  $\check{\beta}$  in terms of asymptotic efficiency.

### 2.2 Long run variance for AR(1)

Often one needs to estimate the long-run variance

$$V_{ze} \equiv \lim_{T \rightarrow \infty} \mathbb{V} \left( \frac{1}{\sqrt{T}} \sum_{t=1}^T z_t e_t \right)$$

of a stationary sequence  $z_t e_t$  that satisfies the restriction  $\mathbb{E}[e_t | z_t] = 0$ . Derive a compact expression for  $V_{ze}$  in the case when  $e_t$  and  $z_t$  follow independent scalar  $AR(1)$  processes. For this example, propose a method to consistently estimate  $V_{ze}$ , and show your estimator's consistency.

### 2.3 Asymptotics of averages of AR(1) and MA(1)

Let  $x_t$  be a martingale difference sequence with respect to its own past, and let all conditions for the CLT be satisfied:  $\sqrt{T} \bar{x}_T = T^{-1/2} \sum_{t=1}^T x_t \xrightarrow{d} \mathcal{N}(0, \sigma^2)$ . Let now  $y_t = \rho y_{t-1} + x_t$  and  $z_t = x_t + \theta x_{t-1}$ , where  $|\rho| < 1$  and  $|\theta| < 1$ . Consider time averages  $\bar{y}_T = T^{-1} \sum_{t=1}^T y_t$  and  $\bar{z}_T = T^{-1} \sum_{t=1}^T z_t$ .

1. Are  $y_t$  and  $z_t$  martingale difference sequences relative to their own past?
2. Find the asymptotic distributions of  $\bar{y}_T$  and  $\bar{z}_T$ .

3. How would you estimate the asymptotic variances of  $\bar{y}_T$  and  $\bar{z}_T$ ?
4. Repeat what you did in parts 1–3 when  $\mathbf{x}_t$  is a  $k \times 1$  vector, and we have  $\sqrt{T}\bar{\mathbf{x}}_T = T^{-1/2} \sum_{t=1}^T \mathbf{x}_t \xrightarrow{d} \mathcal{N}(\mathbf{0}, \Sigma)$ ,  $\mathbf{y}_t = \mathbf{P}\mathbf{y}_{t-1} + \mathbf{x}_t$ ,  $\mathbf{z}_t = \mathbf{x}_t + \Theta\mathbf{x}_{t-1}$ , where  $\mathbf{P}$  and  $\Theta$  are  $k \times k$  matrices with eigenvalues inside the unit circle.

## 2.4 Asymptotics for impulse response functions

A stationary and ergodic process  $z_t$  that admits the representation

$$z_t = \mu + \sum_{j=0}^{\infty} \phi_j \varepsilon_{t-j},$$

where  $\sum_{j=0}^{\infty} |\phi_j| < \infty$  and  $\varepsilon_t$  is zero mean IID, is called *linear*. The function  $IRF(j) = \phi_j$  is called *impulse response function* of  $z_t$ , reflecting the fact that  $\phi_j = \partial z_t / \partial \varepsilon_{t-j}$ , a response of  $z_t$  to its unit shock  $j$  periods ago.

1. Show that the strong zero mean AR(1) and ARMA(1,1) processes

$$y_t = \rho y_{t-1} + \varepsilon_t, \quad |\rho| < 1$$

and

$$z_t = \rho z_{t-1} + \varepsilon_t - \theta \varepsilon_{t-1}, \quad |\rho| < 1, |\theta| < 1, \theta \neq \rho,$$

are linear, and derive their impulse response functions.

2. Suppose a sample  $z_1, \dots, z_T$  is given. For the AR(1) process, construct an estimator of the IRF on the basis of the OLS estimator of  $\rho$ . Derive the asymptotic distribution of your IRF estimator for fixed horizon  $j$  as the sample size  $T \rightarrow \infty$ .
3. Suppose that for the ARMA(1,1) process one estimates  $\rho$  from the sample  $z_1, \dots, z_T$  by

$$\hat{\rho} = \frac{\sum_{t=3}^T z_t z_{t-2}}{\sum_{t=3}^T z_{t-1} z_{t-2}},$$

and  $\hat{\theta}$  – by an appropriate root of the quadratic equation

$$-\frac{\hat{\theta}}{1 + \hat{\theta}^2} = \frac{\sum_{t=2}^T \hat{\varepsilon}_t \hat{\varepsilon}_{t-1}}{\sum_{t=2}^T \hat{\varepsilon}_t^2}, \quad \hat{\varepsilon}_t = z_t - \hat{\rho} z_{t-1}.$$

On the basis of these estimates, construct an estimator of the impulse response function you derived. Outline the steps (no need to show all math) which you would undertake in order to derive its asymptotic distribution for fixed  $j$  as  $T \rightarrow \infty$ .

## 3. BOOTSTRAP

---

### 3.1 Brief and exhaustive

Evaluate the following claims.

1. The only difference between Monte–Carlo and the bootstrap is possibility and impossibility, respectively, of sampling from the true population.
2. When one does bootstrap, there is no reason to raise the number of bootstrap repetition too high: there is a level when making it larger does not yield any improvement in precision.
3. The bootstrap estimator of the parameter of interest is preferable to the asymptotic one, since its rate of convergence to the true parameter is often larger.

### 3.2 Bootstrapping $t$ -ratio

Consider the following bootstrap procedure. Using the nonparametric bootstrap, generate bootstrap samples and calculate  $\frac{\hat{\theta}_b^* - \hat{\theta}}{s(\hat{\theta})}$  at each bootstrap repetition. Find the quantiles  $q_{\alpha/2}^*$  and  $q_{1-\alpha/2}^*$  from this bootstrap distribution, and construct

$$CI = [\hat{\theta} - s(\hat{\theta})q_{1-\alpha/2}^*, \hat{\theta} + s(\hat{\theta})q_{\alpha/2}^*].$$

Show that  $CI$  is exactly the same as the percentile interval, and *not* the percentile- $t$  interval.

### 3.3 Bootstrap bias correction

1. Consider a random variable  $x$  with mean  $\mu$ . A random sample  $\{x_i\}_{i=1}^n$  is available. One estimates  $\mu$  by  $\bar{x}_n$  and  $\mu^2$  by  $\bar{x}_n^2$ . Find out what the bootstrap bias corrected estimators of  $\mu$  and  $\mu^2$  are.
2. Suppose we have a sample of two independent observations  $z_1 = 0$  and  $z_2 = 3$  from the same distribution. Let us be interested in  $\mathbb{E}[z^2]$  and  $(\mathbb{E}[z])^2$  which are natural to estimate by  $\bar{z}^2 = \frac{1}{2}(z_1^2 + z_2^2)$  and  $\bar{z}^2 = \frac{1}{4}(z_1 + z_2)^2$ . Compute the bootstrap-bias-corrected estimates of the quantities of interest.

## 3.4 Bootstrap in linear model

1. Suppose one has a random sample of  $n$  observations from the linear regression model

$$y = x'\beta + e, \quad \mathbb{E}[e|x] = 0.$$

Is the nonparametric bootstrap valid or invalid in the presence of heteroskedasticity? Explain.

2. Let the model be

$$y = x'\beta + e,$$

but  $\mathbb{E}[ex] \neq 0$ , i.e. the regressors are endogenous. The OLS estimator  $\hat{\beta}$  of the parameter  $\beta$  is biased. We know that the bootstrap is a good way to estimate bias, so the idea is to estimate the bias of  $\hat{\beta}$  and construct a bias-adjusted estimate of  $\beta$ . Explain whether or not the non-parametric bootstrap can be used to implement this idea.

3. Take the linear regression

$$y = x'\beta + e, \quad \mathbb{E}[e|x] = 0.$$

For a particular value of  $x$ , the object of interest is the conditional mean  $g(x) = \mathbb{E}[y|x]$ . Describe how you would use the percentile-t bootstrap to construct a confidence interval for  $g(x)$ .

## 3.5 Bootstrap for impulse response functions

Recall the formulation of Problem 2.4.

1. Describe in detail how to construct 95% error bands around the IRF estimates for the AR(1) process using the bootstrap that attains asymptotic refinement.
2. It is well known that in spite of their asymptotic unbiasedness, usual estimates of impulse response functions are significantly biased in samples typically encountered in practice. Propose a bootstrap algorithm to construct a bias corrected impulse response function for the above ARMA(1,1) process.



# 4. REGRESSION AND PROJECTION

---

## 4.1 Regressing and projecting dice

Let  $y$  be a random variable that denotes the number of dots obtained when a fair six sided die is rolled. Let

$$x = \begin{cases} y & \text{if } y \text{ is even,} \\ 0 & \text{otherwise.} \end{cases}$$

- (i) Find the joint distribution of  $(x, y)$ .
- (ii) Find the best predictor of  $y$  given  $x$ .
- (iii) Find the best linear predictor,  $\mathbb{BLP}[y|x]$ , of  $y$  conditional on  $x$ .
- (iv) Calculate  $\mathbb{E}[U_{BP}^2]$  and  $\mathbb{E}[U_{BLP}^2]$ , the mean square prediction errors for cases (ii) and (iii) respectively, and show that  $\mathbb{E}[U_{BP}^2] \leq \mathbb{E}[U_{BLP}^2]$ .

## 4.2 Mixture of normals

Suppose that  $n$  pairs  $(x_i, y_i)$ ,  $i = 1, \dots, n$ , are independently drawn from the following *mixture of normals* distribution:

$$\begin{pmatrix} x \\ y \end{pmatrix} \sim \begin{cases} \mathcal{N}\left(\begin{pmatrix} 0 \\ 4 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}\right) & \text{with probability } p, \\ \mathcal{N}\left(\begin{pmatrix} 4 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}\right) & \text{with probability } 1 - p, \end{cases}$$

where  $0 < p < 1$ .

1. Derive the best linear predictor  $\mathbb{BLP}[y|x]$  of  $y$  given  $x$ .
2. Argue that the conditional expectation function  $\mathbb{E}[y|x]$  is nonlinear. Provide a step-by-step algorithm allowing one to derive  $\mathbb{E}[y|x]$ , and derive it if you can.

## 4.3 Bernoulli regressor

Let  $x$  be distributed Bernoulli, and, conditional on  $x$ ,  $y$  be distributed as

$$y|x \sim \begin{cases} \mathcal{N}(\mu_0, \sigma_0^2), & x = 0, \\ \mathcal{N}(\mu_1, \sigma_1^2), & x = 1. \end{cases}$$

Write out  $\mathbb{E}[y|x]$  and  $\mathbb{E}[y^2|x]$  as linear functions of  $x$ . Why are these expectations linear in  $x$ ?

## 4.4 Best polynomial approximation

Given jointly distributed random variables  $x$  and  $y$ , a best  $k^{\text{th}}$  order polynomial approximation  $\mathbb{BPA}_k[y|x]$  to  $\mathbb{E}[y|x]$ , in the MSE sense, is a solution to the problem

$$\min_{\alpha_0, \alpha_1, \dots, \alpha_k} \mathbb{E} \left[ \left( \mathbb{E}[y|x] - \alpha_0 - \alpha_1 x - \dots - \alpha_k x^k \right)^2 \right].$$

Assuming that  $\mathbb{BPA}_k[y|x]$  exists, find its characterization and derive the properties of the associated prediction error  $U_k = y - \mathbb{BPA}_k[y|x]$ .

## 4.5 Handling conditional expectations

1. Consider the following situation. The vector  $(y, x, z, w)$  is a random quadruple. It is known that

$$\mathbb{E}[y|x, z, w] = \alpha + \beta x + \gamma z.$$

It is also known that  $\mathbb{C}[x, z] = 0$  and that  $\mathbb{C}[w, z] > 0$ . The parameters  $\alpha$ ,  $\beta$  and  $\gamma$  are not known. A random sample of observations on  $(y, x, w)$  is available;  $z$  is not observable. In this setting, a researcher weighs two options for estimating  $\beta$ . One is a linear least squares fit of  $y$  on  $x$ . The other is a linear least squares fit of  $y$  on  $(x, w)$ . Compare these options.

2. Let  $(x, y, z)$  be a random triple. For a given real constant  $\gamma$ , a researcher wants to estimate  $\mathbb{E}[y|\mathbb{E}[x|z] = \gamma]$ . The researcher knows that  $\mathbb{E}[x|z]$  and  $\mathbb{E}[y|z]$  are strictly increasing and continuous functions of  $z$ , and is given consistent estimates of these functions. Show how the researcher can use them to obtain a consistent estimate of the quantity of interest.

# 5. LINEAR REGRESSION AND OLS

---

## 5.1 Fixed and random regressors

1. Comment on: “Treating regressors  $x$  in a mean regression as random variables rather than fixed numbers simplifies further analysis, since then the observations  $(x_i, y_i)$  may be treated as IID across  $i$ ”.
2. A labor economist argues: “It is more plausible to think of my regressors as random rather than fixed. Look at *education*, for example. A person chooses her level of education, thus it is random. *Age* may be misreported, so it is random too. Even *gender* is random, because one can get a sex change operation done.” Comment on this pearl.
3. Consider a linear mean regression  $y = x'\beta + e$ ,  $\mathbb{E}[e|x] = 0$ , where  $x$ , instead of being IID across  $i$ , depends on  $i$  through an unknown function  $\varphi$  as  $x_i = \varphi(i) + u_i$ , where  $u_i$  are IID independent of  $e_i$ . Show that the OLS estimator of  $\beta$  is still unbiased.

## 5.2 Consistency of OLS under serially correlated errors

<sup>1</sup>Let  $\{y_t\}_{t=-\infty}^{+\infty}$  be a strictly stationary and ergodic stochastic process with zero mean and finite variance.

- (i) Define

$$\beta = \frac{\mathbb{C}[y_t, y_{t-1}]}{\mathbb{V}[y_t]}, \quad u_t = y_t - \beta y_{t-1},$$

so that we can write

$$y_t = \beta y_{t-1} + u_t.$$

Show that the error  $u_t$  satisfies  $\mathbb{E}[u_t] = 0$  and  $\mathbb{C}[u_t, y_{t-1}] = 0$ .

- (ii) Show that the OLS estimator  $\hat{\beta}$  from the regression of  $y_t$  on  $y_{t-1}$  is consistent for  $\beta$ .
- (iii) Show that, without further assumptions,  $u_t$  is serially correlated. Construct an example with serially correlated  $u_t$ .
- (iv) A 1994 paper in the *Journal of Econometrics* leads with the statement: “It is well known that in linear regression models with lagged dependent variables, ordinary least squares (OLS) estimators are inconsistent if the errors are autocorrelated”. This statement, or a slight variation of it, appears in virtually all econometrics textbooks. Reconcile this statement with your findings from parts (ii) and (iii).

---

<sup>1</sup>This problem closely follows J.M. Wooldridge (1998) Consistency of OLS in the Presence of Lagged Dependent Variable and Serially Correlated Errors. *Econometric Theory* 14, Problem 98.2.1.

## 5.3 Estimation of linear combination

Suppose one has a random sample of  $n$  observations from the linear regression model

$$y = \alpha + \beta x + \gamma z + e,$$

where  $e$  has mean zero and variance  $\sigma^2$  and is independent of  $(x, z)$ .

1. What is the conditional variance of the best linear conditionally (on the  $x$  and  $z$  samples) unbiased estimator  $\hat{\theta}$  of

$$\theta = \alpha + \beta c_x + \gamma c_z,$$

where  $c_x$  and  $c_z$  are some given constants?

2. Obtain the limiting distribution of

$$\sqrt{n}(\hat{\theta} - \theta).$$

Write your answer as a function of the means, variances and correlations of  $x$ ,  $z$  and  $e$  and of the constants  $\alpha, \beta, \gamma, c_x, c_z$ , assuming that all moments are finite.

3. For which value of the correlation coefficient between  $x$  and  $z$  is the asymptotic variance minimized for given variances of  $e$  and  $x$ ?
4. Discuss the relationship of the result of part 3 with the problem of multicollinearity.

## 5.4 Incomplete regression

Consider the linear regression

$$y = x'\beta + e, \quad \mathbb{E}[e|x] = 0, \quad \mathbb{E}[e^2|x] = \sigma^2,$$

where  $x$  is  $k_1 \times 1$ . Suppose that some component of the error  $e$  is observable, so that

$$e = z'\gamma + \eta,$$

where  $z_i$  is a  $k_2 \times 1$  vector of observables such that  $\mathbb{E}[\eta|z] = 0$  and  $\mathbb{E}[xz'] \neq 0$ . A researcher wants to estimate  $\beta$  and  $\gamma$  and considers two alternatives:

1. Run the regression of  $y$  on  $x$  and  $z$  to find the OLS estimates  $\hat{\beta}$  and  $\hat{\gamma}$  of  $\beta$  and  $\gamma$ .
2. Run the regression of  $y$  on  $x$  to get the OLS estimate  $\hat{\beta}$  of  $\beta$ , compute the OLS residuals  $\hat{e} = y - x'\hat{\beta}$  and run the regression of  $\hat{e}$  on  $z$  to retrieve the OLS estimate  $\hat{\gamma}$  of  $\gamma$ .

Which of the two methods would you recommend from the point of view of consistency of  $\hat{\beta}$  and  $\hat{\gamma}$ ? For the method(s) that yield(s) consistent estimates, find the limiting distribution of  $\sqrt{n}(\hat{\gamma} - \gamma)$ .

## 5.5 Generated coefficient

Consider the following regression model:

$$y = \beta x + \alpha z + u,$$

where  $\alpha$  and  $\beta$  are scalar unknown parameters,  $u$  has zero mean and unit variance, pair  $(x, z)$  are independent of  $u$  with  $\mathbb{E}[x^2] = \gamma_x^2 \neq 0$ ,  $\mathbb{E}[z^2] = \gamma_z^2 \neq 0$ ,  $\mathbb{E}[xz] = \gamma_{xz} \neq 0$ . A collection of triples  $\{(x_i, z_i, y_i)\}_{i=1}^n$  is a random sample. Suppose we are given an estimator  $\hat{\alpha}$  of  $\alpha$  independent of all  $u_i$ 's, and the limiting distribution of  $\sqrt{n}(\hat{\alpha} - \alpha)$  is  $\mathcal{N}(0, 1)$  as  $n \rightarrow \infty$ . Define the estimator  $\hat{\beta}$  of  $\beta$  as

$$\hat{\beta} = \left( \sum_{i=1}^n x_i^2 \right)^{-1} \sum_{i=1}^n x_i (y_i - \hat{\alpha} z_i).$$

Obtain the asymptotic distribution of  $\hat{\beta}$  as  $n \rightarrow \infty$ .

## 5.6 OLS in nonlinear model

Consider the equation  $y = (\alpha + \beta x)e$ , where  $y$  and  $x$  are scalar observables,  $e$  is unobservable. Let  $\mathbb{E}[e|x] = 1$  and  $\mathbb{V}[e|x] = 1$ . How would you estimate  $(\alpha, \beta)$  by OLS? How would you construct standard errors?

## 5.7 Long and short regressions

Take the true model  $y = x_1' \beta_1 + x_2' \beta_2 + e$ ,  $\mathbb{E}[e|x_1, x_2] = 0$ , and assume random sampling. Suppose that  $\beta_1$  is estimated by regressing  $y$  on  $x_1$  only. Find the probability limit of this estimator. What are the conditions when it is consistent for  $\beta_1$ ?

## 5.8 Ridge regression

In the standard linear mean regression model, one estimates  $k \times 1$  parameter  $\beta$  by

$$\tilde{\beta} = (\mathcal{X}'\mathcal{X} + \lambda I_k)^{-1} \mathcal{X}'\mathcal{Y},$$

where  $\lambda > 0$  is a fixed scalar,  $I_k$  is a  $k \times k$  identity matrix,  $\mathcal{X}$  is  $n \times k$  and  $\mathcal{Y}$  is  $n \times 1$  matrices of data.

1. Find  $\mathbb{E}[\tilde{\beta}|\mathcal{X}]$ . Is  $\tilde{\beta}$  conditionally unbiased? Is it unbiased?
2. Find the probability limit of  $\tilde{\beta}$  as  $n \rightarrow \infty$ . Is  $\tilde{\beta}$  consistent?
3. Find the asymptotic distribution of  $\tilde{\beta}$ .
4. From your viewpoint, why may one want to use  $\tilde{\beta}$  instead of the OLS estimator  $\hat{\beta}$ ? Give conditions under which  $\tilde{\beta}$  is preferable to  $\hat{\beta}$  according to your criterion, and vice versa.

## 5.9 Inconsistency under alternative

Suppose that

$$y = \alpha + \beta x + u,$$

where  $u$  is distributed  $\mathcal{N}(0, \sigma^2)$  independently of  $x$ . The variable  $x$  is unobserved. Instead we observe  $z = x + v$ , where  $v$  is distributed  $\mathcal{N}(0, \eta^2)$  independently of  $x$  and  $u$ . Given a sample of size  $n$ , it is proposed to run the linear regression of  $y$  on  $z$  and use a conventional  $t$ -test to test the null hypothesis  $\beta = 0$ . Critically evaluate this proposal.

## 5.10 Returns to schooling

A researcher presents his research on returns to schooling at a lunchtime seminar. He runs OLS, using a random sample of individuals, on a Mincer-type linear regression, where the left side variable is a logarithm of hourly wage rate. The results are shown in the table.

Regressor	Point estimate	Standard error	t-statistic
Constant	-0.30	2.16	-0.14
Male ( $g$ )	0.39	0.07	5.54
Age ( $a$ )	0.14	0.12	1.16
Experience ( $e$ )	0.019	0.036	0.52
Completed schooling ( $s$ )	0.027	0.023	1.15
Ability ( $f$ )	0.081	0.124	0.65
Schooling-ability interaction ( $sf$ )	-0.001	0.014	-0.06

1. The presenter says: “Our model yields a 2.7 percentage point return per additional year of schooling (I allow returns to schooling to vary by ability by introducing ability-schooling interaction, but the corresponding estimate is essentially zero). At the same time, the estimated coefficient on ability is 8.1 (although it is statistically insignificant). This implies that one would have to acquire three additional years of education to compensate for one standard deviation lower innate ability in terms of labor market returns.” A person from the audience argues: “So you have just divided one insignificant estimate by another insignificant estimate. This is like dividing zero by zero. You can get any answer by dividing zero by zero, so your number ‘3’ is as good as any other number.” How would you professionally respond to this argument?
2. Another person from the audience argues: “Your dummy variable ‘Male’ enters the regression only as a separate variable, so the gender influences only the intercept. But the corresponding estimate is statistically very significant (in fact, it is the only significant variable in your regression). This makes me think that it must enter the regression also in interactions with the other variables. If I were you, I would run two regressions, one for males and one for females, and test for differences in coefficients across the two using a sup-Wald test. In any case, I would compute bootstrap standard errors to replace your asymptotic standard errors hoping that most of parameters would become statistically significant with more precise standard errors.” How would you professionally respond to these arguments?

# 6. HETEROSKEDASTICITY AND GLS

## 6.1 Conditional variance estimation

Econometrician A claims: “In an IID context, to run OLS and GLS I don’t need to know the skedastic function. See, I can estimate the conditional variance matrix  $\Omega$  of the error vector by  $\hat{\Omega} = \text{diag}\{\hat{e}_i^2\}_{i=1}^n$ , where  $\hat{e}_i$  for  $i = 1, \dots, n$  are OLS residuals. When I run OLS, I can estimate the variance matrix by  $(\mathcal{X}'\mathcal{X})^{-1} \mathcal{X}'\hat{\Omega}\mathcal{X}(\mathcal{X}'\mathcal{X})^{-1}$ ; when I run feasible GLS, I use the formula  $\hat{\beta} = (\mathcal{X}'\hat{\Omega}^{-1}\mathcal{X})^{-1}\mathcal{X}'\hat{\Omega}^{-1}\mathcal{Y}$ .” Econometrician B argues: “That ain’t right. In both cases you are using only one observation,  $\hat{e}_i^2$ , to estimate the value of the skedastic function,  $\sigma^2(x_i)$ . Hence, your estimates will be inconsistent and inference wrong.” Resolve this dispute.

## 6.2 Exponential heteroskedasticity

Let  $y$  be scalar and  $x$  be  $k \times 1$  vector random variables. Observations  $(y_i, x_i)$  are drawn at random from the population of  $(y, x)$ . You are told that  $\mathbb{E}[y|x] = x'\beta$  and that  $\mathbb{V}[y|x] = \exp(x'\beta + \alpha)$ , with  $(\beta, \alpha)$  unknown. You are asked to estimate  $\beta$ .

1. Propose an estimation method that is asymptotically equivalent to GLS that would be computable were  $\mathbb{V}[y|x]$  fully known.
2. In what sense is the feasible GLS estimator of part 1 efficient? In which sense is it inefficient?

## 6.3 OLS and GLS are identical

Let  $\mathcal{Y} = \mathcal{X}(\beta + v) + \mathcal{U}$ , where  $\mathcal{X}$  is  $n \times k$ ,  $\mathcal{Y}$  and  $\mathcal{U}$  are  $n \times 1$ , and  $\beta$  and  $v$  are  $k \times 1$ . The parameter of interest is  $\beta$ . The properties of  $(\mathcal{Y}, \mathcal{X}, \mathcal{U}, v)$  are:  $\mathbb{E}[\mathcal{U}|\mathcal{X}] = 0$ ,  $\mathbb{E}[v|\mathcal{X}] = 0$ ,  $\mathbb{E}[\mathcal{U}\mathcal{U}'|\mathcal{X}] = \sigma^2 I_n$ ,  $\mathbb{E}[vv'|\mathcal{X}] = \Gamma$ ,  $\mathbb{E}[\mathcal{U}v'|\mathcal{X}] = 0$ .  $\mathcal{Y}$  and  $\mathcal{X}$  are observable, while  $\mathcal{U}$  and  $v$  are not.

1. What are  $\mathbb{E}[\mathcal{Y}|\mathcal{X}]$  and  $\mathbb{V}[\mathcal{Y}|\mathcal{X}]$ ? Denote the latter by  $\Sigma$ . Is the environment homo- or heteroskedastic?
2. Write out the OLS and GLS estimators  $\hat{\beta}$  and  $\tilde{\beta}$  of  $\beta$ . Prove that in this model they are identical. Hint: First prove that  $\mathcal{X}'\hat{\mathcal{E}} = 0$ , where  $\hat{\mathcal{E}}$  is the  $n \times 1$  vector of OLS residuals. Next prove that  $\mathcal{X}'\Sigma^{-1}\hat{\mathcal{E}} = 0$ . Then conclude. Alternatively, use formulae for the inverse of a sum of two matrices. The first method is preferable, being more “econometric”.
3. Discuss benefits of using both estimators in this model.

## 6.4 OLS and GLS are equivalent

Let us have a regression written in a matrix form:  $\mathcal{Y} = \mathcal{X}\beta + \mathcal{U}$ , where  $\mathcal{X}$  is  $n \times k$ ,  $\mathcal{Y}$  and  $\mathcal{U}$  are  $n \times 1$ , and  $\beta$  is  $k \times 1$ . The parameter of interest is  $\beta$ . The properties of  $u$  are:  $\mathbb{E}[\mathcal{U}|\mathcal{X}] = 0$ ,  $\mathbb{E}[\mathcal{U}\mathcal{U}'|\mathcal{X}] = \Sigma$ . Let it be also known that  $\Sigma\mathcal{X} = \mathcal{X}\Theta$  for some  $k \times k$  nonsingular matrix  $\Theta$ .

1. Prove that in this model the OLS and GLS estimators  $\hat{\beta}$  and  $\tilde{\beta}$  of  $\beta$  have the same finite sample conditional variance.
2. Apply this result to the following regression on a constant:

$$y_i = \alpha + u_i,$$

where the disturbances are equicorrelated, that is,  $\mathbb{E}[u_i] = 0$ ,  $\mathbb{V}[u_i] = \sigma^2$  and  $\mathbb{C}[u_i, u_j] = \rho\sigma^2$  for  $i \neq j$ .

## 6.5 Equicorrelated observations

Suppose  $x_i = \theta + u_i$ , where  $\mathbb{E}[u_i] = 0$  and

$$\mathbb{E}[u_i u_j] = \begin{cases} 1 & \text{if } i = j \\ \gamma & \text{if } i \neq j \end{cases}$$

with  $i, j = 1, \dots, n$ . Is  $\bar{x}_n = \frac{1}{n}(x_1 + \dots + x_n)$  the best linear unbiased estimator of  $\theta$ ? Investigate  $\bar{x}_n$  for consistency.

## 6.6 Unbiasedness of certain FGLS estimators

Show that

- (a) for a random variable  $z$ , if  $z$  and  $-z$  have the same distribution, then  $\mathbb{E}[z] = 0$ ;
- (b) for a random vector  $\varepsilon$  and a vector function  $q(\varepsilon)$  of  $\varepsilon$ , if  $\varepsilon$  and  $-\varepsilon$  have the same distribution and  $q(-\varepsilon) = -q(\varepsilon)$  for all  $\varepsilon$ , then  $\mathbb{E}[q(\varepsilon)] = 0$ .

Consider the linear regression model written in matrix form:

$$\mathcal{Y} = \mathcal{X}\beta + \mathcal{E}, \quad \mathbb{E}[\mathcal{E}|\mathcal{X}] = 0, \quad \mathbb{E}[\mathcal{E}\mathcal{E}'|\mathcal{X}] = \Sigma.$$

Let  $\hat{\Sigma}$  be an estimate of  $\Sigma$  which is a function of products of least squares residuals, i.e.  $\hat{\Sigma} = F(\mathcal{M}\mathcal{E}\mathcal{E}'\mathcal{M}) = H(\mathcal{E}\mathcal{E}')$  for  $\mathcal{M} = I - \mathcal{X}(\mathcal{X}'\mathcal{X})^{-1}\mathcal{X}'$ . Show that if  $\mathcal{E}$  and  $-\mathcal{E}$  have the same conditional distribution (e.g. if  $\mathcal{E}$  is conditionally normal), then the feasible GLS estimator

$$\tilde{\beta}_F = \left(\mathcal{X}'\hat{\Sigma}^{-1}\mathcal{X}\right)^{-1}\mathcal{X}'\hat{\Sigma}^{-1}\mathcal{Y}$$

is unbiased.



# 7. VARIANCE ESTIMATION

---

## 7.1 White estimator

Evaluate the following claims.

1. When one suspects heteroskedasticity, one should use the White formula

$$Q_{xx}^{-1} Q_{xxe^2} Q_{xx}^{-1}$$

instead of good old  $\sigma^2 Q_{xx}^{-1}$ , since under heteroskedasticity the latter does not make sense, because  $\sigma^2$  is different for each observation.

2. Since for the OLS estimator

$$\hat{\beta} = (\mathcal{X}'\mathcal{X})^{-1} \mathcal{X}'\mathcal{Y}$$

we have

$$\mathbb{E}[\hat{\beta}|\mathcal{X}] = \beta$$

and

$$\mathbb{V}[\hat{\beta}|\mathcal{X}] = (\mathcal{X}'\mathcal{X})^{-1} \mathcal{X}'\Omega\mathcal{X} (\mathcal{X}'\mathcal{X})^{-1},$$

we can estimate the finite sample variance by

$$\widehat{\mathbb{V}[\hat{\beta}|\mathcal{X}]} = (\mathcal{X}'\mathcal{X})^{-1} \sum_{i=1}^n x_i x_i' \hat{e}_i^2 (\mathcal{X}'\mathcal{X})^{-1}$$

(which, apart from the factor  $n$ , is the same as the White estimator of the asymptotic variance) and construct  $t$  and Wald statistics using it. Thus, we do not need asymptotic theory to do OLS estimation and inference.

## 7.2 HAC estimation under homoskedasticity

We look for a simplification of HAC variance estimators under conditional homoskedasticity. Suppose that the regressors  $x_t$  and left side variable  $y_t$  in a linear time series regression

$$y_t = x_t' \beta + e_t, \quad \mathbb{E}[e_t | x_t, x_{t-1}, \dots] = 0$$

are jointly stationary and ergodic. The error  $e_t$  is serially correlated of unknown order, but let it be known that it is conditionally homoskedastic, i.e.  $\mathbb{E}[e_t e_{t-j} | x_t, x_{t-1}, \dots] = \gamma_j$  is constant (i.e. does not depend on  $x_t, x_{t-1}, \dots$ ) for all  $j \geq 0$ . Develop a Newey–West-type HAC estimator of the long-run variance of  $x_t e_t$  that would take advantage of conditional homoskedasticity.

### 7.3 Expectations of White and Newey–West estimators in IID setting

Suppose one has a random sample of  $n$  observations from the linear conditionally homoskedastic regression model

$$y_i = x_i' \beta + e_i, \quad \mathbb{E}[e_i | x_i] = 0, \quad \mathbb{E}[e_i^2 | x_i] = \sigma^2.$$

Let  $\hat{\beta}$  be the OLS estimator of  $\beta$ , and let  $\hat{V}_{\hat{\beta}}$  and  $\check{V}_{\hat{\beta}}$  be the White and Newey–West estimators of the asymptotic variance matrix of  $\hat{\beta}$ . Find  $\mathbb{E}[\hat{V}_{\hat{\beta}} | \mathcal{X}]$  and  $\mathbb{E}[\check{V}_{\hat{\beta}} | \mathcal{X}]$ , where  $\mathcal{X}$  is the matrix of stacked regressors for all observations.

# 8. NONLINEAR REGRESSION

---

## 8.1 Local and global identification

Consider the nonlinear regression  $\mathbb{E}[y|x] = \beta_1 + \beta_2^2 x$ , where  $\beta_2 \neq 0$  and  $\mathbb{V}[x] \neq 0$ . Which identification condition for  $(\beta_1, \beta_2)'$  fails and which does not?

## 8.2 Identification when regressor is nonrandom

Suppose we regress  $y$  on scalar  $x$ , but  $x$  is distributed only at one point (that is,

$$\Pr\{x = a\} = 1$$

for some  $a$ ). When does the identification condition hold and when does it fail if the regression is linear and has no intercept? If the regression is nonlinear? Provide both algebraic and intuitive/graphical explanations.

## 8.3 Cobb–Douglas production function

Suppose we have a random sample of  $n$  firms with data on output  $Q$ , capital  $K$  and labor  $L$ , and want to estimate the Cobb–Douglas production function

$$Q = \alpha K^\theta L^{1-\theta} \varepsilon,$$

where  $\varepsilon$  has the property  $\mathbb{E}[\varepsilon|K, L] = 1$ . Evaluate the following suggestions of estimation of  $\theta$ :

1. Run a linear regression of  $\log Q - \log L$  on a constant and  $\log K - \log L$
2. For various values of  $\theta$  on a grid, run a linear regression of  $Q$  on  $K^\theta L^{1-\theta}$  without a constant, and select the value of  $\theta$  that minimizes a sum of squared OLS errors.

## 8.4 Exponential regression

Suppose you have the homoskedastic nonlinear regression

$$y = \exp(\alpha + \beta x) + e, \quad \mathbb{E}[e|x] = 0, \quad \mathbb{E}[e^2|x] = \sigma^2$$

and random sample  $\{(x_i, y_i)\}_{i=1}^n$ . Let the true  $\beta$  be 0, and  $x$  be distributed as standard normal. Investigate the problem for local identifiability, and derive the asymptotic distribution of the NLLS estimator of  $(\alpha, \beta)$ . Describe a concentration method algorithm giving all formulas (including standard errors that you would use in practice) in explicit forms.

## 8.5 Power regression

Suppose you have the nonlinear regression

$$y = \alpha(1 + x^\beta) + e, \quad \mathbb{E}[e|x] = 0$$

and IID data  $\{(x_i, y_i)\}_{i=1}^n$ . How would you test  $H_0 : \alpha = 0$  properly?

## 8.6 Transition regression

Given the random sample  $\{(x_i, y_i)\}_{i=1}^n$ , consider the nonlinear regression

$$y = \beta_1 + \frac{\beta_2}{1 + \beta_3 x} + e, \quad \mathbb{E}[e|x] = 0.$$

1. Describe how to test using the  $t$ -statistic if the marginal influence of  $x$  on the conditional mean of  $y$ , evaluated at  $x = 0$ , equals 1.
2. Describe how to test using the Wald statistic that the regression function does not depend on  $x$ .

## 8.7 Nonlinear consumption function

Consider the model

$$\mathbb{E}[c_t | y_{t-1}, y_{t-2}, y_{t-3}, \dots] = \alpha + \beta \mathbb{I}\{y_{t-1} > \Delta\} + \delta y_{t-1}^\gamma,$$

where  $c_t$  is consumption at  $t$  and  $y_t$  is income at  $t$ . The pair  $(c_t, y_t)$  is continuously distributed, stationary and ergodic. The parameter  $\Delta$  represents a “normal” income level, and is known. Suppose you are given a long quarterly series of length  $T$  on  $c_t$  and  $y_t$ .

1. Describe at least three different situations when parameter identification will fail.
2. Describe in detail how you will run the NLLS estimation employing the concentration method, including construction of standard errors for model parameters.
3. Describe how you will test the hypothesis  $H_0 : \gamma = 1$  against  $H_a : \gamma < 1$ 
  - (a) by employing the asymptotic approach,
  - (b) by employing the bootstrap approach without using standard errors.

# 9. EXTREMUM ESTIMATORS

## 9.1 Regression on constant

Consider the following model:

$$y = \beta + e,$$

where all variables are scalars. Assume that  $\{y_i\}_{i=1}^n$  is a random sample, and  $\mathbb{E}[e] = 0$ ,  $\mathbb{E}[e^2] = \beta^2$ ,  $\mathbb{E}[e^3] = 0$  and  $\mathbb{E}[e^4] = \kappa$ . Consider the following three estimators of  $\beta$ :

$$\hat{\beta}_1 = \frac{1}{n} \sum_{i=1}^n y_i,$$

$$\hat{\beta}_2 = \arg \min_b \left\{ \log b^2 + \frac{1}{nb^2} \sum_{i=1}^n (y_i - b)^2 \right\},$$

$$\hat{\beta}_3 = \frac{1}{2} \arg \min_b \sum_{i=1}^n \left( \frac{y_i}{b} - 1 \right)^2.$$

Derive the asymptotic distributions of these three estimators. Which of them would you prefer most on the asymptotic basis? What is the idea behind each of the three estimators?

## 9.2 Quadratic regression

Consider a nonlinear regression model

$$y = (\beta_0 + x)^2 + u,$$

where we assume:

- (A) The parameter space is  $\mathbb{B} = [-\frac{1}{2}, +\frac{1}{2}]$ .
- (B) The error  $u$  has properties  $\mathbb{E}[u] = 0$ ,  $\mathbb{V}[u] = \sigma_0^2$ .
- (C) The regressor  $x$  has is distributed uniformly over  $[1, 2]$  independently of  $u$ . In particular, this implies  $\mathbb{E}[x^{-1}] = \ln 2$  and  $\mathbb{E}[x^r] = \frac{1}{1+r}(2^{r+1} - 1)$  for integer  $r \neq -1$ .

A random sample  $\{(x_i, y_i)\}_{i=1}^n$  is available. Define two estimators of  $\beta_0$ :

1.  $\hat{\beta}$  minimizes  $\mathcal{S}_n(\beta) = \sum_{i=1}^n (y_i - (\beta + x_i)^2)^2$  over  $\mathbb{B}$ .
2.  $\tilde{\beta}$  minimizes  $\mathcal{W}_n(\beta) = \sum_{i=1}^n \left\{ \frac{y_i}{(\beta + x_i)^2} + \ln(\beta + x_i)^2 \right\}$  over  $\mathbb{B}$ .

For the case  $\beta_0 = 0$ , obtain asymptotic distributions of  $\hat{\beta}$  and  $\tilde{\beta}$ . Which one of the two do you prefer on the asymptotic basis?

### 9.3 Nonlinearity at left hand side

A random sample  $\{(x_i, y_i)\}_{i=1}^n$  is available for the nonlinear model

$$(y + \alpha)^2 = \beta x + e, \quad \mathbb{E}[e|x] = 0, \quad \mathbb{E}[e^2|x] = \sigma^2,$$

where the parameters  $\alpha$  and  $\beta$  are scalars.

1. Show that the NLLS estimator of  $\alpha$  and  $\beta$

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \arg \min_{a,b} \sum_{i=1}^n \left( (y_i + a)^2 - bx_i \right)^2$$

is in general inconsistent. What feature makes the model differ from a nonlinear regression where the NLLS estimator is consistent?

2. Propose a consistent CMM estimator of  $\alpha$  and  $\beta$  and derive its asymptotic distribution.

### 9.4 Least fourth powers

Suppose  $y = \beta x + e$ , where all variables are scalars,  $x$  and  $e$  are independent, and the distribution of  $e$  is symmetric around 0. For a random sample  $\{(x_i, y_i)\}_{i=1}^n$ , consider the following extremum estimator of  $\beta$ :

$$\hat{\beta} = \arg \min_b \sum_{i=1}^n (y_i - bx_i)^4.$$

Derive the asymptotic properties of  $\hat{\beta}$ , paying special attention to the identification condition. Compare this estimator with the OLS estimator in terms of asymptotic efficiency for the case when  $x$  and  $e$  are normally distributed.

### 9.5 Asymmetric loss

Suppose that  $\{(x_i, y_i)\}_{i=1}^n$  is a random sample from a population satisfying

$$y = \alpha + x'\beta + e,$$

where  $e$  is independent of  $x$ , a  $k \times 1$  vector. Suppose also that all moments of  $x$  and  $e$  are finite and that  $\mathbb{E}[xx']$  is nonsingular. Suppose that  $\hat{\alpha}$  and  $\hat{\beta}$  are defined to be the values of  $\alpha$  and  $\beta$  that minimize

$$\frac{1}{n} \sum_{i=1}^n \rho(y_i - \alpha - x_i'\beta)$$

over some set  $\Theta \subset \mathbb{R}^{k+1}$ , where for some  $0 < \gamma < 1$

$$\rho(u) = \begin{cases} \gamma u^3 & \text{if } u \geq 0, \\ -(1 - \gamma)u^3 & \text{if } u < 0. \end{cases}$$

Describe the asymptotic behavior of the estimators  $\hat{\alpha}$  and  $\hat{\beta}$  as  $n \rightarrow \infty$ . If you need to make additional assumptions be sure to specify what these are and why they are needed.

# 10. MAXIMUM LIKELIHOOD ESTIMATION

## 10.1 Normal distribution

Let  $x_1, \dots, x_n$  be a random sample from  $\mathcal{N}(\mu, \mu^2)$ . Derive the ML estimator  $\hat{\mu}$  of  $\mu$  and prove its consistency.

## 10.2 Pareto distribution

A random variable  $X$  is said to have a Pareto distribution with parameter  $\lambda$ , denoted  $X \sim \text{Pareto}(\lambda)$ , if it is continuously distributed with density

$$f_X(x|\lambda) = \begin{cases} \lambda x^{-(\lambda+1)}, & \text{if } x > 1, \\ 0, & \text{otherwise.} \end{cases}$$

A random sample  $x_1, \dots, x_n$  from the  $\text{Pareto}(\lambda)$  population is available.

- (a) Derive the ML estimator  $\hat{\lambda}$  of  $\lambda$ , prove its consistency and find its asymptotic distribution.
- (b) Derive the Wald, Likelihood Ratio and Lagrange Multiplier test statistics for testing the null hypothesis  $H_0 : \lambda = \lambda_0$  against the alternative hypothesis  $H_a : \lambda \neq \lambda_0$ . Do any of these statistics coincide?

## 10.3 Comparison of ML tests

<sup>1</sup>Berndt and Savin in 1977 showed that  $\mathcal{W} \geq \mathcal{LR} \geq \mathcal{LM}$  for the case of a multivariate regression model with normal disturbances. Ullah and Zinde-Walsh in 1984 showed that this inequality is not robust to non-normality of the disturbances. In the spirit of the latter article, this problem considers simple examples from non-normal distributions and illustrates how this conflict among criteria is affected.

1. Consider a random sample  $x_1, \dots, x_n$  from a Poisson distribution with parameter  $\lambda$ . Show that testing  $\lambda = 3$  versus  $\lambda \neq 3$  yields  $\mathcal{W} \geq \mathcal{LM}$  for  $\bar{x} \leq 3$  and  $\mathcal{W} \leq \mathcal{LM}$  for  $\bar{x} \geq 3$ .
2. Consider a random sample  $x_1, \dots, x_n$  from an exponential distribution with parameter  $\theta$ . Show that testing  $\theta = 3$  versus  $\theta \neq 3$  yields  $\mathcal{W} \geq \mathcal{LM}$  for  $0 < \bar{x} \leq 3$  and  $\mathcal{W} \leq \mathcal{LM}$  for  $\bar{x} \geq 3$ .
3. Consider a random sample  $x_1, \dots, x_n$  from a Bernoulli distribution with parameter  $\theta$ . Show that for testing  $\theta = \frac{1}{2}$  versus  $\theta \neq \frac{1}{2}$ , we always get  $\mathcal{W} \geq \mathcal{LM}$ . Show also that for testing  $\theta = \frac{2}{3}$  versus  $\theta \neq \frac{2}{3}$ , we get  $\mathcal{W} \leq \mathcal{LM}$  for  $\frac{1}{3} \leq \bar{x} \leq \frac{2}{3}$  and  $\mathcal{W} \geq \mathcal{LM}$  for  $0 < \bar{x} \leq \frac{1}{3}$  or  $\frac{2}{3} \leq \bar{x} \leq 1$ .

<sup>1</sup>This problem closely follows Badi H. Baltagi (2000) Conflict Among Criteria for Testing Hypotheses: Examples from Non-Normal Distributions. *Econometric Theory* 16, Problem 00.2.4.

## 10.4 Invariance of ML tests to reparameterizations of null

<sup>2</sup>Consider the hypothesis

$$H_0 : h(\theta) = 0,$$

where  $h : \mathbb{R}^k \rightarrow \mathbb{R}^q$ . It is possible to recast the hypothesis  $H_0$  in an equivalent form

$$H_0 : g(\theta) = 0,$$

where  $g : \mathbb{R}^k \rightarrow \mathbb{R}^q$  is such that  $g(\theta) = f(h(\theta)) - f(0)$  for some one-to-one function  $f : \mathbb{R}^q \rightarrow \mathbb{R}^q$ .

1. Show that the  $\mathcal{LR}$  statistic is invariant to such reparameterization.
2. Show that the  $\mathcal{W}$  statistic is invariant to such reparameterization when  $f$  is linear, but may not be when  $f$  is nonlinear.
3. Suppose that  $\theta \in \mathbb{R}^2$  and reparameterize  $H_0 : \theta_1 = \theta_2$  as  $(\theta_1 - \alpha) / (\theta_2 - \alpha) = 1$  for some  $\alpha$ . Show that the  $\mathcal{W}$  statistic may be made as close to zero as desired by manipulating  $\alpha$ . What value of  $\alpha$  gives the largest possible value to the  $\mathcal{W}$  statistic?

## 10.5 Misspecified maximum likelihood

1. Suppose that the nonlinear regression model

$$\mathbb{E}[y|x] = g(x, \beta)$$

is estimated by maximum likelihood based on the conditional homoskedastic normal distribution, although the true conditional distribution is from a different family. Provide a simple argument why the ML estimator of  $\beta$  is nevertheless consistent.

2. Suppose we know the true density  $f(z|\theta)$  up to the parameter  $\theta$ , but instead of using  $\log f(z|q)$  in the objective function of the extremum problem which would give the ML estimate, we use  $f(z|q)$  itself. What asymptotic properties do you expect from the resulting estimator of  $\theta$ ? Will it be consistent? Will it be asymptotically normal?
3. Suppose that a stationary and ergodic conditionally heteroskedastic time series regression

$$\mathbb{E}[y_t | I_{t-1}] = x_t' \beta,$$

where  $x_t$  contains various lagged  $y_t$ 's, is estimated by maximum likelihood based on a conditionally normal distribution  $\mathcal{N}(x_t' \beta, \sigma_t^2)$ , with  $\sigma_t^2$  depending on past data whose parametric form, however, does not match the true conditional variance  $\mathbb{V}[y_t | I_{t-1}]$ . Determine whether or not the resulting estimator provides a consistent estimate of  $\beta$ .

---

<sup>2</sup>This problem closely follows discussion in the book Ruud, Paul (2000) *An Introduction to Classical Econometric Theory*; Oxford University Press.



## 10.6 Individual effects

Suppose  $\{(x_i, y_i)\}_{i=1}^n$  is a serially independent sample from a sequence of jointly normal distributions with  $\mathbb{E}[x_i] = \mathbb{E}[y_i] = \mu_i$ ,  $\mathbb{V}[x_i] = \mathbb{V}[y_i] = \sigma^2$ , and  $\mathbb{C}[x_i, y_i] = 0$  (i.e.,  $x_i$  and  $y_i$  are independent with common but varying means and a constant common variance). All parameters are unknown. Derive the maximum likelihood estimate of  $\sigma^2$  and show that it is inconsistent. Explain why. Find an estimator of  $\sigma^2$  which would be consistent.

## 10.7 Irregular confidence interval

Let  $x_1, \dots, x_n$  be a random sample from a population of  $x$  distributed uniformly on  $[0, \theta]$ . Construct an asymptotic confidence interval for  $\theta$  with significance level 5% by employing a maximum likelihood approach.

## 10.8 Trivial parameter space

Consider a parametric model with density  $f(X|\theta_0)$ , known up to a parameter  $\theta_0$ , but with  $\Theta = \{\theta_1\}$ , i.e. the parameter space is reduced to only one element. What is an ML estimator of  $\theta_0$ , and what are its asymptotic properties?

## 10.9 Nuisance parameter in density

Let random vector  $Z \equiv (Y, X)'$  have a joint density of the form

$$f(Z|\theta_0) = f_c(Y|X, \gamma_0, \delta_0) f_m(X|\delta_0),$$

where  $\theta_0 \equiv (\gamma_0, \delta_0)$ , both  $\gamma_0$  and  $\delta_0$  are scalar parameters, and  $f_c$  and  $f_m$  denote the conditional and marginal distributions, respectively. Let  $\hat{\theta}_c \equiv (\hat{\gamma}_c, \hat{\delta}_c)$  be the conditional ML estimators of  $\gamma_0$  and  $\delta_0$ , and  $\hat{\delta}_m$  be the marginal ML estimator of  $\delta_0$ . Now define

$$\tilde{\gamma} \equiv \arg \max_{\gamma} \sum_i \ln f_c(y_i|x_i, \gamma, \hat{\delta}_m),$$

a two-step estimator of subparameter  $\gamma_0$  which uses marginal ML to obtain a preliminary estimator of the “nuisance parameter”  $\delta_0$ . Find the asymptotic distribution of  $\tilde{\gamma}$ . How does it compare to that for  $\hat{\gamma}_c$ ? You may assume all the needed regularity conditions for consistency and asymptotic normality to hold.

## 10.10 MLE versus OLS

Consider the model where  $y$  is regressed only on a constant:

$$y = \alpha + e,$$

where  $e$  conditioned on  $x$  is distributed as  $\mathcal{N}(0, x^2\sigma^2)$ , the random variable  $x$  is not present in the regression,  $\sigma^2$  is unknown,  $y$  and  $x$  are observable,  $e$  is unobservable. The collection of pairs  $\{(y_i, x_i)\}_{i=1}^n$  is a random sample.

1. Find the OLS estimator  $\hat{\alpha}_{OLS}$  of  $\alpha$ . Is it unbiased? Consistent? Obtain its asymptotic distribution. Is  $\hat{\alpha}_{OLS}$  the best linear unbiased estimator for  $\alpha$ ?
2. Find the ML estimator  $\hat{\alpha}_{ML}$  of  $\alpha$  and derive its asymptotic distribution. Is  $\hat{\alpha}_{ML}$  unbiased? Is  $\hat{\alpha}_{ML}$  asymptotically more efficient than  $\hat{\alpha}_{OLS}$ ? Does your conclusion contradict your answer to the last question of part 1? Why or why not?

## 10.11 MLE versus GLS

Consider the following normal linear regression model with conditional heteroskedasticity of known form. Conditional on  $x$ , the dependent variable  $y$  is normally distributed with

$$\mathbb{E}[y|x] = x'\beta, \quad \mathbb{V}[y|x] = \sigma^2 (x'\beta)^2.$$

Available is a random sample  $(x_1, y_1), \dots, (x_n, y_n)$ . Describe a feasible generalized least squares estimator for  $\beta$  based on the OLS estimator for  $\beta$ . Show that this GLS estimator is asymptotically less efficient than the maximum likelihood estimator. Explain the source of inefficiency.

## 10.12 MLE in heteroskedastic time series regression

Assume that data  $(y_t, x_t)$ ,  $t = 1, 2, \dots, T$ , are stationary and ergodic and generated by

$$y_t = \alpha + \beta x_t + u_t,$$

where  $u_t|x_t, y_{t-1}, x_{t-1}, y_{t-2}, \dots \sim \mathcal{N}(0, \sigma_t^2)$  and  $x_t|y_{t-1}, x_{t-1}, y_{t-2}, x_{t-2}, \dots \sim \mathcal{N}(0, v)$ . Explain, without going into deep math, how to find estimates and their standard errors *for all parameters* when:

1. The entire  $\sigma_t^2$  as a function of  $x_t$  is fully known.
2. The values of  $\sigma_t^2$  at  $t = 1, 2, \dots, T$  are known.
3. It is known that  $\sigma_t^2 = (\theta + \delta x_t)^2$ , but the parameters  $\theta$  and  $\delta$  are unknown.
4. It is known that  $\sigma_t^2 = \theta + \delta u_{t-1}^2$ , but the parameters  $\theta$  and  $\delta$  are unknown.
5. It is only known that  $\sigma_t^2$  is stationary.

## 10.13 Does the link matter?

<sup>3</sup>Consider a binary random variable  $y$  and a scalar random variable  $x$  such that

$$\mathbb{P}\{y = 1|x\} = F(\alpha + \beta x),$$

where the link  $F(\cdot)$  is a continuous distribution function. Show that when  $x$  assumes only two different values, the value of the log-likelihood function evaluated at the maximum likelihood estimates of  $\alpha$  and  $\beta$  is independent of the form of the link function. What are the maximum likelihood estimates of  $\alpha$  and  $\beta$ ?

## 10.14 Maximum likelihood and binary variables

Suppose  $z$  and  $y$  are discrete random variables taking values 0 or 1. The distribution of  $z$  and  $y$  is given by

$$\mathbb{P}\{z = 1\} = \alpha, \quad \mathbb{P}\{y = 1|z\} = \frac{e^{\gamma z}}{1 + e^{\gamma z}}, \quad z = 0, 1.$$

Here  $\alpha$  and  $\gamma$  are scalar parameters of interest. Find the ML estimator of  $(\alpha, \gamma)$  from a random sample giving explicit formulas whenever possible, and derive its asymptotic distribution.

## 10.15 Maximum likelihood and binary dependent variable

Suppose  $y$  is a discrete random variable taking values 0 or 1 representing some choice of an individual. The distribution of  $y$  given the individual's characteristic  $x$  is

$$\mathbb{P}\{y = 1|x\} = \frac{e^{\gamma x}}{1 + e^{\gamma x}},$$

where  $\gamma$  is the scalar parameter of interest. The data  $\{(y_i, x_i)\}_{i=1}^n$  are a random sample. When deriving various estimators, try to make the formulas as explicit as possible.

1. Derive the ML estimator of  $\gamma$  and its asymptotic distribution.
2. Find the (nonlinear) regression function by regressing  $y$  on  $x$ . Derive the NLLS estimator of  $\gamma$  and its asymptotic distribution.
3. Show that the regression you obtained in part 2 is heteroskedastic. Setting weights  $\omega(x)$  equal to the variance of  $y$  conditional on  $x$ , derive the WNLLS estimator of  $\gamma$  and its asymptotic distribution.
4. Write out the systems of moment conditions implied by the ML, NLLS and WNLLS problems of parts 1–3.
5. Rank the three estimators in terms of asymptotic efficiency. Do any of your findings appear unexpected? Give intuitive explanation for anything unusual.

---

<sup>3</sup>This problem closely follows Joao M.C. Santos Silva (1999) Does the link matter? *Econometric Theory* 15, Problem 99.5.3.

## 10.16 Poisson regression

Some random variables called counts (like number of patent applications by a firm, or number of doctor visits by a patient) take non-negative integer values  $0, 1, 2, \dots$ . Suppose that  $y$  is a count variable, and that, conditional on scalar random variable  $x$ , its density is Poisson:

$$f(y|x, \alpha, \beta) = \frac{\exp(-\lambda(x, \alpha, \beta))\lambda(x, \alpha, \beta)^y}{y!}, \quad \text{where } \lambda(x, \alpha, \beta) = \exp(\alpha + \beta x).$$

Here  $\alpha$  and  $\beta$  are unknown scalar parameters. Assume that  $x$  has a nondegenerate distribution (i.e., is not a constant). The data  $\{(x_i, y_i)\}_{i=1}^n$  is a random sample.

1. Derive the three asymptotic ML tests ( $\mathcal{W}$ ,  $\mathcal{LR}$ ,  $\mathcal{LM}$ ) for the null hypothesis  $H_0 : \beta = 0$ , giving fully implementable formulas (in explicit form whenever possible) depending only on the data.
2. Suppose the researcher has misspecified the conditional mean and uses

$$\lambda(x, \alpha, \beta) = \alpha + \exp(\beta x).$$

Will the coefficient  $\beta$  be consistently estimated?

## 10.17 Bootstrapping ML tests

1. For the likelihood ratio test of  $H_0 : g(\theta) = 0$ , we use the statistic

$$\mathcal{LR} = 2 \left( \max_{q \in \Theta} \ell_n(q) - \max_{q \in \Theta, g(q)=0} \ell_n(q) \right).$$

Write out the formula for the bootstrap statistic  $\mathcal{LR}^*$ .

2. For the Lagrange Multiplier test of  $H_0 : g(\theta) = 0$ , we use the statistic

$$\mathcal{LM} = \frac{1}{n} \sum_i s(z_i, \hat{\theta}_{ML}^R)' \hat{\mathcal{I}}^{-1} \sum_i s(z_i, \hat{\theta}_{ML}^R).$$

Write out the formula for the bootstrap statistic  $\mathcal{LM}^*$ .

# 11. INSTRUMENTAL VARIABLES

---

## 11.1 Invalid 2SLS

Consider the model

$$y = \alpha z^2 + u, \quad z = \pi x + v,$$

where  $\mathbb{E}[(u, v)' | x] = 0$  and  $\mathbb{V}[(u, v)' | x] = \Sigma$ , with  $\Sigma$  unknown. The triples  $\{(x_i, z_i, y_i)\}_{i=1}^n$  constitute a random sample.

1. Show that  $\alpha$ ,  $\pi$  and  $\Sigma$  are identified. Suggest analog estimators for these parameters.
2. Consider the following two stage estimation method. In the first stage, regress  $z$  on  $x$  and define  $\hat{z} = \hat{\pi}x$ , where  $\hat{\pi}$  is the OLS estimator. In the second stage, regress  $y$  in  $\hat{z}^2$  to obtain the least squares estimate of  $\alpha$ . Show that the resulting estimator of  $\alpha$  is inconsistent.
3. Suggest a method in the spirit of 2SLS for estimating  $\alpha$  consistently.

## 11.2 Consumption function

Consider the consumption function

$$C_t = \alpha + \lambda Y_t + e_t, \tag{11.1}$$

where  $C_t$  is aggregate consumption at  $t$ , and  $Y_t$  is aggregate income at  $t$ . The ordinary least squares (OLS) estimation applied to (11.1) may give an inconsistent estimate of the marginal propensity to consume (MPC)  $\lambda$ . The remedy suggested by Haavelmo lies in treating the aggregate income as endogenous:

$$Y_t = C_t + I_t + G_t, \tag{11.2}$$

where  $I_t$  is aggregate investment at  $t$ , and  $G_t$  is government consumption at  $t$ , and both variables are exogenous. Assume that the shock  $e_t$  is mean zero IID across time, and all variables are jointly stationary and ergodic. A sample of size  $T$  containing  $Y_t$ ,  $C_t$ ,  $I_t$ , and  $G_t$  is available.

1. Show that the OLS estimator of  $\lambda$  is indeed inconsistent. Compute the amount and direction of this inconsistency.
2. Econometrician A intends to estimate  $(\alpha, \lambda)'$  by running 2SLS on (11.1) using the instrumental vector  $(1, I_t, G_t)'$ . Econometrician B argues that it is not necessary to use this relatively complicated estimator since running simple IV on (11.1) using the instrumental vector  $(1, I_t + G_t)'$  will do the same. Is econometrician B right?
3. Econometrician C regresses  $Y_t$  on a constant and  $C_t$ , and obtains corresponding OLS estimates  $(\hat{\theta}_0, \hat{\theta}_C)'$ . Econometrician D regresses  $Y_t$  on a constant,  $C_t$ ,  $I_t$ , and  $G_t$  and obtains corresponding OLS estimates  $(\hat{\phi}_0, \hat{\phi}_C, \hat{\phi}_I, \hat{\phi}_G)'$ . What values do parameters  $\hat{\theta}_C$  and  $\hat{\phi}_C$  consistently estimate?

## 11.3 Optimal combination of instruments

Suppose you have the following specification, where  $e$  may be correlated with  $x$ :

$$y = \beta x + e.$$

1. You have instruments  $z$  and  $\zeta$  which are mutually uncorrelated. What are their necessary properties to provide consistent IV estimators  $\hat{\beta}_z$  and  $\hat{\beta}_\zeta$ ? Derive the asymptotic distributions of these estimators.
2. Calculate the optimal IV estimator as a linear combination of  $\hat{\beta}_z$  and  $\hat{\beta}_\zeta$ .
3. You notice that  $\hat{\beta}_z$  and  $\hat{\beta}_\zeta$  are not that close together. Give a test statistic which allows you to decide if they are estimating the same parameter. If the test rejects, what assumptions are you rejecting?

## 11.4 Trade and growth

In the paper “Does Trade Cause Growth?” (*American Economic Review*, June 1999), Jeffrey Frankel and David Romer study the effect of trade on income. Their simple specification is

$$\log Y_i = \alpha + \beta T_i + \gamma W_i + \varepsilon_i, \quad (11.3)$$

where  $Y_i$  is per capita income,  $T_i$  is international trade,  $W_i$  is within-country trade, and  $\varepsilon_i$  reflects other influences on income. Since the latter is likely to be correlated with the trade variables, Frankel and Romer decide to use instrumental variables to estimate the coefficients in (11.3). As instruments, they use a country’s proximity to other countries  $P_i$  and its size  $S_i$ , so that

$$T_i = \psi + \phi P_i + \delta_i \quad (11.4)$$

and

$$W_i = \eta + \lambda S_i + \nu_i, \quad (11.5)$$

where  $\delta_i$  and  $\nu_i$  are the best linear prediction errors.

1. As the key identifying assumption, the authors use the fact that countries’ geographical characteristics  $P_i$  and  $S_i$  are uncorrelated with the error term in (11.3). Provide an economic rationale for this assumption and a detailed explanation how to estimate (11.3) when one has data on  $Y$ ,  $T$ ,  $W$ ,  $P$  and  $S$  for a list of countries.
2. Unfortunately, data on within-country trade are not available. Determine if it is possible to estimate any of the coefficients in (11.3) without further assumptions. If it is, provide all the details on how to do it.
3. In order to be able to estimate key coefficients in (11.3), the authors add another identifying assumption that  $P_i$  is uncorrelated with the error term in (11.5). Provide a detailed explanation how to estimate (11.3) when one has data on  $Y$ ,  $T$ ,  $P$  and  $S$  for a list of countries.
4. The authors estimated an equation similar to (11.3) by OLS and IV and found out that the IV estimates are greater than the OLS estimates. One explanation may be that the discrepancy is due to a sampling error. Provide another, more econometric, explanation why there is a discrepancy and what the reason is that the IV estimates are larger.

# 12. GENERALIZED METHOD OF MOMENTS

## 12.1 Nonlinear simultaneous equations

Let

$$y = \beta x + u, \quad x = \gamma y^2 + v,$$

where  $x$  and  $y$  are observable, but  $u$  and  $v$  are not. The data  $\{(x_i, y_i)\}_{i=1}^n$  is a random sample.

1. Suppose we know that  $\mathbb{E}[u] = \mathbb{E}[v] = 0$ . When are  $\beta$  and  $\gamma$  identified? Propose analog estimators for these parameters.
2. Let also be known that  $\mathbb{E}[uv] = 0$ . Propose a method to estimate  $\beta$  and  $\gamma$  as efficiently as possible given the above information. What is the asymptotic distribution of your estimator?

## 12.2 Improved GMM

Consider GMM estimation with the use of the moment function

$$m(x, y, q) = \begin{pmatrix} x - q \\ y \end{pmatrix}.$$

Determine under what conditions the second restriction helps in reducing the asymptotic variance of the GMM estimator of  $\theta$ .

## 12.3 Minimum Distance estimation

Consider a similar to GMM procedure called the Minimum Distance (MD) estimation. Suppose we want to estimate a parameter  $\gamma_0 \in \Gamma$  implicitly defined by  $\theta_0 = s(\gamma_0)$ , where  $s : \mathbb{R}^k \rightarrow \mathbb{R}^\ell$  with  $\ell \geq k$ , and available is an estimator  $\hat{\theta}$  of  $\theta_0$  with asymptotic properties

$$\hat{\theta} \xrightarrow{p} \theta_0, \quad \sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N}(0, V_{\hat{\theta}}).$$

Also suppose that available is a symmetric and positive definite estimator  $\hat{V}_{\hat{\theta}}$  of  $V_{\hat{\theta}}$ . The MD estimator is defined as

$$\hat{\gamma}_{MD} = \arg \min_{\gamma \in \Gamma} \left( \hat{\theta} - s(\gamma) \right)' \hat{W} \left( \hat{\theta} - s(\gamma) \right),$$

where  $\hat{W}$  is some symmetric positive definite data-dependent matrix consistent for a symmetric positive definite weight matrix  $W$ . Assume that  $\Gamma$  is compact,  $s(\gamma)$  is continuously differentiable with full rank matrix of derivatives  $S(\gamma) = \partial s(\gamma) / \partial \gamma'$  on  $\Gamma$ ,  $\gamma_0$  is unique and all needed moments exist.

1. Give an informal argument for consistency of  $\hat{\gamma}_{MD}$ . Derive the asymptotic distribution of  $\hat{\gamma}_{MD}$ .
2. Find the optimal choice for the weight matrix  $W$  and suggest its consistent estimator.
3. Develop a specification test, i.e. of the hypothesis  $H_0 : \exists \gamma_0$  such that  $\theta_0 = s(\gamma_0)$ .
4. Apply parts 1–3 to the following problem. Suppose that we have an autoregression of order 2 without a constant term:

$$(1 - \rho L)^2 y_t = \varepsilon_t,$$

where  $|\rho| < 1$ ,  $L$  is the lag operator, and  $\varepsilon_t$  is IID(0,  $\sigma^2$ ). Written in another form, the model is

$$y_t = \theta_1 y_{t-1} + \theta_2 y_{t-2} + \varepsilon_t,$$

and  $(\theta_1, \theta_2)'$  may be efficiently estimated by OLS. The target, however, is to estimate  $\rho$  and verify that both autoregressive roots are indeed equal.

## 12.4 Formation of moment conditions

1. Let  $z$  be a scalar random variable. Let it be known that  $z$  has mean  $\mu$  and that its fourth central moment equals three times its squared variance (like for a normal random variable). Formulate a system of moment conditions for GMM estimation of  $\mu$ .
2. Consider the AR(1)–ARCH(1) model

$$y_t = \rho y_{t-1} + e_t, \quad \mathbb{E}[e_t | I_{t-1}] = 0, \quad \mathbb{E}[e_t^2 | I_{t-1}] = \omega + \gamma e_{t-1}^2,$$

where  $I_{t-1}$  embeds information on the history from  $y_{t-1}$  to the past. Such models are typically estimated by ML, but may also be estimated by GMM, if desired. Suppose you have data on  $y_t$  for  $t = 1, \dots, T$ . Construct an overidentifying set of moment conditions to be used by GMM.

## 12.5 What CMM estimates

Let  $g(z, q)$  be a function such that dimensions of  $g$  and  $q$  are identical, and let  $z_1, \dots, z_n$  be a random sample. Note that nothing is said about moment conditions. Define  $\hat{\theta}$  as the solution to

$$\sum_{i=1}^n g(z_i, q) = 0.$$

What is the probability limit of  $\hat{\theta}$ ? What is the asymptotic distribution of  $\hat{\theta}$ ?



## 12.6 Trinity for GMM

Derive the three classical tests ( $\mathcal{W}$ ,  $\mathcal{LR}$ ,  $\mathcal{LM}$ ) for the composite null

$$H_0 : \theta \in \Theta_0 \equiv \{\theta : h(\theta) = 0\},$$

where  $h : \mathbb{R}^k \rightarrow \mathbb{R}^q$ , for the efficient GMM case. The analog for the Likelihood Ratio test will be called the *Distance Difference test*. Hint: treat the GMM objective function as the “normalized loglikelihood”, and its derivative as the “sample score”.

## 12.7 All about J

1. Show that the J-test statistic diverges to infinity when the system of moment conditions is misspecified.
2. Provide an example showing that the J-test statistic need not be asymptotically chi-squared with degrees of freedom equalling the degree of overidentification under valid moment restrictions, if one uses a non-efficient GMM.
3. Suppose an econometrician estimates parameters of a time series regression by GMM after having chosen an overidentifying vector of instrumental variables. He performs the overidentification test and claims: “A big value of the  $\mathcal{J}$ -statistic is an evidence against validity of the chosen instruments”. Comment on this claim.

## 12.8 Interest rates and future inflation

Frederic Mishkin in early 90’s investigated whether the term structure of current nominal interest rates can give information about future path of inflation. He specified the following econometric model:

$$\pi_t^m - \pi_t^n = \alpha_{m,n} + \beta_{m,n} (i_t^m - i_t^n) + \eta_t^{m,n}, \quad \mathbb{E}_t [\eta_t^{m,n}] = 0, \quad (12.1)$$

where  $\pi_t^k$  is  $k$ -periods-into-the-future inflation rate,  $i_t^k$  is the current nominal interest rate for  $k$ -periods-ahead maturity, and  $\eta_t^{m,n}$  is the prediction error.

1. Show how (12.1) can be obtained from the conventional econometric model that tests the hypothesis of conditional unbiasedness of interest rates as predictors of inflation. What restriction on the parameters in (12.1) implies that the term structure provides *no* information about future shifts in inflation? Determine the autocorrelation structure of  $\eta_t^{m,n}$ .
2. Describe in detail how you would test the hypothesis that the term structure provides no information about future shifts in inflation, by using overidentifying GMM and asymptotic theory. Make sure that you discuss such issues as selection of instruments, construction of the optimal weighting matrix, construction of the GMM objective function, estimation of asymptotic variance, etc.

3. Mishkin obtained the following results (standard errors are in parentheses):

$m, n$ (months)	$\alpha_{m,n}$	$\beta_{m,n}$	$t$ -test of $\beta_{m,n} = 0$	$t$ -test of $\beta_{m,n} = 1$
3, 1	0.1421 (0.1851)	-0.3127 (0.4498)	-0.70	2.92
6, 3	0.0379 (0.1427)	0.1813 (0.5499)	0.33	1.49
9, 6	0.0826 (0.0647)	0.0014 (0.2695)	0.01	3.71

Discuss and interpret the estimates and results of hypotheses tests.

## 12.9 Spot and forward exchange rates

Consider a simple problem of prediction of spot exchange rates by forward rates:

$$s_{t+1} - s_t = \alpha + \beta(f_t - s_t) + e_{t+1}, \quad \mathbb{E}_t[e_{t+1}] = 0, \quad \mathbb{E}_t[e_{t+1}^2] = \sigma^2,$$

where  $s_t$  is the spot rate at  $t$ ,  $f_t$  is the forward rate for one-month forwards at  $t$ , and  $\mathbb{E}_t$  denotes expectation conditional on time  $t$  information. The current spot rate is subtracted to achieve stationarity. Suppose the researcher decides to use ordinary least squares to estimate  $\alpha$  and  $\beta$ . Recall that the moment conditions used by the OLS estimator are

$$\mathbb{E}[e_{t+1}] = 0, \quad \mathbb{E}[(f_t - s_t)e_{t+1}] = 0. \quad (12.2)$$

1. Beside (12.2), there are other moment conditions that can be used in estimation:

$$\mathbb{E}[(f_{t-k} - s_{t-k})e_{t+1}] = 0,$$

because  $f_{t-k} - s_{t-k}$  belongs to information at time  $t$  for any  $k \geq 1$ . Consider the case  $k = 1$  and show that such moment condition is redundant.

2. Beside (12.2), there is another moment condition that can be used in estimation:

$$\mathbb{E}[(f_t - s_t)(f_{t+1} - f_t)] = 0,$$

because information at time  $t$  should be unable to predict future movements in forward rates. Although this moment condition does not involve  $\alpha$  or  $\beta$ , its use may improve efficiency of estimation. Under what condition is the efficient GMM estimator using both moment conditions as efficient as the OLS estimator? Is this condition likely to be satisfied in practice?

## 12.10 Returns from financial market

Suppose you have the following parametric model supported by one of financial theories, for the daily dynamics of 3-month Treasury bills return  $r_t$ :

$$r_{t+1} - r_t = \beta_0 + \beta_1 r_t + \beta_2 r_t^2 + \beta_3 r_t^{-1} + \varepsilon_{t+1},$$

where

$$\varepsilon_{t+1}|I_t \sim \mathcal{N}\left(0, \sigma^2 r_t^{2\gamma}\right),$$

where  $I_t$  contains information on  $r_t, r_{t-1}, \dots$ . Such model arises from a discretization of a continuous time process defined by a diffusion model for returns. Suppose you have a long stationary and ergodic sample  $\{r_t\}_{t=1}^T$ .

In answering the following questions, first of all write out the mentioned estimators, giving explicit formulas whenever possible.

1. Derive the asymptotic distributions of the OLS and feasible GLS estimators of the regression coefficients. Write out moment conditions that correspond to these estimators.
2. Derive the asymptotic distributions of the ML estimator of all parameters. Write out a system of moment conditions that corresponds to this estimator.
3. Construct an exactly identifying system of moment conditions from (a) the moment conditions implicitly used by OLS estimation of the regression function, and (b) the moment conditions implicitly used by NLLS estimation of the skedastic function. Derive the asymptotic distribution of the resulting method of moments estimator of all parameters.
4. Compare the asymptotic efficiency of the estimators from parts 1–3 for the regression parameters, and of the estimators from parts 2–3 for  $\sigma^2$  and  $\gamma$ .

## 12.11 Instrumental variables in ARMA models

1. Consider an  $AR(1)$  model  $x_t = \rho x_{t-1} + e_t$  with  $\mathbb{E}[e_t|I_{t-1}] = 0$ ,  $\mathbb{E}[e_t^2|I_{t-1}] = \sigma^2$ , and  $|\rho| < 1$ . We can look at this as an instrumental variables regression that implies, among others, instruments  $x_{t-1}, x_{t-2}, \dots$ . Find the asymptotic variance of the instrumental variables estimator that uses instrument  $x_{t-j}$ , where  $j = 1, 2, \dots$ . What does your result suggest on what the optimal instrument must be?
2. Consider an  $ARMA(1, 1)$  model  $y_t = \alpha y_{t-1} + e_t - \theta e_{t-1}$  with  $|\alpha| < 1$ ,  $|\theta| < 1$  and  $\mathbb{E}[e_t|I_{t-1}] = 0$ . Suppose you want to estimate  $\alpha$  by just-identifying IV. What instrument would you use and why?

## 12.12 Hausman may not work

1. Suppose we want to perform a Hausman test for validity of instruments  $z$  for a linear conditionally homoskedastic mean regression model. To this end, we compute OLS estimator  $\hat{\beta}_{OLS}$  and 2SLS estimator  $\hat{\beta}_{2SLS}$  of  $\beta$  using  $x$  and  $z$  as instruments. Explain carefully why the Hausman test based on comparison of  $\hat{\beta}_{OLS}$  and  $\hat{\beta}_{2SLS}$  will not work.
2. Suppose that in a context of a linear model with (overidentifying) instrumental variables a researcher intends to test for conditional homoskedasticity using a Hausman test based on the difference between the 2SLS and GMM estimators of parameters. Explain carefully why this is not a good idea.

## 12.13 Testing moment conditions

In the linear model

$$y = x'\beta + u$$

under random sampling and the unconditional moment restriction  $\mathbb{E}[xu] = 0$ , suppose you wanted to test the additional moment restriction  $\mathbb{E}[xu^3] = 0$ , which might be implied by conditional symmetry of the error terms  $u$ .

A natural way to test for the validity of this extra moment condition would be to efficiently estimate the parameter vector  $\beta$  both with and without the additional restriction, and then to check whether the corresponding estimates differ significantly. Devise such a test and give step-by-step instructions for carrying it out.

## 12.14 Bootstrapping OLS

We know that one should use recentering when bootstrapping a GMM estimator. We also know that the OLS estimator is one of GMM estimators. However, when we bootstrap the OLS estimator, we calculate

$$\hat{\beta}^* = (\mathcal{X}^{*\prime}\mathcal{X}^*)^{-1}\mathcal{X}^{*\prime}\mathcal{Y}^*$$

at each bootstrap repetition, and do not recenter. Resolve the contradiction.

## 12.15 Bootstrapping DD

The *Distance Difference test* statistic for testing the composite null  $H_0 : h(\theta) = 0$  is defined as

$$\mathcal{DD} = n \left[ \min_{q:h(q)=0} \mathcal{Q}_n(q) - \min_q \mathcal{Q}_n(q) \right],$$

where  $\mathcal{Q}_n(q)$  is the GMM objective function

$$\mathcal{Q}_n(q) = \left( \frac{1}{n} \sum_{i=1}^n m(z_i, q) \right)' \hat{Q}_{mm}^{-1} \left( \frac{1}{n} \sum_{i=1}^n m(z_i, q) \right),$$

where  $\hat{Q}_{mm}$  consistently estimates  $Q_{mm} = \mathbb{E}[m(z, \theta)m(z, \theta)']$ . It is known that, as the sample size  $n$  tends to infinity,  $\mathcal{DD} \xrightarrow{d} \chi_{\dim(q)}^2$ . Write out a detailed formula for the bootstrap statistic  $\mathcal{DD}^*$ .

# 13. PANEL DATA

---

## 13.1 Alternating individual effects

Suppose that the unobservable individual effects in a one-way error component model are different across odd and even periods:

$$\begin{aligned} y_{it} &= \mu_i^O + x'_{it}\beta + v_{it} & \text{for odd } t, \\ y_{it} &= \mu_i^E + x'_{it}\beta + v_{it} & \text{for even } t, \end{aligned} \quad (*)$$

where  $t = 1, 2, \dots, 2T$ ,  $i = 1, \dots, n$ . Note that there are  $2T$  observations for each individual. We will call (\*) “alternating effects” specification. As usual, we assume that  $v_{it}$  are  $IID(0, \sigma_v^2)$  independent of  $x$ 's.

1. There are two ways to arrange the observations: (a) in the usual way, first by individual, then by time for each individual; (b) first all “odd” observations in the usual order, then all “even” observations, so it is as though there are  $2N$  “individuals” each having  $T$  observations. Find out the  $Q$ -matrices that wipe out individual effects for both arrangements and explain how they transform the original equations. For the rest of the problem, choose the  $Q$ -matrix to your liking.
2. Treating individual effects as fixed, describe the Within estimator and its properties. Develop an  $F$ -test for individual effects, allowing heterogeneity across odd and even periods.
3. Treating individual effects as random and assuming their independence of  $x$ 's,  $v$ 's and each other, propose a feasible GLS procedure. Consider two cases: (a) when the variance of “alternating effects” is the same:  $\mathbb{V}[\mu_i^O] = \mathbb{V}[\mu_i^E] = \sigma_\mu^2$ , (b) when the variance of “alternating effects” is different:  $\mathbb{V}[\mu_i^O] = \sigma_O^2$ ,  $\mathbb{V}[\mu_i^E] = \sigma_E^2$ ,  $\sigma_O^2 \neq \sigma_E^2$ .

## 13.2 Time invariant regressors

Consider a panel data model

$$y_{it} = x'_{it}\beta + z_i\gamma + \mu_i + v_{it}, \quad i = 1, 2, \dots, n, \quad t = 1, 2, \dots, T,$$

where  $n$  is large and  $T$  is small. One wants to estimate  $\beta$  and  $\gamma$ .

1. Explain how to efficiently estimate  $\beta$  and  $\gamma$  under (a) fixed effects, (b) random effects, whenever it is possible. State clearly all assumptions that you will need.
2. Consider the following proposal to estimate  $\gamma$ . At the first step, estimate the model  $y_{it} = x'_{it}\beta + \pi_i + v_{it}$  by the least squares dummy variables approach. At the second step, take these estimates  $\hat{\pi}_i$  and estimate the coefficient of the regression of  $\hat{\pi}_i$  on  $z_i$ . Investigate the resulting estimator of  $\gamma$  for consistency. Can you suggest a better estimator of  $\gamma$ ?

## 13.3 Within and Between

Recall the standard one-way static error component model. For simplicity, assume that all data are in deviations from their means so that there is no intercept. Denote by  $\hat{\beta}_W$  and  $\hat{\beta}_B$  the Within and Between estimators of structural parameters, respectively. Denote the matrix of right side variables by  $X$ . Show that under random effects,  $\hat{\beta}_W$  and  $\hat{\beta}_B$  are uncorrelated conditionally on  $X$ . Develop an asymptotic test for random effects based on the difference between  $\hat{\beta}_W$  and  $\hat{\beta}_B$ . In particular, derive the asymptotic distribution of your test statistic when the random effects model is valid, and show that it asymptotically diverges to infinity when the random effects are inappropriate.

## 13.4 Panels and instruments

Recall the standard one-way static error component model with random effects where the regressors are denoted by  $x_{it}$ . Suppose that these regressors are endogenous. Additional data  $z_{it}$  for each individual at each time period are given, and let these additional variables be independent of the errors and correlated with the regressors. Hence, one can use these variables as instrumental variables (IV). Using knowledge of linear panel regressions and IV theory, answer the following questions about various estimators, not worrying about standard errors.

1. Develop the IV-Within estimator and show that it will be identical irrespective of whether one uses original instruments or Within-transformed instruments.
2. Develop the IV-Between estimator and show that it will be identical irrespective of whether one uses original instruments or Between-transformed instruments.
3. What will happen if we use Within-transformed instruments in the Between regression? If we use Between-transformed instruments in the Within regression?
4. Develop the IV-GLS estimator and show that now it matters whether one uses original instruments or GLS-transformed instruments. Intuitively, which way would you prefer?
5. What will happen if we use Within-transformed instruments in the GLS-transformed regression? If we use Between-transformed instruments in the GLS-transformed regression?

## 13.5 Differencing transformations

Evaluate the following proposals.

1. In a one-way error component model with fixed effects, instead of using individual dummies, one can alternatively eliminate individual effects by taking the first differencing (FD) transformation. After this procedure one has  $n(T - 1)$  equations without individual effects, so the vector  $\beta$  of structural parameters can be estimated by OLS.
2. Recall the standard dynamic panel data model. The individual heterogeneity may be removed not only by first differencing, but also by (for example) subtracting the equation corresponding to  $t = 2$  from each other equation for the same individual.

## 13.6 Nonlinear panel data model

We know (see Problem 9.3) that the NLLS estimator of  $\alpha$  and  $\beta$

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \arg \min_{a,b} \sum_{i=1}^n \left( (y_i + a)^2 - bx_i \right)^2$$

in the nonlinear model

$$(y + \alpha)^2 = \beta x + e,$$

where  $e$  is independent of  $x$ , is in general inconsistent.

Suppose that there is a panel  $\{x_{it}, y_{it}\}_{i=1}^n \}_{t=1}^T$ , where  $n$  is large and  $T$  is small, so that there is an opportunity to control individual heterogeneity. Write out a one-way error component model assuming the same functional form but allowing for individual heterogeneity in the form of random effects. Using an analogy with the theory of a linear panel regression, propose a multistep procedure of estimating  $\alpha$  and  $\beta$  adapting the estimator you used in Problem 9.3 to a panel data environment.

## 13.7 Durbin–Watson statistic and panel data

<sup>1</sup>Consider the standard one-way error component model with random effects:

$$y_{it} = x'_{it}\beta + \mu_i + v_{it}, \quad i = 1, \dots, n, \quad t = 1, \dots, T, \quad (13.1)$$

where  $\beta$  is  $k \times 1$ ,  $\mu_i$  are random individual effects,  $\mu_i \sim IID(0, \sigma_\mu^2)$ ,  $v_{it}$  are idiosyncratic shocks,  $v_{it} \sim IID(0, \sigma_v^2)$ , and  $\mu_i$  and  $v_{it}$  are independent of  $x_{it}$  for all  $i$  and  $t$  and mutually. The equations are arranged so that the index  $t$  is faster than the index  $i$ . Consider running OLS on the original regression (13.1) and running OLS on the GLS-transformed regression

$$y_{it} - \hat{\pi}\bar{y}_i = (x'_{it} - \hat{\pi}\bar{x}_i)\beta + (1 - \hat{\pi})\mu_i + v_{it} - \hat{\pi}\bar{v}_i, \quad i = 1, \dots, n, \quad t = 1, \dots, T, \quad (13.2)$$

where  $\hat{\pi}$  is a consistent (as  $n \rightarrow \infty$  and  $T$  stays fixed) estimate of  $\pi = 1 - \sigma_v / \sqrt{\sigma_v^2 + T\sigma_\mu^2}$ . When each OLS estimate is obtained using a typical regression package, the Durbin–Watson (DW) statistic is provided among the regression output. Recall that if  $\hat{e}_1, \hat{e}_2, \dots, \hat{e}_{N-1}, \hat{e}_N$  is a series of regression residuals, then the DW statistic is

$$\mathcal{DW} = \frac{\sum_{j=2}^N (\hat{e}_j - \hat{e}_{j-1})^2}{\sum_{j=1}^N \hat{e}_j^2}.$$

1. Derive the probability limits of the two DW statistics, as  $n \rightarrow \infty$  and  $T$  stays fixed.
2. Using the obtained result, propose an asymptotic test for individual effects based on the DW statistic [Hint: That the errors are estimated does not affect the asymptotic distribution of the DW statistic. Take this for granted.]

<sup>1</sup>This problem is a part of S. Anatolyev (2002, 2003) Durbin–Watson statistic and random individual effects. *Econometric Theory* 18, Problem 02.5.1, 1273–1274, and 19, Solution 02.5.2, 882–883.

## 13.8 Higher-order dynamic panel

Formulate a linear dynamic panel regression with a single weakly exogenous regressor, and AR(2) feedback in place of AR(1) feedback (i.e. when two most recent lags of the left side variable are present at the right side). Describe the algorithm of estimation of this model.



# 14. NONPARAMETRIC ESTIMATION

## 14.1 Nonparametric regression with discrete regressor

Let  $(x_i, y_i)$ ,  $i = 1, \dots, n$  be a random sample from the population of  $(x, y)$ , where  $x$  has a discrete distribution with the support  $a_{(1)}, \dots, a_{(k)}$ , where  $a_{(1)} < \dots < a_{(k)}$ . Having written the conditional expectation  $\mathbb{E}[y|x = a_{(j)}]$  in the form that allows to apply the analogy principle, propose an analog estimator  $\hat{g}_j$  of  $g_j = \mathbb{E}[y|x = a_{(j)}]$  and derive its asymptotic distribution.

## 14.2 Nonparametric density estimation

Suppose we have a random sample  $\{x_i\}_{i=1}^n$  and let

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}[x_i \leq x]$$

denote the empirical distribution function of  $x_i$ , where  $\mathbb{I}(\cdot)$  is an indicator function. Consider two density estimators:

◦ one-sided estimator:

$$\hat{f}_1(x) = \frac{\hat{F}(x+h) - \hat{F}(x)}{h}$$

◦ two-sided estimator:

$$\hat{f}_2(x) = \frac{\hat{F}(x+h/2) - \hat{F}(x-h/2)}{h}$$

Show that:

- $\hat{F}(x)$  is an unbiased estimator of  $F(x)$ . Hint: recall that  $F(x) = \mathbb{P}\{x_i \leq x\} = \mathbb{E}[\mathbb{I}[x_i \leq x]]$ .
- The bias of  $\hat{f}_1(x)$  is  $O(h^a)$ . Find the value of  $a$ . Hint: take a second-order Taylor series expansion of  $F(x+h)$  around  $x$ .
- The bias of  $\hat{f}_2(x)$  is  $O(h^b)$ . Find the value of  $b$ . Hint: take a second-order Taylor series expansion of  $F(x + \frac{h}{2})$  and  $F(x - \frac{h}{2})$  around  $x$ .

## 14.3 Nadaraya–Watson density estimator

Derive the asymptotic distribution of the Nadaraya–Watson estimator of the density of a scalar random variable  $x$  having a continuous distribution, similarly to how the asymptotic distribution of the Nadaraya–Watson estimator of the regression function is derived, under similar conditions. Give interpretation to how your expressions for asymptotic bias and asymptotic variance depend on the shape of the density.

## 14.4 First difference transformation and nonparametric regression

This problem illustrates the use of a difference operator in nonparametric estimation with IID data. Suppose that there is a scalar variable  $z$  that takes values on a bounded support. For simplicity, let  $z$  be deterministic and compose a uniform grid on the unit interval  $[0, 1]$ . The other variables are IID. Assume that for the function  $g(\cdot)$  below the following *Lipschitz condition* is satisfied:

$$|g(u) - g(v)| \leq G|u - v|$$

for some constant  $G$ .

1. Consider a nonparametric regression of  $y$  on  $z$ :

$$y_i = g(z_i) + e_i, \quad i = 1, \dots, n, \quad (14.1)$$

where  $\mathbb{E}[e_i|z_i] = 0$ . Let the data  $\{(z_i, y_i)\}_{i=1}^n$  be ordered so that the  $z$ 's are in increasing order. A *first difference transformation* results in the following set of equations:

$$y_i - y_{i-1} = g(z_i) - g(z_{i-1}) + e_i - e_{i-1}, \quad i = 2, \dots, n. \quad (14.2)$$

The target is to estimate  $\sigma^2 \equiv \mathbb{E}[e_i^2]$ . Propose its consistent estimator based on the FD-transformed regression (14.2). Prove consistency of your estimator.

2. Consider the following partially linear regression of  $y$  on  $x$  and  $z$ :

$$y_i = x_i' \beta + g(z_i) + e_i, \quad i = 1, \dots, n, \quad (14.3)$$

where  $\mathbb{E}[e_i|x_i, z_i] = 0$ . Let the data  $\{(x_i, z_i, y_i)\}_{i=1}^n$  be ordered so that the  $z$ 's are in increasing order. The target is to nonparametrically estimate  $g$ . Propose its consistent estimator using the FD-transformation of (14.3). [Hint: on the first step, consistently estimate  $\beta$  from the FD-transformed regression.] Prove consistency of your estimator.

## 14.5 Unbiasedness of kernel estimates

Recall the Nadaraya–Watson kernel estimator  $\hat{g}(x)$  of the conditional mean  $g(x) \equiv \mathbb{E}[y|x]$  constructed for a random sample. Show that if  $g(x) = c$ , where  $c$  is some constant, then  $\hat{g}(x)$  is unbiased, and provide intuition behind this result. Find out under what circumstance will the local linear estimator of  $g(x)$  be unbiased under random sampling. Finally, investigate the kernel estimator of the density  $f(x)$  of  $x$  for unbiasedness under random sampling.

## 14.6 Shape restriction

Firms produce some product using technology  $f(l, k)$ . The functional form of  $f$  is unknown, although we know that it exhibits constant returns to scale. For a firm  $i$ , we observe labor  $l_i$ , capital  $k_i$ , and output  $y_i$ , and the data generating process takes the form  $y_i = f(l_i, k_i) + \varepsilon_i$ , where  $\mathbb{E}[\varepsilon_i] = 0$  and  $\varepsilon_i$  is independent of  $(l_i, k_i)$ . Using random sample  $\{y_i, l_i, k_i\}_{i=1}^n$ , suggest a nonparametric estimator of  $f(l, k)$  which also exhibits constant returns to scale.

## 14.7 Nonparametric hazard rate

Let  $z_1, \dots, z_n$  be scalar IID random variables with unknown PDF  $f(\cdot)$  and CDF  $F(\cdot)$ . Assume that the distribution of  $z$  has support  $\mathbb{R}$ . Pick  $t \in \mathbb{R}$  such that  $0 < F(t) < 1$ . The objective is to estimate the *hazard rate*  $H(t) = f(t)/(1 - F(t))$ .

(i) Suggest a nonparametric estimator for  $F(t)$ . Denote this estimator by  $\hat{F}(t)$ .

(ii) Let

$$\hat{f}(t) = \frac{1}{nh_n} \sum_{j=1}^n k\left(\frac{z_j - t}{h_n}\right)$$

denote the Nadaraya–Watson estimate of  $f(t)$  where the bandwidth  $h_n$  is chosen so that  $nh_n^5 \rightarrow 0$ , and  $k(\cdot)$  is a symmetric kernel. Find the asymptotic distribution of  $\hat{f}(t)$ . Do not worry about regularity conditions.

(iii) Use  $\hat{f}(t)$  and  $\hat{F}(t)$  to suggest an estimator for  $H(t)$ . Denote this estimator by  $\hat{H}(t)$ . Find the asymptotic distribution of  $\hat{H}(t)$ .

## 14.8 Nonparametrics and perfect fit

Analyze carefully the asymptotic properties of the Nadaraya–Watson estimator of a regression function with perfect fit, i.e. when the variance of the error is zero.

## 14.9 Nonparametrics and extreme observations

Discuss the behavior of a Nadaraya–Watson mean regression estimate when one of IID observations, say  $(x_1, y_1)$ , tends to assume extreme values. Specifically, discuss

(a) the case when  $x_1$  is very big,

(b) the case when  $y_1$  is very big.



# 15. CONDITIONAL MOMENT RESTRICTIONS

## 15.1 Usefulness of skedastic function

Suppose that for the following linear regression model

$$y = x'\beta + e, \quad \mathbb{E}[e|x] = 0$$

the form of a skedastic function is

$$\mathbb{E}[e^2|x] = h(x, \beta, \pi),$$

where  $h(\cdot)$  is a known smooth function, and  $\pi$  is an additional parameter vector. Compare asymptotic variances of optimal GMM estimators of  $\beta$  when only the first restriction or both restrictions are employed. Under what conditions does including the second restriction into a set of moment restrictions reduce asymptotic variance? What if the function  $h(\cdot)$  does not depend on  $\beta$ ? What if in addition the distribution of  $e$  conditional on  $x$  is symmetric?

## 15.2 Symmetric regression error

Consider the regression

$$y = \alpha x + e, \quad \mathbb{E}[e|x] = 0,$$

where all variables are scalars. The random sample  $\{y_i, x_i\}_{i=1}^n$  is available.

1. The researcher also suspects that  $y$ , conditional on  $x$ , is distributed symmetrically around the conditional mean. Devise a Hausman specification test for this symmetry. Be specific and give all details at all stages when constructing the test.
2. Suppose that even though the Hausman test rejects symmetry, the researcher uses the assumption that  $e|x \sim \mathcal{N}(0, \sigma^2)$ . Derive the asymptotic properties of the QML estimator of  $\alpha$ .

## 15.3 Optimal instrumentation of consumption function

Consider the model

$$c_t = \alpha + \delta y_t^\gamma + e_t,$$

where  $c_t$  is consumption at  $t$ ,  $y_t$  is income at  $t$ , and all variables are jointly stationary. There is endogeneity in income, however, so that  $e_t$  is not mean independent of  $y_t$ . However, the lagged values of income are predetermined, so  $e_t$  is mean independent of the past history of  $y_t$  and  $c_t$ :

$$\mathbb{E}[e_t | y_{t-1}, c_{t-1}, y_{t-2}, c_{t-2}, \dots] = 0.$$

Suppose a long time series on  $c_t$  and  $y_t$  is available. Outline how you would estimate the model parameters most efficiently.

## 15.4 Optimal instrument in AR-ARCH model

Consider an  $AR(1) - ARCH(1)$  model:  $x_t = \rho x_{t-1} + \varepsilon_t$  where the distribution of  $\varepsilon_t$  conditional on  $I_{t-1}$  is symmetric around 0 with  $\mathbb{E}[\varepsilon_t^2 | I_{t-1}] = (1 - \alpha) + \alpha \varepsilon_{t-1}^2$ , where  $0 < \rho, \alpha < 1$  and  $I_t = \{x_t, x_{t-1}, \dots\}$ .

1. Let the space of admissible instruments for estimation of the  $AR(1)$  part be

$$\mathcal{Z}_t = \left\{ \sum_{i=1}^{\infty} \phi_i x_{t-i}, \text{ s.t. } \sum_{i=1}^{\infty} \phi_i^2 < \infty \right\}.$$

Using the optimality condition, find the optimal instrument as a function of the model parameters  $\rho$  and  $\alpha$ . Outline how to construct its feasible version.

2. Use your intuition to speculate on relative efficiency of the optimal instrument you found in part 1 versus the optimal instrument based on the conditional moment restriction  $\mathbb{E}[\varepsilon_t | I_{t-1}] = 0$ .

## 15.5 Optimal instrument in AR with nonlinear error

Consider an  $AR(1)$  model  $x_t = \rho x_{t-1} + \varepsilon_t$ , where the disturbance  $\varepsilon_t$  is generated as  $\varepsilon_t = \eta_t \eta_{t-1}$ , where  $\eta_t$  is an IID sequence with mean zero and variance one.

1. Show that  $\mathbb{E}[\varepsilon_t x_{t-j}] = 0$  for all  $j \geq 1$ . Find the optimal instrument based on this system of unconditional moment restrictions. How many lags of  $x_t$  does it employ? Outline how to construct its feasible version.
2. Show that  $\mathbb{E}[\varepsilon_t | I_{t-1}] = 0$ , where  $I_t = \{x_t, x_{t-1}, \dots\}$ . Find the optimal instrument based on this conditional moment restriction. Outline how to construct its feasible version, or, if that is impossible, explain why.

## 15.6 Optimal IV estimation of a constant

Consider the following  $MA(p)$  data generating mechanism:

$$y_t = \alpha + \Theta(L)\varepsilon_t,$$

where  $\varepsilon_t$  is a mean zero IID sequence, and  $\Theta(L)$  is lag polynomial of finite order  $p$ . Derive the optimal instrument for estimation of  $\alpha$  based on the conditional moment restriction

$$\mathbb{E}[y_t | y_{t-p-1}, y_{t-p-2}, \dots] = \alpha.$$

## 15.7 Negative binomial distribution and PML

Determine if using the negative binomial distribution having the density

$$f(u, m) = \frac{\Gamma(a+u)}{\Gamma(a)\Gamma(1+u)} \left(\frac{m}{a}\right)^u \left(1 + \frac{m}{a}\right)^{-(a+u)},$$

where  $m$  is its mean and  $a$  is an arbitrary known constant, leads to consistent estimation of  $\theta_0$  in the mean regression

$$\mathbb{E}[y|x] = m(x, \theta_0)$$

when the true conditional distribution is heteroskedastic normal, with the skedastic function

$$\mathbb{V}[y|x] = 2m(x, \theta_0).$$

## 15.8 Nesting and PML

Consider the regression model  $\mathbb{E}[y|x] = m(x, \theta_0)$ . Suppose that a PML1 estimator based on the density  $f(z, \mu)$  parameterized by the mean  $\mu$  consistently estimates the true parameter  $\theta_0$ . Consider another density  $h(z, \mu, \varsigma)$  parameterized by two parameters, the mean  $\mu$  and some other parameter  $\varsigma$ , which nests  $f(z, \mu)$  (i.e.,  $f$  is a special case of  $h$ ). Use the example of the Weibull distribution having the density

$$h(z, \mu, \varsigma) = \varsigma \left(\frac{\Gamma(1+\varsigma^{-1})}{\mu}\right)^\varsigma z^{\varsigma-1} \exp\left(-\left(\frac{\Gamma(1+\varsigma^{-1})}{\mu}z\right)^\varsigma\right) \cdot \mathbb{I}[z \geq 0], \quad \varsigma > 0,$$

to show that the PML1 estimator based on  $h(z, \mu, \varsigma)$  does not necessarily consistently estimates  $\theta_0$ . What is the econometric explanation of this perhaps counter-intuitive result?

Hint: you may use the information that  $\mathbb{E}[z^2] = \mu^2 \Gamma(2+\varsigma^{-1}) / \Gamma(1+\varsigma^{-1})^2$ ,  $\Gamma(x) = (x-1)\Gamma(x-1)$ , and  $\Gamma(1) = 1$ .

## 15.9 Misspecification in variance

Consider the regression model  $\mathbb{E}[y|x] = m(x, \theta_0)$ . Suppose that this regression is conditionally normal and homoskedastic. A researcher, however, uses the following conditional density to construct a PML1 estimator of  $\theta_0$ :

$$(y|x, \theta) \sim \mathcal{N}\left(m(x, \theta), m(x, \theta)^2\right).$$

Establish if such estimator is consistent for  $\theta_0$ .

## 15.10 Modified Poisson regression and PML estimators

<sup>1</sup>Let the observable random variable  $y$  be distributed, conditionally on observable  $x$  and unobservable  $\varepsilon$  as Poisson with the parameter  $\lambda(x) = \exp(x'\beta + \varepsilon)$ , where  $\mathbb{E}[\exp \varepsilon | x] = 1$  and  $\mathbb{V}[\exp \varepsilon | x] = \sigma^2$ . Suppose that vector  $x$  is distributed as multivariate standard normal.

1. Find the regression and skedastic functions, where the conditional information involves only  $x$ .
2. Find the asymptotic variances of the Nonlinear Least Squares (NLLS) and Weighted Nonlinear Least Squares (WNLLS) estimators of  $\beta$ .
3. Find the asymptotic variances of the Pseudo-Maximum Likelihood (PML) estimators of  $\beta$  based on
  - (a) the normal distribution;
  - (b) the Poisson distribution;
  - (c) the Gamma distribution.
4. Rank the five estimators in terms of asymptotic efficiency.

## 15.11 Optimal instrument and regression on constant

Consider the following model:

$$y_i = \alpha + e_i, \quad i = 1, \dots, n,$$

where unobservable  $e_i$  conditionally on  $x_i$  is distributed *symmetrically* with mean zero and variance  $x_i^2 \sigma^2$  with unknown  $\sigma^2$ . The data  $(y_i, x_i)$  are IID.

1. Construct a pair of conditional moment restrictions from the information about the conditional mean and conditional variance. Derive the optimal unconditional moment restrictions, corresponding to (a) the conditional restriction associated with the conditional mean; (b) the conditional restrictions associated with both the conditional mean and conditional variance.
2. Describe in detail the GMM estimators that correspond to the two optimal sets of unconditional moment restrictions of part 1. Note that in part 1(a) the parameter  $\sigma^2$  is not identified, therefore propose your own estimator of  $\sigma^2$  that differs from the one implied by part 1(b). All estimators that you construct should be fully feasible. If you use nonparametric estimation, give all the details. Your description should also contain estimation of asymptotic variances.
3. Compare the asymptotic properties of the GMM estimators that you designed.
4. Derive the Pseudo-Maximum Likelihood estimator of  $\alpha$  and  $\sigma^2$  of order 2 (PML2) that is based on the normal distribution. Derive its asymptotic properties. How does this estimator relate to the GMM estimators you obtained in part 2?

---

<sup>1</sup>The idea of this problem is borrowed from Gourieroux, C. and Monfort, A. "Statistics and Econometric Models", Cambridge University Press, 1995.



# 16. EMPIRICAL LIKELIHOOD

## 16.1 Common mean

Suppose we have the following moment restrictions:  $\mathbb{E}[x] = \mathbb{E}[y] = \theta$ .

1. Find the system of equations that yields the empirical likelihood (EL) estimator  $\hat{\theta}$  of  $\theta$ , the associated Lagrange multipliers  $\hat{\lambda}$  and the implied probabilities  $\hat{p}_i$ . Derive the asymptotic variances of  $\hat{\theta}$  and  $\hat{\lambda}$  and show how to estimate them.
2. Reduce the number of parameters by eliminating the redundant ones. Then linearize the system of equations with respect to the Lagrange multipliers that are left, around their population counterparts of zero. This will help to find an approximate, but explicit solution for  $\hat{\theta}$ ,  $\hat{\lambda}$  and  $\hat{p}_i$ . Derive that solution and interpret it.
3. Instead of defining the objective function

$$\frac{1}{n} \sum_{i=1}^n \log p_i$$

as in the EL approach, let the objective function be

$$-\frac{1}{n} \sum_{i=1}^n p_i \log p_i.$$

This gives rise to the *exponential tilting* (ET) estimator of  $\theta$ . Find the system of equations that yields the ET estimator of  $\hat{\theta}$ , the associated Lagrange multipliers  $\hat{\lambda}$  and the implied probabilities  $\hat{p}_i$ . Derive the asymptotic variances of  $\hat{\theta}$  and  $\hat{\lambda}$  and show how to estimate them.

## 16.2 Kullback–Leibler Information Criterion

The Kullback–Leibler Information Criterion (KLIC) measures the distance between distributions, say  $g(z)$  and  $h(z)$ :

$$\mathcal{KLIC}(g : h) = \mathbb{E}_g \left[ \log \frac{g(z)}{h(z)} \right],$$

where  $\mathbb{E}_g[\cdot]$  denotes mathematical expectation according to  $g(z)$ .

Suppose we have the following moment condition:

$$\mathbb{E} \left[ m(z_i, \theta_0) \right] = \mathbf{0}_{\ell \times 1}, \quad \ell \geq k,$$

and a random sample  $z_1, \dots, z_n$  with no elements equal to each other. Denote by  $e$  the empirical distribution function (EDF), i.e. the one that assigns probability  $\frac{1}{n}$  to each sample point. Denote by  $\pi$  a discrete distribution that assigns probability  $\pi_i$  to the sample point  $z_i$ ,  $i = 1, \dots, n$ .

1. Show that minimization of  $\mathcal{K}\mathcal{L}\mathcal{I}\mathcal{C}(e : \pi)$  subject to  $\sum_{i=1}^n \pi_i = 1$  and  $\sum_{i=1}^n \pi_i m(z_i, \theta) = 0$  yields the Empirical Likelihood (EL) value of  $\theta$  and corresponding implied probabilities.
2. Now we switch the roles of  $e$  and  $\pi$  and consider minimization of  $\mathcal{K}\mathcal{L}\mathcal{I}\mathcal{C}(\pi : e)$  subject to the same constraints. What familiar estimator emerges as the solution to this optimization problem?
3. Now suppose that we have *a priori* knowledge about the distribution of the data. So, instead of using the EDF, we use the distribution  $p$  that assigns known probability  $p_i$  to the sample point  $z_i$ ,  $i = 1, \dots, n$  (of course,  $\sum_{i=1}^n p_i = 1$ ). Analyze how the solutions to the optimization problems in parts 1 and 2 change.
4. Now suppose that we have postulated a family of densities  $f(z, \theta)$  which is compatible with the moment condition. Interpret the value of  $\theta$  that minimizes  $\mathcal{K}\mathcal{L}\mathcal{I}\mathcal{C}(e : f)$ .

## 16.3 Empirical likelihood as IV estimation

Consider a linear model with instrumental variables:

$$y = x'\beta + e, \quad \mathbb{E}[ze] = 0,$$

where  $x$  is  $k \times 1$ ,  $z$  is  $\ell \times 1$ , and  $\ell \geq k$ . Write down the EL estimator of  $\beta$  in a matrix form of a (not completely feasible) instrumental variables estimator. Also write down the efficient GMM estimator, and explain intuitively why the former is expected to exhibit better finite sample properties than the latter.

# 17. ADVANCED ASYMPTOTIC THEORY

---

## 17.1 Maximum likelihood and asymptotic bias

Derive the second order bias of the Maximim Likelihood (ML) estimator  $\hat{\lambda}$  of the parameter  $\lambda > 0$  of the exponential distribution

$$f(y, \lambda) = \begin{cases} \lambda \exp(-\lambda y), & y \geq 0 \\ 0, & y < 0 \end{cases}$$

obtained from random sample  $y_1, \dots, y_T$ ,

- (a) using an explicit formula for  $\hat{\lambda}$ ;
- (b) using the expression for the second order bias of extremum estimators.

Construct the bias corrected ML estimator of  $\lambda$ .

## 17.2 Empirical likelihood and asymptotic bias

Consider estimation of a scalar parameter  $\theta$  on the basis of the moment function

$$m(x, y, \theta) = \begin{pmatrix} x - \theta \\ y - \theta \end{pmatrix}$$

and IID data  $(x_i, y_i)$ ,  $i = 1, \dots, n$ . Show that the second order asymptotic bias of the empirical likelihood estimator of  $\theta$  equals 0.

## 17.3 Asymptotically irrelevant instruments

Consider the linear model

$$y = \beta x + e,$$

where scalar random variables  $x$  and  $e$  are correlated with the correlation coefficient  $\rho$ . Available are data for an  $\ell \times 1$  vector of instruments  $z$ . These instruments, however, are asymptotically irrelevant, i.e.  $\mathbb{E}[zx] = 0$ . The data  $(x_i, y_i, z_i)$ ,  $i = 1, \dots, n$ , are IID. You may additionally assume that both  $x$  and  $e$  are homoskedastic conditional on  $z$ , and that they are homocorrelated conditional on  $z$ .

1. Find the probability limit of the 2SLS estimator of  $\beta$  from the first principles (i.e. without using the weak instruments theory).

2. Verify that your result in part 1 conforms to the weak instruments theory being its special case.
3. Find the expected value of the probability limit of the 2SLS estimator. How does it relate to the probability limit of the OLS estimator?

## 17.4 Weakly endogenous regressors

Consider the regression

$$\mathcal{Y} = \mathcal{X}\beta + \mathcal{E},$$

where the regressors  $\mathcal{X}$  are correlated with the error  $\mathcal{E}$ , but this correlation is weak. Consider the decomposition of  $\mathcal{E}$  to its projection on  $\mathcal{X}$  and the orthogonal component  $\mathcal{U}$ :

$$\mathcal{E} = \mathcal{X}\pi + \mathcal{U}.$$

Assume that  $(n^{-1}\mathcal{X}'\mathcal{X}, n^{-1/2}\mathcal{X}'\mathcal{U}) \xrightarrow{p} (Q, \xi)$ , where  $\xi \sim \mathcal{N}(0, \sigma_u^2 Q)$  and  $Q$  has full rank. Show that under the assumption of the drifting parameter DGP  $\pi = c/\sqrt{n}$ , where  $n$  is the sample size and  $c$  is fixed, the OLS estimator of  $\beta$  is consistent and asymptotically noncentral normal, and derive the asymptotic distribution of the Wald test statistic for testing the set of linear restriction  $R\beta = r$ , where  $R$  has full rank  $q$ .

## 17.5 Weakly invalid instruments

Consider a linear model with IID data

$$y = \beta x + e,$$

where all variables are scalars.

1. Suppose that  $x$  and  $e$  are correlated, but there is an  $\ell \times 1$  strong “instrument”  $z$  weakly correlated with  $e$ . Derive the asymptotic (as  $n \rightarrow \infty$ ) distributions of the 2SLS estimator of  $\beta$ , its  $t$  ratio, and the overidentification test statistic

$$\mathcal{J} = n \frac{\widehat{\mathcal{U}}' \mathcal{Z} (\mathcal{Z}' \mathcal{Z})^{-1} \mathcal{Z}' \widehat{\mathcal{U}}}{\widehat{\mathcal{U}}' \widehat{\mathcal{U}}},$$

where  $\widehat{\mathcal{U}} \equiv \mathcal{Y} - \widehat{\beta}\mathcal{X}$  are the vector of 2SLS residuals and  $\mathcal{Z}$  is the matrix of instruments, under the drifting DGP  $\omega = c_\omega/\sqrt{n}$ , where  $\omega$  is the vector of coefficients on  $z$  in the linear projection of  $e$  on  $z$ . Also, specialize to the case  $\ell = 1$ .

2. Suppose that  $x$  and  $e$  are correlated, but there is an  $\ell \times 1$  weak “instrument”  $z$  weakly correlated with  $e$ . Derive the asymptotic (as  $n \rightarrow \infty$ ) distributions of the 2SLS estimator of  $\beta$ , its  $t$  ratio, and the overidentification test statistic  $\mathcal{J}$ , under the drifting DGP  $\omega = c_\omega/\sqrt{n}$  and  $\pi = c_\pi/\sqrt{n}$ , where  $\omega$  is the vector of coefficients on  $z$  in the linear projection of  $e$  on  $z$ , and  $\pi$  is the vector of coefficients on  $z$  in the linear projection of  $x$  on  $z$ . Also, specialize to the case  $\ell = 1$ .

**Part II**  
**Solutions**



# 1. ASYMPTOTIC THEORY: GENERAL AND INDEPENDENT DATA

---

## 1.1 Asymptotics of transformations

1. There are at least two ways to solve the problem. An easier way is using the “second-order Delta-method”, see Problem 1.8. The second way is using trigonometric identities and the regular Delta-method:

$$T(1 - \cos \hat{\phi}) = 2T \sin^2 \frac{\hat{\phi}}{2} = 2 \left( \sqrt{T} \sin \frac{\hat{\phi}}{2} \right)^2 \xrightarrow{d} 2 \left( \frac{\cos \pi}{2} \cdot \mathcal{N}(0, 1) \right)^2 \sim \frac{1}{2} \chi_1^2.$$

2. Recalling how the Delta-method is proved,

$$T \sin \hat{\psi} = T(\sin \hat{\psi} - \sin 2\pi) = T \left. \frac{\partial \sin \psi}{\partial \psi} \right|_{\psi \rightarrow 2\pi} (\hat{\psi} - 2\pi) \xrightarrow{d} \mathcal{N}(0, 1).$$

3. By the Mann–Wald theorem,

$$\log T + \log \hat{\theta} \xrightarrow{d} \log \chi_1^2 \Rightarrow \log \hat{\theta} \xrightarrow{d} -\infty \Rightarrow T \log \hat{\theta} \xrightarrow{d} -\infty.$$

## 1.2 Asymptotics of rotated logarithms

Use the Delta-method for

$$\sqrt{n} \left( \begin{pmatrix} U_n \\ V_n \end{pmatrix} - \begin{pmatrix} \mu_u \\ \mu_v \end{pmatrix} \right) \xrightarrow{d} \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma \right)$$

and

$$g \left( \begin{pmatrix} x \\ y \end{pmatrix} \right) = \begin{pmatrix} \ln x - \ln y \\ \ln x + \ln y \end{pmatrix}.$$

We have

$$\frac{\partial g}{\partial(x \ y)} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 1/x & -1/y \\ 1/x & 1/y \end{pmatrix}, \quad G = \frac{\partial g}{\partial(x \ y)} \begin{pmatrix} \mu_u \\ \mu_v \end{pmatrix} = \begin{pmatrix} 1/\mu_u & -1/\mu_v \\ 1/\mu_u & 1/\mu_v \end{pmatrix},$$

so

$$\sqrt{n} \left( \begin{pmatrix} \ln U_n - \ln V_n \\ \ln U_n + \ln V_n \end{pmatrix} - \begin{pmatrix} \ln \mu_u - \ln \mu_v \\ \ln \mu_u + \ln \mu_v \end{pmatrix} \right) \xrightarrow{d} \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, G\Sigma G' \right),$$

where

$$G\Sigma G' = \begin{pmatrix} \frac{\omega_{uu}}{\mu_u^2} - \frac{2\omega_{uv}}{\mu_u\mu_v} + \frac{\omega_{vv}}{\mu_v^2} & \frac{\omega_{uu}}{\mu_u^2} - \frac{\omega_{vv}}{\mu_v^2} \\ \frac{\omega_{uu}}{\mu_u^2} - \frac{\omega_{vv}}{\mu_v^2} & \frac{\omega_{uu}}{\mu_u^2} + \frac{2\omega_{uv}}{\mu_u\mu_v} + \frac{\omega_{vv}}{\mu_v^2} \end{pmatrix}.$$

It follows that  $\ln U_n - \ln V_n$  and  $\ln U_n + \ln V_n$  are asymptotically independent when  $\frac{\omega_{uu}}{\mu_u^2} = \frac{\omega_{vv}}{\mu_v^2}$ .

## 1.3 Escaping probability mass

(i) The expectation is

$$\mathbb{E}[\hat{\mu}_n] = \mathbb{E}[\bar{x}_n|A_n] \mathbb{P}\{A_n\} + \mathbb{E}[n|\bar{A}_n] \mathbb{P}\{\bar{A}_n\} = \mu \left(1 - \frac{1}{n}\right) + 1,$$

so the bias is

$$\mathbb{E}[\hat{\mu}_n] - \mu = -\frac{\mu}{n} + 1 \rightarrow 1$$

as  $n \rightarrow \infty$ . The expectation of the square is

$$\mathbb{E}[\hat{\mu}_n^2] = \mathbb{E}[\bar{x}_n^2|A_n] \mathbb{P}\{A_n\} + \mathbb{E}[n^2|\bar{A}_n] \mathbb{P}\{\bar{A}_n\} = \left(\mu^2 + \frac{\sigma^2}{n}\right) \left(1 - \frac{1}{n}\right) + n,$$

so the variance is

$$\mathbb{V}[\hat{\mu}_n] = \mathbb{E}[\hat{\mu}_n^2] - (\mathbb{E}[\hat{\mu}_n])^2 = \frac{1}{n} \left(1 - \frac{1}{n}\right) \left((n - \mu)^2 + \sigma^2\right) \rightarrow \infty$$

as  $n \rightarrow \infty$ . As a consequence, the MSE of  $\hat{\mu}_n$  is

$$\text{MSE}[\hat{\mu}_n] = \mathbb{V}[\hat{\mu}_n] + (\mathbb{E}[\hat{\mu}_n] - \mu)^2 = \frac{1}{n} \left((n - \mu)^2 + \left(1 - \frac{1}{n}\right) \sigma^2\right),$$

and also tends to infinity as  $n \rightarrow \infty$ .

(ii) Despite the MSE of  $\hat{\mu}_n$  goes to infinity,  $\hat{\mu}_n$  is consistent: for any  $\varepsilon > 0$ ,

$$\mathbb{P}\{|\hat{\mu}_n - \mu| > \varepsilon\} = \mathbb{P}\{|\bar{x}_n - \mu| > \varepsilon\} \left(1 - \frac{1}{n}\right) + \mathbb{P}\{|n - \mu| > \varepsilon\} \frac{1}{n} \rightarrow 0$$

by consistency of  $\bar{x}_n$  and boundedness of probabilities. The CDF of  $\sqrt{n}(\hat{\mu}_n - \mu)$  is

$$\begin{aligned} F_{\sqrt{n}(\hat{\mu}_n - \mu)}(t) &= \mathbb{P}\{\sqrt{n}(\hat{\mu}_n - \mu) \leq t\} \\ &= \mathbb{P}\{\sqrt{n}(\bar{x}_n - \mu) \leq t\} \left(1 - \frac{1}{n}\right) + \mathbb{P}\{\sqrt{n}(n - \mu) \leq t\} \frac{1}{n} \\ &\rightarrow F_{\mathcal{N}(0, \sigma^2)}(t) \end{aligned}$$

by the asymptotic normality of  $\bar{x}_n$  and boundedness of probabilities.

(iii) Since the asymptotic distributions of  $\bar{x}_n$  and  $\hat{\mu}_n$  are the same, the approximate confidence intervals for  $\mu$  will be identical except that they center at  $\bar{x}_n$  and  $\hat{\mu}_n$ , respectively.

## 1.4 Asymptotics of $t$ -ratios

The solution is straightforward once we determine to which vector the LLN or CLT should be applied.

(a) When  $\mu = 0$ , we have:  $\bar{x} \xrightarrow{p} 0$ ,  $\sqrt{n}\bar{x} \xrightarrow{d} \mathcal{N}(0, \sigma^2)$ , and  $\hat{\sigma}^2 \xrightarrow{p} \sigma^2$ . Therefore

$$\sqrt{n}T_n = \frac{\sqrt{n}\bar{x}}{\hat{\sigma}} \xrightarrow{d} \frac{1}{\sigma} \mathcal{N}(0, \sigma^2) \sim \mathcal{N}(0, 1).$$



(b) Consider the vector

$$W_n \equiv \begin{pmatrix} \bar{x} \\ \hat{\sigma}^2 \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} x_i \\ (x_i - \mu)^2 \end{pmatrix} - \begin{pmatrix} 0 \\ (\bar{x} - \mu)^2 \end{pmatrix}.$$

By the LLN, the last term goes in probability to the zero vector, while the first term, and thus the whole  $W_n$ , converges in probability to

$$\text{plim}_{n \rightarrow \infty} W_n = \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix}.$$

Moreover, because  $\sqrt{n}(\bar{x} - \mu) \xrightarrow{d} \mathcal{N}(0, \sigma^2)$ , we have  $\sqrt{n}(\bar{x} - \mu)^2 \xrightarrow{d} 0$ .

Next, let  $W_i \equiv (x_i, (x_i - \mu)^2)'$ . Then  $\sqrt{n} \left( W_n - \text{plim}_{n \rightarrow \infty} W_n \right) \xrightarrow{d} \mathcal{N}(0, V)$ , where  $V \equiv \mathbb{V}[W_i]$ .

Let us calculate  $V$ . First,  $\mathbb{V}[x_i] = \sigma^2$  and  $\mathbb{V}[(x_i - \mu)^2] = \mathbb{E}[(x_i - \mu)^2 - \sigma^2]^2 = \tau - \sigma^4$ . Second,  $\mathbb{C}[x_i, (x_i - \mu)^2] = \mathbb{E}[(x_i - \mu)((x_i - \mu)^2 - \sigma^2)] = 0$ . Therefore,

$$\sqrt{n} \left( W_n - \text{plim}_{n \rightarrow \infty} W_n \right) \xrightarrow{d} \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma^2 & 0 \\ 0 & \tau - \sigma^4 \end{pmatrix} \right)$$

Now use the Delta-method with the transformation

$$g \left( \begin{pmatrix} t_1 \\ t_2 \end{pmatrix} \right) = \frac{t_1}{\sqrt{t_2}} \quad \Rightarrow \quad g' \left( \begin{pmatrix} t_1 \\ t_2 \end{pmatrix} \right) = \frac{1}{\sqrt{t_2}} \begin{pmatrix} 1 \\ -\frac{t_1}{2t_2} \end{pmatrix}$$

to get

$$\sqrt{n} \left( T_n - \text{plim}_{n \rightarrow \infty} T_n \right) \xrightarrow{d} \mathcal{N} \left( 0, 1 + \frac{\mu^2(\tau - \sigma^4)}{4\sigma^6} \right).$$

Indeed, this reduces to  $\mathcal{N}(0, 1)$  when  $\mu = 0$ .

(c) Similarly, consider the vector

$$W_n \equiv \begin{pmatrix} \bar{x} \\ \bar{\sigma}^2 \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} x_i \\ x_i^2 \end{pmatrix}.$$

By the LLN,  $W_n$  converges in probability to

$$\text{plim}_{n \rightarrow \infty} W_n = \begin{pmatrix} \mu \\ \mu^2 + \sigma^2 \end{pmatrix}.$$

Next,  $\sqrt{n} \left( W_n - \text{plim}_{n \rightarrow \infty} W_n \right) \xrightarrow{d} \mathcal{N}(0, V)$ , where  $V \equiv \mathbb{V}[W_i]$ ,  $W_i \equiv (x_i, x_i^2)'$ . Let us calculate

$V$ . First,  $\mathbb{V}[x_i] = \sigma^2$  and  $\mathbb{V}[x_i^2] = \mathbb{E}[(x_i^2 - \mu^2 - \sigma^2)^2] = \tau + 4\mu^2\sigma^2 - \sigma^4$ . Second,  $\mathbb{C}[x_i, x_i^2] = \mathbb{E}[(x_i - \mu)(x_i^2 - \mu^2 - \sigma^2)] = 2\mu\sigma^2$ . Therefore,

$$\sqrt{n} \left( W_n - \text{plim}_{n \rightarrow \infty} W_n \right) \xrightarrow{d} \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma^2 & 2\mu\sigma^2 \\ 2\mu\sigma^2 & \tau + 4\mu^2\sigma^2 - \sigma^4 \end{pmatrix} \right).$$

Now use the Delta-method with

$$g \left( \begin{pmatrix} t_1 \\ t_2 \end{pmatrix} \right) = \frac{t_1}{\sqrt{t_2}}$$

to get

$$\sqrt{n} \left( R_n - \text{plim}_{n \rightarrow \infty} R_n \right) \xrightarrow{d} \mathcal{N} \left( 0, \frac{\mu^2\tau - \mu^2\sigma^4 + 4\sigma^6}{4(\mu^2 + \sigma^2)^3} \right).$$

This reduces to the answer in part (b) if and only if  $\mu = 0$ . Under this condition,  $T_n$  and  $R_n$  are asymptotically equivalent.

## 1.5 Creeping bug on simplex

Since  $x_k$  and  $y_k$  are perfectly correlated, it suffices to consider either one,  $x_k$  say. Note that at each step  $x_k$  increases by  $\frac{1}{k}$  with probability  $p$ , or stays the same. That is,  $x_k = x_{k-1} + \frac{1}{k}\xi_k$ , where  $\xi_k$  is IID  $\mathcal{B}(p)$ . This means that  $x_k = \frac{1}{k} \sum_{i=1}^k \xi_i$  which by the LLN converges in probability to  $\mathbb{E}[\xi_i] = p$  as  $k \rightarrow \infty$ . Therefore,  $\text{plim}(x_k, y_k) = (p, 1-p)$ . Next, due to the CLT,

$$\sqrt{n}(x_k - \text{plim } x_k) \xrightarrow{d} \mathcal{N}(0, p(1-p)).$$

Therefore, the rate of convergence is  $\sqrt{n}$ , as usual, and

$$\sqrt{n} \left( \begin{pmatrix} x_k \\ y_k \end{pmatrix} - \text{plim} \begin{pmatrix} x_k \\ y_k \end{pmatrix} \right) \xrightarrow{d} \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} p(1-p) & -p(1-p) \\ -p(1-p) & p(1-p) \end{pmatrix} \right).$$

## 1.6 Asymptotics of sample variance

According to the LLN,  $a$  must be  $\mathbb{E}[x]$ , and  $b$  must be  $\mathbb{E}[x^2]$ . According to the CLT,  $c(n)$  must be  $\sqrt{n}$ , and

$$\sqrt{n} \begin{pmatrix} \bar{x}_n - \mathbb{E}[x] \\ \bar{x}_n^2 - \mathbb{E}[x^2] \end{pmatrix} \xrightarrow{d} \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \mathbb{V}[x] & \mathbb{C}[x, x^2] \\ \mathbb{C}[x, x^2] & \mathbb{V}[x^2] \end{pmatrix} \right).$$

Using the Delta-method with

$$g \left( \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \right) = u_2 - u_1^2$$

whose derivative is

$$G \left( \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \right) = \begin{pmatrix} -2u_1 \\ 1 \end{pmatrix},$$

we get

$$\sqrt{n} \left( \bar{x}_n^2 - (\bar{x}_n)^2 - (\mathbb{E}[x^2] - \mathbb{E}[x]^2) \right) \xrightarrow{d} N \left( 0, \begin{pmatrix} -2\mathbb{E}[x] \\ 1 \end{pmatrix}' \begin{pmatrix} \mathbb{V}[x] & \mathbb{C}[x, x^2] \\ \mathbb{C}[x, x^2] & \mathbb{V}[x^2] \end{pmatrix} \begin{pmatrix} -2\mathbb{E}[x] \\ 1 \end{pmatrix} \right),$$

or

$$\sqrt{n} \left( \bar{x}_n^2 - (\bar{x}_n)^2 - \mathbb{V}[x] \right) \xrightarrow{d} \mathcal{N} \left( 0, \mathbb{V}[x^2] + 4\mathbb{E}[x]^2\mathbb{V}[x] - 4\mathbb{E}[x]\mathbb{C}[x, x^2] \right).$$

## 1.7 Asymptotics of roots

More precisely, we need to assume that  $F$  is continuously differentiable, at least at  $(a, \theta)$ , and that it possesses a continuously differentiable inverse at  $(a, \theta)$ . Suppose that

$$\sqrt{n}(\hat{a} - a) \xrightarrow{d} \mathcal{N}(0, V_{\hat{a}}).$$

The consistency of  $\hat{\theta}$  follows from application of the Mann–Wald theorem. Next, the first-order Taylor expansion of

$$F(\hat{a}, \hat{\theta}) = 0$$

around  $(a, \theta)$  yields

$$F(a, \theta) + \frac{\partial F(a^*, \theta^*)}{\partial a'} (\hat{a} - a) + \frac{\partial F(a^*, \theta^*)}{\partial \theta'} (\hat{\theta} - \theta) = 0,$$

where  $(a^*, \theta^*)$  lies between  $(a, \theta)$  and  $(\hat{a}, \hat{\theta})$  componentwise, and hence is consistent for  $(a, \theta)$ .

Because  $F(a, \theta) = 0$ , we have

$$\frac{\partial F(a^*, \theta^*)}{\partial a'} \sqrt{n} (\hat{a} - a) + \frac{\partial F(a^*, \theta^*)}{\partial \theta'} \sqrt{n} (\hat{\theta} - \theta) = 0,$$

or

$$\begin{aligned} \sqrt{n} (\hat{\theta} - \theta) &= - \left( \frac{\partial F(a^*, \theta^*)}{\partial \theta'} \right)^{-1} \frac{\partial F(a^*, \theta^*)}{\partial a'} \sqrt{n} (\hat{a} - a) \\ &\xrightarrow{d} \mathcal{N} \left( 0, \left( \frac{\partial F(a, \theta)}{\partial \theta'} \right)^{-1} \frac{\partial F(a, \theta)}{\partial a'} V_{\hat{a}} \frac{\partial F(a, \theta)'}{\partial a} \left( \frac{\partial F(a, \theta)'}{\partial \theta} \right)^{-1} \right), \end{aligned}$$

where it is assumed that  $\partial F(a, \theta) / \partial \theta'$  is of full rank. Alternatively, one could use the theorem on differentiation of an implicit function.

## 1.8 Second-order Delta-method

- (a) From the CLT,  $\sqrt{n}S_n \xrightarrow{d} \mathcal{N}(0, 1)$ . Using the Mann–Wald theorem for  $g(x) = x^2$  we get  $nS_n^2 \xrightarrow{d} \chi_1^2$ .
- (b) The Taylor expansion around  $\cos(0) = 1$  yields  $\cos(S_n) = 1 - \frac{1}{2} \cos(S_n^*) S_n^2$ , where  $S_n^* \in [0, S_n]$ . From the LLN and Mann–Wald theorem,  $\cos(S_n^*) \xrightarrow{p} 1$ , and from the Slutsky theorem,  $2n(1 - \cos(S_n)) \xrightarrow{d} \chi_1^2$ .
- (c) Let  $z_n \xrightarrow{p} z = \text{const}$  and  $\sqrt{n}(z_n - z) \xrightarrow{d} \mathcal{N}(0, \sigma^2)$ . Let  $g$  be twice continuously differentiable at  $z$  with  $g'(z) = 0$  and  $g''(z) \neq 0$ . Then

$$\frac{2n g(z_n) - g(z)}{\sigma^2 g''(z)} \xrightarrow{d} \chi_1^2.$$

*Proof.* Indeed, as  $g'(z) = 0$ , from the second-order Taylor expansion,

$$g(z_n) = g(z) + \frac{1}{2} g''(z^*) (z_n - z)^2,$$

and, since  $g''(z^*) \xrightarrow{p} g''(z)$  and  $\frac{\sqrt{n}(z_n - z)}{\sigma} \xrightarrow{d} \mathcal{N}(0, 1)$ , we have

$$\frac{2n g(z_n) - g(z)}{\sigma^2 g''(z)} = \left[ \frac{\sqrt{n}(z_n - z)}{\sigma} \right]^2 \xrightarrow{d} \chi_1^2.$$

*QED*

## 1.9 Asymptotics with shrinking regressor

The OLS estimates are

$$\hat{\beta} = \frac{n^{-1} \sum_i y_i x_i - (n^{-1} \sum_i y_i)(n^{-1} \sum_i x_i)}{n^{-1} \sum_i x_i^2 - (n^{-1} \sum_i x_i)^2}, \quad \hat{\alpha} = \frac{1}{n} \sum_i y_i - \hat{\beta} \cdot \frac{1}{n} \sum_i x_i \quad (1.1)$$

and

$$\hat{\sigma}^2 = \frac{1}{n} \sum_i \hat{e}_i^2.$$

Let us consider  $\hat{\beta}$  first. From (1.1) it follows that

$$\begin{aligned} \hat{\beta} &= \frac{n^{-1} \sum_i (\alpha + \beta x_i + u_i) x_i - n^{-2} \sum_i (\alpha + \beta x_i + u_i) \sum_i x_i}{n^{-1} \sum_i x_i^2 - n^{-2} (\sum_i x_i)^2} \\ &= \beta + \frac{n^{-1} \sum_i \rho^i u_i - n^{-2} \sum_i u_i \sum_i \rho^i}{n^{-1} \sum_i \rho^{2i} - n^{-2} (\sum_i \rho^i)^2} = \beta + \frac{\sum_i \rho^i u_i - \frac{\rho(1-\rho^n)}{1-\rho} (n^{-1} \sum_i u_i)}{\frac{\rho^2(1-\rho^{2n})}{1-\rho^2} - n^{-1} \left( \frac{\rho(1-\rho^n)}{1-\rho} \right)^2}, \end{aligned}$$

which converges to

$$\beta + \frac{1-\rho^2}{\rho^2} \text{plim}_{n \rightarrow \infty} \sum_{i=1}^n \rho^i u_i,$$

if  $\xi \equiv \text{plim} \sum_i \rho^i u_i$  exists and is a well-defined random variable. Its moments are  $\mathbb{E}[\xi] = 0$ ,  $\mathbb{E}[\xi^2] = \sigma^2 \frac{\rho^2}{1-\rho^2}$  and  $\mathbb{E}[\xi^3] = \nu \frac{\rho^3}{1-\rho^3}$ . Hence,

$$\hat{\beta} - \beta \xrightarrow{d} \frac{1-\rho^2}{\rho^2} \xi. \quad (1.2)$$

Now let us look at  $\hat{\alpha}$ . Again, from (1.1) we see that

$$\hat{\alpha} = \alpha + (\beta - \hat{\beta}) \cdot \frac{1}{n} \sum_i \rho^i + \frac{1}{n} \sum_i u_i \xrightarrow{p} \alpha,$$

where we used (1.2) and the LLN for an average of  $u_i$ . Next,

$$\sqrt{n}(\hat{\alpha} - \alpha) = \frac{1}{\sqrt{n}}(\beta - \hat{\beta}) \frac{\rho(1-\rho^{1+n})}{1-\rho} + \frac{1}{\sqrt{n}} \sum_i u_i = U_n + V_n.$$

Because of (1.2),  $U_n \xrightarrow{p} 0$ . From the CLT it follows that  $V_n \xrightarrow{d} \mathcal{N}(0, \sigma^2)$ . Taken together,

$$\sqrt{n}(\hat{\alpha} - \alpha) \xrightarrow{d} \mathcal{N}(0, \sigma^2).$$

Lastly, let us look at  $\hat{\sigma}^2$ :

$$\hat{\sigma}^2 = \frac{1}{n} \sum_i \hat{e}_i^2 = \frac{1}{n} \sum_i \left( (\alpha - \hat{\alpha}) + (\beta - \hat{\beta})x_i + u_i \right)^2. \quad (1.3)$$

Using the facts that: (1)  $(\alpha - \hat{\alpha})^2 \xrightarrow{p} 0$ , (2)  $(\beta - \hat{\beta})^2/n \xrightarrow{p} 0$ , (3)  $n^{-1} \sum_i u_i^2 \xrightarrow{p} \sigma^2$ , (4)  $n^{-1} \sum_i u_i \xrightarrow{p} 0$ , (5)  $n^{-1/2} \sum_i \rho^i u_i \xrightarrow{p} 0$ , we conclude that

$$\hat{\sigma}^2 \xrightarrow{p} \sigma^2.$$

The rest of this solution is optional and is usually not meant when the asymptotics of  $\hat{\sigma}^2$  is concerned. Before proceeding to deriving its asymptotic distribution, we would like to mark out that  $n^{-\delta}(\beta - \hat{\beta}) \xrightarrow{p} 0$  and  $n^{-\delta} \sum_i \rho^i u_i \xrightarrow{p} 0$  for any  $\delta > 0$ . Using the same algebra as before we get

$$\sqrt{n}(\hat{\sigma}^2 - \sigma^2) \stackrel{A}{\approx} \frac{1}{\sqrt{n}} \sum_i (u_i^2 - \sigma^2),$$

since the other terms converge in probability to zero. Using the CLT, we get

$$\sqrt{n}(\hat{\sigma}^2 - \sigma^2) \xrightarrow{d} \mathcal{N}(0, m_4),$$

where  $m_4 = \mathbb{E}[u_i^4] - \sigma^4$ , provided that it is finite.

## 1.10 Power trends

1. The OLS estimator satisfies

$$\hat{\beta} - \beta = \left( \sum_{i=1}^n x_i^2 \right)^{-1} \sum_{i=1}^n x_i \sigma_i \varepsilon_i = \sqrt{\delta} \left( \sum_{i=1}^n i^{2\lambda} \right)^{-1} \sum_{i=1}^n i^{\lambda+\mu/2} \varepsilon_i.$$

We can see that  $\mathbb{E}[\hat{\beta} - \beta] = 0$  and

$$\mathbb{V}[\hat{\beta} - \beta] = \delta \left( \sum_{i=1}^n i^{2\lambda} \right)^{-2} \sum_{i=1}^n i^{2\lambda+\mu}.$$

If  $\mathbb{V}[\hat{\beta} - \beta] \rightarrow 0$ , the estimator  $\hat{\beta}$  will be consistent. This will occur when  $\mu < 2\lambda + 1$  (in this case the continuous analog  $\left( \int^T t^{2\lambda} dt \right)^2 \sim T^{2(2\lambda+1)}$  of the first sum squared diverges faster or converges slower than the continuous analog  $\int^T t^{2\lambda+\mu} dt \sim T^{2\lambda+\mu+1}$  of the second sum, as  $2(2\lambda + 1) < 2\lambda + \mu + 1$  if and only if  $\mu < 2\lambda + 1$ ). In this case the asymptotic distribution is

$$\begin{aligned} n^{\lambda+(1-\mu)/2}(\hat{\beta} - \beta) &= \sqrt{\delta} n^{\lambda+(1-\mu)/2} \left( \sum_{i=1}^n i^{2\lambda} \right)^{-1} \sum_{i=1}^n i^{\lambda+\mu/2} \varepsilon_i \\ &\xrightarrow{d} \mathcal{N} \left( 0, \delta \lim_{n \rightarrow \infty} n^{2\lambda+1-\mu} \left( \sum_{i=1}^n i^{2\lambda} \right)^{-2} \sum_{i=1}^n i^{2\lambda+\mu} \right) \end{aligned}$$

by Lyapunov's CLT for independent heterogeneous observations, provided that

$$\frac{\left( \sum_{i=1}^n i^{3(2\lambda+\mu)/2} \right)^{1/3}}{\left( \sum_{i=1}^n i^{2\lambda+\mu} \right)^{1/2}} \rightarrow 0$$

as  $n \rightarrow \infty$ , which is satisfied as  $\left( \int^T t^{2\lambda+\mu} dt \right)^{1/2} \sim T^{(2\lambda+\mu+1)/2}$  diverges faster or converges slower than  $\left( \int^T t^{3(2\lambda+\mu)/2} dt \right)^{1/3} \sim T^{(3(2\lambda+\mu)/2+1)/3}$ .

2. The GLS estimator satisfies

$$\tilde{\beta} - \beta = \left( \sum_{i=1}^n \frac{x_i^2}{i^\mu} \right)^{-1} \sum_{i=1}^n \frac{x_i \sigma_i \varepsilon_i}{i^\mu} = \sqrt{\delta} \left( \sum_{i=1}^n i^{2\lambda-\mu} \right)^{-1} \sum_{i=1}^n i^{\lambda-\mu/2} \varepsilon_i.$$

Again,  $\mathbb{E}[\tilde{\beta} - \beta] = 0$  and

$$\mathbb{V}[\tilde{\beta} - \beta] = \delta \left( \sum_{i=1}^n i^{2\lambda-\mu} \right)^{-2} \sum_{i=1}^n i^{2\lambda-\mu}.$$

The estimator  $\tilde{\beta}$  will be consistent under the same condition, i.e. when  $\mu < 2\lambda + 1$ . In this case the asymptotic distribution is

$$\begin{aligned} n^{\lambda+(1-\mu)/2}(\tilde{\beta} - \beta) &= \sqrt{\delta} n^{\lambda+(1-\mu)/2} \left( \sum_{i=1}^n i^{2\lambda-\mu} \right)^{-1} \sum_{i=1}^n i^{\lambda-\mu/2} \varepsilon_i \\ &\xrightarrow{d} \mathcal{N} \left( 0, \delta \lim_{n \rightarrow \infty} n^{2\lambda+1-\mu} \left( \sum_{i=1}^n i^{2\lambda-\mu} \right)^{-2} \sum_{i=1}^n i^{2\lambda-\mu} \right) \end{aligned}$$

by Lyapunov's CLT for independent heterogeneous observations.

## 2. ASYMPTOTIC THEORY: TIME SERIES

### 2.1 Trended vs. differenced regression

1. The OLS estimator  $\hat{\beta}$  in this case is

$$\hat{\beta} = \frac{\sum_{t=1}^T (y_t - \bar{y})(t - \bar{t})}{\sum_{t=1}^T (t - \bar{t})^2}.$$

Now,

$$\hat{\beta} - \beta = \left( \frac{1}{T^{-3} \sum_t t^2 - (T^{-2} \sum_t t)^2}, -\frac{T^{-2} \sum_t t}{T^{-3} \sum_t t^2 - (T^{-2} \sum_t t)^2} \right) \begin{pmatrix} T^{-3} \sum_t \varepsilon_t t \\ T^{-2} \sum_t \varepsilon_t \end{pmatrix}.$$

Because

$$\sum_{t=1}^T t = \frac{T(T+1)}{2}, \quad \sum_{t=1}^T t^2 = \frac{T(T+1)(2T+1)}{6},$$

it is easy to see that the first vector converges to  $(12, -6)$ . Therefore,

$$T^{3/2}(\hat{\beta} - \beta) = (12, -6) \frac{1}{\sqrt{T}} \sum_t \begin{pmatrix} \frac{t}{T} \varepsilon_t \\ \varepsilon_t \end{pmatrix}.$$

Assuming that all conditions for the CLT for heterogenous martingale difference sequences (e.g., Hamilton (1994) "Time series analysis", proposition 7.8) hold, we find that

$$\frac{1}{\sqrt{T}} \sum_{t=1}^T \begin{pmatrix} \frac{t}{T} \varepsilon_t \\ \varepsilon_t \end{pmatrix} \xrightarrow{d} \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \sigma^2 \begin{pmatrix} \frac{1}{3} & \frac{1}{2} \\ \frac{1}{2} & 1 \end{pmatrix} \right),$$

because

$$\begin{aligned} \lim \frac{1}{T} \sum_{t=1}^T \mathbb{V} \left[ \begin{pmatrix} \frac{t}{T} \varepsilon_t \\ \varepsilon_t \end{pmatrix} \right] &= \sigma^2 \lim \frac{1}{T} \sum_{t=1}^T \begin{pmatrix} \frac{t}{T} \\ 1 \end{pmatrix}^2 = \frac{1}{3}, \\ \lim \frac{1}{T} \sum_{t=1}^T \mathbb{V} [\varepsilon_t] &= \sigma^2, \\ \lim \frac{1}{T} \sum_{t=1}^T \mathbb{C} \left[ \begin{pmatrix} \frac{t}{T} \varepsilon_t \\ \varepsilon_t \end{pmatrix} \right] &= \sigma^2 \lim \frac{1}{T} \sum_{t=1}^T \frac{t}{T} = \frac{1}{2}. \end{aligned}$$

Consequently,

$$T^{3/2}(\hat{\beta} - \beta) \rightarrow (12, -6) \cdot \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \sigma^2 \begin{pmatrix} \frac{1}{3} & \frac{1}{2} \\ \frac{1}{2} & 1 \end{pmatrix} \right) = \mathcal{N}(0, 12\sigma^2).$$

2. For the regression in differences, the OLS estimator is

$$\check{\beta} = \frac{1}{T} \sum_{t=1}^T (y_t - y_{t-1}) = \beta + \frac{\varepsilon_T - \varepsilon_0}{T}.$$

So,  $T(\check{\beta} - \beta) = \varepsilon_T - \varepsilon_0 \sim \mathcal{D}(0, 2\sigma^2)$ .

3. When  $T$  is sufficiently large,  $\hat{\beta} \stackrel{A}{\sim} \mathcal{N}\left(\beta, \frac{12\sigma^2}{T^3}\right)$ , and  $\check{\beta} \sim \mathcal{D}\left(\beta, \frac{2\sigma^2}{T^2}\right)$ . It is easy to see that for large  $T$ , the (approximate) variance of the first estimator is less than that of the second. Intuitively, it is much easier to identify the trend of a trending sequence than the drift of a drifting sequence.

## 2.2 Long run variance for AR(1)

The long run variance is  $V_{ze} = \sum_{j=-\infty}^{+\infty} \mathbb{C}[z_t e_t, z_{t-j} e_{t-j}]$ . Because  $e_t$  and  $z_t$  are scalars, independent at all lags and leads, and  $\mathbb{E}[e_t] = 0$ , we have  $\mathbb{C}[z_t e_t, z_{t-j} e_{t-j}] = \mathbb{E}[z_t z_{t-j}] \mathbb{E}[e_t e_{t-j}]$ . Let for simplicity  $z_t$  also have zero mean. Then for  $j \geq 0$ ,  $\mathbb{E}[z_t z_{t-j}] = \rho_z^j (1 - \rho_z^2)^{-1} \sigma_z^2$  and  $\mathbb{E}[e_t e_{t-j}] = \rho_e^j (1 - \rho_e^2)^{-1} \sigma_e^2$ , where  $\rho_z, \sigma_z^2, \rho_e, \sigma_e^2$  are AR(1) parameters. To sum up,

$$V_{ze} = \frac{\sigma_z^2}{1 - \rho_z^2} \frac{\sigma_e^2}{1 - \rho_e^2} \sum_{j=-\infty}^{+\infty} \rho_z^{|j|} \rho_e^{|j|} = \frac{1 + \rho_z \rho_e}{(1 - \rho_z \rho_e)(1 - \rho_z^2)(1 - \rho_e^2)} \sigma_z^2 \sigma_e^2.$$

To estimate  $V_{ze}$ , find the OLS estimates  $\hat{\rho}_z, \hat{\sigma}_z^2, \hat{\rho}_e, \hat{\sigma}_e^2$  from AR(1) regressions and plug them in. The resulting  $\hat{V}_{ze}$  will be consistent by the Continuous Mapping Theorem.

## 2.3 Asymptotics of averages of AR(1) and MA(1)

Note that  $y_t$  can be rewritten as  $y_t = \sum_{j=0}^{+\infty} \rho^j x_{t-j}$ .

1. (i)  $y_t$  is not an MDS relative to own past  $\{y_{t-1}, y_{t-2}, \dots\}$  as it is correlated with older  $y_t$ 's; (ii)  $z_t$  is an MDS relative to  $\{x_{t-2}, x_{t-3}, \dots\}$ , but is not an MDS relative to own past  $\{z_{t-1}, z_{t-2}, \dots\}$ , because  $z_t$  and  $z_{t-1}$  are correlated through  $x_{t-1}$ .
2. (i) By the CLT for the general stationary and ergodic case,  $\sqrt{T} \bar{y}_T \xrightarrow{d} \mathcal{N}(0, q_{yy})$ , where  $q_{yy} = \sum_{j=-\infty}^{+\infty} \underbrace{\mathbb{C}[y_t, y_{t-j}]}_{\gamma_j}$ . It can be shown that for an AR(1) process,  $\gamma_0 = \frac{\sigma^2}{1 - \rho^2}$ ,  $\gamma_j = \frac{\sigma^2}{1 - \rho^2} \rho^j$ . Therefore,  $q_{yy} = \sum_{j=-\infty}^{+\infty} \gamma_j = \frac{\sigma^2}{(1 - \rho)^2}$ . (ii) By the CLT for the general stationary and ergodic case,  $\sqrt{T} \bar{z}_T \xrightarrow{d} \mathcal{N}(0, q_{zz})$ , where  $q_{zz} = \gamma_0 + 2\gamma_1 + 2 \underbrace{\sum_{j=2}^{+\infty} \gamma_j}_{=0} = (1 + \theta^2)\sigma^2 + 2\theta\sigma^2 = \sigma^2(1 + \theta)^2$ .
3. If we have consistent estimates  $\hat{\sigma}^2, \hat{\rho}, \hat{\theta}$  of  $\sigma^2, \rho, \theta$ , we can estimate  $q_{yy}$  and  $q_{zz}$  consistently by  $\frac{\hat{\sigma}^2}{(1 - \hat{\rho})^2}$  and  $\hat{\sigma}^2(1 + \hat{\theta})^2$ , respectively. Note that these are positive numbers by construction. Alternatively, we could use robust estimators, like the Newey–West nonparametric estimator, ignoring additional information that we have. But under the circumstances this seems to be less efficient.



4. For vectors, (i)  $\sqrt{T}\bar{\mathbf{y}}_T \xrightarrow{d} \mathcal{N}(0, Q_{yy})$ , where  $Q_{yy} = \sum_{j=-\infty}^{+\infty} \underbrace{\mathbb{C}[\mathbf{y}_t, \mathbf{y}_{t-j}]}_{\Gamma_j}$ . But  $\Gamma_0 = \sum_{j=0}^{+\infty} \mathbf{P}^j \Sigma \mathbf{P}'^j$ ,

$\Gamma_j = \mathbf{P}^j \Gamma_0$  if  $j > 0$ , and  $\Gamma_j = \Gamma'_{-j} = \Gamma_0 \mathbf{P}'^{|j|}$  if  $j < 0$ . Hence  $Q_{yy} = \Gamma_0 + \sum_{j=1}^{+\infty} \mathbf{P}^j \Gamma_0 + \sum_{j=1}^{+\infty} \Gamma_0 \mathbf{P}'^j = \Gamma_0 + (\mathbf{I} - \mathbf{P})^{-1} \mathbf{P} \Gamma_0 + \Gamma_0 \mathbf{P}' (\mathbf{I} - \mathbf{P}')^{-1} = (\mathbf{I} - \mathbf{P})^{-1} \Sigma (\mathbf{I} - \mathbf{P}')^{-1}$ ; (ii)  $\sqrt{T}\bar{\mathbf{z}}_T \xrightarrow{d} \mathcal{N}(0, Q_{zz})$ , where  $Q_{zz} = \Gamma_0 + \Gamma_1 + \Gamma_{-1} = \Sigma + \Theta \Sigma \Theta' + \Theta \Sigma + \Sigma \Theta' = (\mathbf{I} + \Theta) \Sigma (\mathbf{I} + \Theta)'$ . As for estimation of asymptotic variances, it is evidently possible to construct consistent estimators of  $Q_{yy}$  and  $Q_{zz}$  that are positive definite by construction.

## 2.4 Asymptotics for impulse response functions

1. For the AR(1) process, we get by repeated substitution

$$y_t = \sum_{j=0}^{\infty} \rho^j \varepsilon_{t-j}.$$

Since the weights decline exponentially, their sum absolutely converges to a finite constant. The IRF is

$$IRF(j) = \rho^j, \quad j \geq 0.$$

The ARMA(1,1) process written via lag polynomials is

$$z_t = \frac{1 - \theta L}{1 - \rho L} \varepsilon_t,$$

of which the MA( $\infty$ ) representation is

$$z_t = \varepsilon_t + (\rho - \theta) \sum_{i=0}^{\infty} \rho^i \varepsilon_{t-i-1}.$$

Since the weights decline exponentially, their sum absolutely converges to a finite constant. The IRF is

$$IRF(0) = 1, \quad IRF(j) = (\rho - \theta) \rho^{j-1}, \quad j > 0.$$

2. The estimate based on the OLS estimator  $\hat{\rho}$  of  $\rho$ , is  $\widehat{IRF}(j) = \hat{\rho}^j$ . Since  $\sqrt{T}(\hat{\rho} - \rho) \xrightarrow{d} \mathcal{N}(0, 1 - \rho^2)$ , we can derive using the Delta-method that

$$\sqrt{T} \left( \widehat{IRF}(j) - IRF(j) \right) \xrightarrow{d} \mathcal{N} \left( 0, j^2 \rho^{2(j-1)} (1 - \rho^2) \right)$$

as  $T \rightarrow \infty$  when  $j \geq 1$ , and  $\widehat{IRF}(0) - IRF(0) = 0$ .

3. Denote  $e_t = \varepsilon_t - \theta \varepsilon_{t-1}$ . Since  $\hat{\rho} \xrightarrow{p} \rho$  and  $\hat{\theta} \xrightarrow{p} \theta$  (this will be shown below), we can construct consistent estimators as

$$\widehat{IRF}(0) = 1, \quad \widehat{IRF}(j) = (\hat{\rho} - \hat{\theta}) \hat{\rho}^{j-1}, \quad j > 0.$$

Evidently,  $\widehat{IRF}(0)$  has a degenerate distribution. To derive the asymptotic distribution of  $\widehat{IRF}(j)$  for  $j > 0$ , let us first derive the asymptotic distribution of  $(\hat{\rho}, \hat{\theta})'$ . Consistency can be shown easily:

$$\hat{\rho} = \rho + \frac{T^{-1} \sum_{t=3}^T z_{t-2} e_t}{T^{-1} \sum_{t=3}^T z_{t-2} z_{t-1}} \xrightarrow{p} \rho + \frac{\mathbb{E}[z_{t-2} e_t]}{\mathbb{E}[z_{t-1} z_t]} = \rho,$$

$$-\frac{\hat{\theta}}{1+\hat{\theta}^2} = \frac{\sum_{t=2}^T \hat{e}_t \hat{e}_{t-1}}{\sum_{t=2}^T \hat{e}_t^2} = \dots \quad \text{expansion of } \hat{e}_t \quad \dots \xrightarrow{p} -\frac{\theta}{1+\theta^2}.$$

Since the solution of a quadratic equation is a continuous function of its coefficients, consistency of  $\hat{\theta}$  obtains. To derive the asymptotic distribution, we need the following component:

$$\frac{1}{\sqrt{T}} \sum_{t=3}^T \begin{pmatrix} e_t e_{t-1} - \mathbb{E}[e_t e_{t-1}] \\ e_t^2 - \mathbb{E}[e_t^2] \\ z_{t-2} e_t \end{pmatrix} \xrightarrow{d} \mathcal{N}(0, \Omega),$$

where  $\Omega$  is a  $3 \times 3$  variance matrix which is a function of  $\theta$ ,  $\rho$ ,  $\sigma_\varepsilon^2$  and  $\kappa = \mathbb{E}[\varepsilon_t^4]$  (derivation is very tedious; one should account for serial correlation in summands). Next, from examining the formula defining  $\hat{\theta}$ , dropping the terms that do not contribute to the asymptotic distribution, we can find that

$$\begin{aligned} \sqrt{T} \left( -\frac{\hat{\theta}}{1+\hat{\theta}^2} - \left( -\frac{\theta}{1+\theta^2} \right) \right) &\stackrel{A}{=} \alpha_1 \sqrt{T} (\hat{\rho} - \rho) \\ &+ \alpha_2 \frac{1}{\sqrt{T}} \sum (e_t e_{t-1} - \mathbb{E}[e_t e_{t-1}]) + \alpha_3 \frac{1}{\sqrt{T}} \sum (e_t^2 - \mathbb{E}[e_t^2]) \end{aligned}$$

for certain constants  $\alpha_1, \alpha_2, \alpha_3$ . It follows by the Delta-method that

$$\begin{aligned} \sqrt{T} (\hat{\theta} - \theta) &\stackrel{A}{=} -\frac{(1+\theta^2)^2}{1-\theta^2} \sqrt{T} \left( -\frac{\hat{\theta}}{1+\hat{\theta}^2} - \left( -\frac{\theta}{1+\theta^2} \right) \right) \\ &\stackrel{A}{=} \beta_1 \sqrt{T} (\hat{\rho} - \rho) + \beta_2 \frac{1}{\sqrt{T}} \sum (e_t e_{t-1} - \mathbb{E}[e_t e_{t-1}]) + \beta_3 \frac{1}{\sqrt{T}} \sum (e_t^2 - \mathbb{E}[e_t^2]) \end{aligned}$$

for certain constants  $\beta_1, \beta_2, \beta_3$ . It follows that

$$\sqrt{T} \begin{pmatrix} \hat{\rho} - \rho \\ \hat{\theta} - \theta \end{pmatrix} \stackrel{A}{=} \Gamma \frac{1}{\sqrt{T}} \sum \begin{pmatrix} z_{t-2} e_t \\ e_t^2 - \mathbb{E}[e_t^2] \\ e_t e_{t-1} - \mathbb{E}[e_t e_{t-1}] \end{pmatrix} \xrightarrow{d} \mathcal{N}(0, \Gamma \Omega \Gamma').$$

for certain  $2 \times 3$  matrix  $\Gamma$ . Finally, applying the Delta-method again, we get

$$\sqrt{T} \left( \widehat{IRF}(j) - IRF(j) \right) \stackrel{A}{=} \gamma' \sqrt{T} \begin{pmatrix} \hat{\rho} - \rho \\ \hat{\theta} - \theta \end{pmatrix} \xrightarrow{d} \mathcal{N}(0, \gamma' \Gamma \Omega \Gamma' \gamma),$$

for certain  $2 \times 1$  vector  $\gamma$ .

# 3. BOOTSTRAP

---

## 3.1 Brief and exhaustive

1. The mentioned difference indeed exists, but it is not the principal one. The two methods have some common features like computer simulations, sampling, etc., but they serve completely different goals. The bootstrap is an alternative to analytical asymptotic theory for making inferences, while Monte–Carlo is used for studying finite-sample properties of estimators or tests.
2. After some point raising  $B$ , the number of bootstrap repetitions, does not help since the bootstrap distribution is intrinsically discrete, and raising  $B$  cannot smooth things out. Even more than that: if we are interested in quantiles (we usually are), the quantile for a discrete distribution is an interval, and the uncertainty about which point to choose to be a quantile does not disappear when we raise  $B$ .
3. There is no such thing as a “bootstrap estimator”. Bootstrapping is a method of inference, not of estimation. The same goes for an “asymptotic estimator”.

## 3.2 Bootstrapping $t$ -ratio

The percentile interval is  $CI_{\%} = [\hat{\theta} - \tilde{q}_{1-\alpha/2}^*, \hat{\theta} - \tilde{q}_{\alpha/2}^*]$ , where  $\tilde{q}_{\alpha}^*$  is a bootstrap  $\alpha$ -quantile of  $\hat{\theta}^* - \hat{\theta}$ , i.e.  $\alpha = \mathbb{P}\{\hat{\theta}^* - \hat{\theta} \leq \tilde{q}_{\alpha}^*\}$ . Then  $\frac{\tilde{q}_{\alpha}^*}{s(\hat{\theta})}$  is the  $\alpha$ -quantile of  $\frac{\hat{\theta}^* - \hat{\theta}}{s(\hat{\theta})} = T_n^*$ , since  $\mathbb{P}\left\{\frac{\hat{\theta}^* - \hat{\theta}}{s(\hat{\theta})} \leq \frac{\tilde{q}_{\alpha}^*}{s(\hat{\theta})}\right\} = \alpha$ . But by construction, the  $\alpha$ -quantile of  $T_n^*$  is  $q_{\alpha}^*$ , hence  $\tilde{q}_{\alpha}^* = s(\hat{\theta})q_{\alpha}^*$ . Substituting this into  $CI_{\%}$  we get the confidence interval as  $CI$  in the problem.

## 3.3 Bootstrap bias correction

1. The bootstrap version  $\bar{x}_n^*$  of  $\bar{x}_n$  has mean  $\bar{x}_n$  with respect to the EDF:  $\mathbb{E}^*[\bar{x}_n^*] = \bar{x}_n$ . Thus the bootstrap version of the bias (which is itself zero) is  $\text{Bias}^*(\bar{x}_n) = \mathbb{E}^*[\bar{x}_n^*] - \bar{x}_n = 0$ . Therefore, the bootstrap bias corrected estimator of  $\mu$  is  $\bar{x}_n - \text{Bias}^*(\bar{x}_n) = \bar{x}_n$ . Now consider the bias of  $\bar{x}_n^2$ :

$$\text{Bias}(\bar{x}_n^2) = \mathbb{E}[\bar{x}_n^2] - \mu^2 = \mathbb{V}[\bar{x}_n] = \frac{1}{n}\mathbb{V}[x].$$

Thus the bootstrap version of the bias is the sample analog of this quantity:

$$\text{Bias}^*(\bar{x}_n^2) = \frac{1}{n}\mathbb{V}^*[x] = \frac{1}{n}\left(\frac{1}{n}\sum x_i^2 - \bar{x}_n^2\right).$$

Therefore, the bootstrap bias corrected estimator of  $\mu^2$  is

$$\bar{x}_n^2 - \text{Bias}^*(\bar{x}_n^2) = \frac{n+1}{n}\bar{x}_n^2 - \frac{1}{n^2}\sum x_i^2.$$

2. Since the sample average is an unbiased estimator of the population mean for any distribution, the bootstrap bias correction for  $\bar{z}^2$  will be zero, and thus the bias-corrected estimator for  $\mathbb{E}[z^2]$  will be

$$\bar{z}^2$$

(cf. part 1). Next note that the bootstrap distribution is 0 with probability  $\frac{1}{2}$  and 3 with probability  $\frac{1}{2}$ , so the bootstrap distribution for  $\bar{z}^2 = \frac{1}{4}(z_1 + z_2)^2$  is 0 with probability  $\frac{1}{4}$ ,  $\frac{9}{4}$  with probability  $\frac{1}{2}$ , and 9 with probability  $\frac{1}{4}$ . Thus the bootstrap bias estimate is

$$\frac{1}{4}\left(0 - \frac{9}{4}\right) + \frac{1}{2}\left(\frac{9}{4} - \frac{9}{4}\right) + \frac{1}{4}\left(9 - \frac{9}{4}\right) = \frac{9}{8},$$

and the bias corrected version is

$$\frac{1}{4}(z_1 + z_2)^2 - \frac{9}{8}.$$

### 3.4 Bootstrap in linear model

1. Due to the assumption of random sampling, there cannot be unconditional heteroskedasticity. If conditional heteroskedasticity is present, it does not invalidate the nonparametric bootstrap. The dependence of the conditional variance on regressors is not destroyed by bootstrap resampling as the data  $(x_i, y_i)$  are resampled in pairs.
2. When we bootstrap an inconsistent estimator, its bootstrap analogs are concentrated more and more around the probability limit of the estimator, and thus the estimate of the bias becomes smaller and smaller as the sample size grows. That is, bootstrapping is able to correct the bias caused by finite sample nonsymmetry of the distribution, but not the asymptotic bias (difference between the probability limit of the estimator and the true parameter value).
3. As a point estimate we take  $\hat{g}(x) = x'\hat{\beta}$ , where  $\hat{\beta}$  is the OLS estimator for  $\beta$ . To pivotize  $\hat{g}(x)$ , observe that

$$\sqrt{n}x'(\hat{\beta} - \beta) \xrightarrow{d} \mathcal{N}(0, x'(\mathbb{E}[xx'])^{-1}\mathbb{E}[xx'e^2](\mathbb{E}[xx'])^{-1}x),$$

so the appropriate statistic to bootstrap is

$$t_g = \frac{x'(\hat{\beta} - \beta)}{\text{se}(\hat{g}(x))},$$

where  $\text{se}(\hat{g}(x)) = \sqrt{x'(\sum_i x_i x_i')^{-1}(\sum_i x_i x_i' \hat{e}_i^2)(\sum_i x_i x_i')^{-1}x}$ . The bootstrap version is

$$t_g^* = \frac{x'(\hat{\beta}^* - \hat{\beta})}{\text{se}(\hat{g}^*(x))},$$

where  $\text{se}(\hat{g}^*(x)) = \sqrt{x'(\sum_i x_i^* x_i^{*'})^{-1}(\sum_i x_i^* x_i^{*'} \hat{e}_i^{*2})(\sum_i x_i^* x_i^{*'})^{-1}x}$ . The rest is standard, and the confidence interval is

$$CI_t = \left[ x'\hat{\beta} - q_{1-\frac{\alpha}{2}}^* \text{se}(\hat{g}(x)); x'\hat{\beta} - q_{\frac{\alpha}{2}}^* \text{se}(\hat{g}(x)) \right],$$

where  $q_{\frac{\alpha}{2}}^*$  and  $q_{1-\frac{\alpha}{2}}^*$  are appropriate bootstrap quantiles for  $t_g^*$ .

### 3.5 Bootstrap for impulse response functions

1. For each  $j \geq 1$  simulate the bootstrap distribution of the absolute value of the  $t$ -statistic:

$$|t_j| = \frac{\sqrt{T}|\hat{\rho}^j - \rho^j|}{j|\hat{\rho}|^{j-1}\sqrt{1 - \hat{\rho}^2}},$$

the bootstrap analog of which is

$$|t_j^*| = \frac{\sqrt{T}|\hat{\rho}^{*j} - \hat{\rho}^j|}{j|\hat{\rho}^*|^{j-1}\sqrt{1 - \hat{\rho}^{*2}}},$$

read off the bootstrap quantiles  $q_{j,1-\alpha}^*$  and construct the symmetric percentile- $t$  confidence interval  $\left[ \hat{\rho}^j \mp q_{j,1-\alpha}^* \cdot j|\hat{\rho}|^{j-1}\sqrt{(1 - \hat{\rho}^2)/T} \right]$ .

2. Most appropriate is the residual bootstrap when bootstrap samples are generated by resampling estimates of innovations  $\varepsilon_t$ . The corrected estimates of the IRFs are

$$\widetilde{IRF}(j) = 2 \left( \hat{\rho} - \hat{\theta} \right) \hat{\rho}^{j-1} - \frac{1}{B} \sum_{b=1}^B \left( \hat{\rho}_b^* - \hat{\theta}_b^* \right) \hat{\rho}_b^{*j-1},$$

where  $\hat{\rho}_b^*, \hat{\theta}_b^*$  are obtained in  $b^{th}$  bootstrap repetition by using the same formulae as used for  $\hat{\rho}, \hat{\theta}$  but computed from the corresponding bootstrap sample.



# 4. REGRESSION AND PROJECTION

## 4.1 Regressing and projecting dice

(i) The joint distribution is

$$(x, y) = \begin{cases} (0, 1) & \text{with probability } \frac{1}{6}, \\ (2, 2) & \text{with probability } \frac{1}{6}, \\ (0, 3) & \text{with probability } \frac{1}{6}, \\ (4, 4) & \text{with probability } \frac{1}{6}, \\ (0, 5) & \text{with probability } \frac{1}{6}, \\ (6, 6) & \text{with probability } \frac{1}{6}. \end{cases}$$

(ii) The best predictor is

$$\mathbb{E}[y|x] = \begin{cases} 3 & x = 0, \\ 2 & x = 2, \\ 4 & x = 4, \\ 6 & x = 6, \end{cases}$$

and undefined for all other  $x$ .

(iii) To find the best linear predictor, we need  $\mathbb{E}[x] = 2$ ,  $\mathbb{E}[y] = \frac{7}{2}$ ,  $\mathbb{E}[xy] = \frac{28}{3}$ ,  $\mathbb{V}[x] = \frac{16}{3}$ ,  $\beta = \mathbb{C}[x, y] / \mathbb{V}[x] = \frac{7}{16}$ ,  $\alpha = \frac{21}{8}$ , so

$$\text{BLP}[y|x] = \frac{21}{8} + \frac{7}{16}x.$$

(iv)

$$U_{BP} = \begin{cases} -2 & \text{with probability } \frac{1}{6}, \\ 0 & \text{with probability } \frac{2}{3}, \\ 2 & \text{with probability } \frac{1}{6}, \end{cases}$$

$$U_{BLP} = \begin{cases} -\frac{13}{8} & \text{with probability } \frac{1}{6}, \\ -\frac{12}{8} & \text{with probability } \frac{1}{6}, \\ \frac{3}{8} & \text{with probability } \frac{1}{6}, \\ -\frac{3}{8} & \text{with probability } \frac{1}{6}, \\ \frac{19}{8} & \text{with probability } \frac{1}{6}, \\ \frac{6}{8} & \text{with probability } \frac{1}{6}. \end{cases}$$

so  $\mathbb{E}[U_{BP}^2] = \frac{4}{3}$ ,  $\mathbb{E}[U_{BLP}^2] \approx 1.9$ . Indeed,  $\mathbb{E}[U_{BP}^2] < \mathbb{E}[U_{BLP}^2]$ .

## 4.2 Mixture of normals

1. We need to compute  $\mathbb{E}[y]$ ,  $\mathbb{E}[x]$ ,  $\mathbb{C}[x, y]$ ,  $\mathbb{V}[x]$ . Let us denote the first event  $A$ , the second event  $\bar{A}$ . Then

$$\begin{aligned}\mathbb{E}[y] &= p \cdot \mathbb{E}[y|A] + (1-p) \cdot \mathbb{E}[y|\bar{A}] = p \cdot 4 + (1-p) \cdot 0 = 4p, \\ \mathbb{E}[x] &= p \cdot \mathbb{E}[x|A] + (1-p) \cdot \mathbb{E}[x|\bar{A}] = p \cdot 0 + (1-p) \cdot 4 = 4 - 4p, \\ \mathbb{E}[xy] &= p \cdot (\mathbb{E}[x|A] \mathbb{E}[y|A] + \mathbb{C}[x, y|A]) + (1-p) \cdot (\mathbb{E}[x|\bar{A}] \mathbb{E}[y|\bar{A}] + \mathbb{C}[x, y|\bar{A}]) = 0, \\ \mathbb{C}[x, y] &= \mathbb{E}[xy] - \mathbb{E}[x] \mathbb{E}[y] = -16p(1-p), \\ \mathbb{E}[x^2] &= p \cdot (\mathbb{E}[x|A]^2 + \mathbb{V}[x|A]) + (1-p) \cdot (\mathbb{E}[x|\bar{A}]^2 + \mathbb{V}[x|\bar{A}]) = 17 - 16p, \\ \mathbb{V}[x] &= \mathbb{E}[x^2] - \mathbb{E}[x]^2 = 1 + 16p - 16p^2.\end{aligned}$$

Summarizing,

$$\begin{aligned}\beta &= \frac{\mathbb{C}[x, y]}{\mathbb{V}[x]} = \frac{1}{1 + 16p - 16p^2} - 1, \quad \alpha = \mathbb{E}[y] - \beta \mathbb{E}[x] = 4 \left( 1 - \frac{1-p}{1 + 16p - 16p^2} \right), \\ \Rightarrow \quad \mathbb{BLP}[y|x] &= 4 \left( 1 - \frac{1-p}{1 + 16p - 16p^2} \right) + \left( \frac{1}{1 + 16p - 16p^2} - 1 \right) x,\end{aligned}$$

2. If  $\mathbb{E}[y|x]$  was linear, it would coincide with  $\mathbb{BLP}[y|x]$ . But take, for example, the value of  $\mathbb{BLP}[y|x]$  for some big  $x$ . This value will be negative. But, given this large  $x$ , the mean of  $y$  cannot be negative, as it should clearly be between 0 and 4. Contradiction. To derive  $\mathbb{E}[y|x]$ , one have to write that

$$\begin{aligned}f\left(\begin{pmatrix} x \\ y \end{pmatrix}\right) &= p \frac{1}{2\pi} \exp\left(-\frac{1}{2}x^2 - \frac{1}{2}(y-4)^2\right) + (1-p) \frac{1}{2\pi} \exp\left(-\frac{1}{2}(x-4)^2 - \frac{1}{2}y^2\right), \\ f(x) &= p \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right) + (1-p) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x-4)^2\right),\end{aligned}$$

so

$$f(y|x) = \frac{1}{\sqrt{2\pi}} \frac{p \exp\left(-\frac{1}{2}x^2 - \frac{1}{2}(y-4)^2\right) + (1-p) \exp\left(-\frac{1}{2}(x-4)^2 - \frac{1}{2}y^2\right)}{p \exp\left(-\frac{1}{2}x^2\right) + (1-p) \exp\left(-\frac{1}{2}(x-4)^2\right)}.$$

The conditional expectation is the integral

$$\mathbb{E}[y|x] = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} y \frac{p \exp\left(-\frac{1}{2}x^2 - \frac{1}{2}(y-4)^2\right) + (1-p) \exp\left(-\frac{1}{2}(x-4)^2 - \frac{1}{2}y^2\right)}{p \exp\left(-\frac{1}{2}x^2\right) + (1-p) \exp\left(-\frac{1}{2}(x-4)^2\right)} dy$$

One can compute this in **Maple**<sup>®</sup> to get

$$\mathbb{E}[y|x] = \frac{4}{1 + (p^{-1} - 1) \exp(4x - 8)}.$$

It has a form of a smooth transition regression function, and has horizontal asymptotes 4 and 0 as  $x \rightarrow -\infty$  and  $x \rightarrow +\infty$ , respectively.



### 4.3 Bernoulli regressor

Note that

$$\mathbb{E}[y|x] = \begin{cases} \mu_0, & x = 0, \\ \mu_1, & x = 1, \end{cases} = \mu_0(1-x) + \mu_1x = \mu_0 + (\mu_1 - \mu_0)x$$

and

$$\mathbb{E}[y^2|x] = \begin{cases} \mu_0^2 + \sigma_0^2, & x = 0, \\ \mu_1^2 + \sigma_1^2, & x = 1, \end{cases} = \mu_0^2 + \sigma_0^2 + (\mu_1^2 - \mu_0^2 + \sigma_1^2 - \sigma_0^2)x.$$

These expectations are linear in  $x$  because the support of  $x$  has only two points, and one can always draw a straight line through two points. The reason is *not* conditional normality!

### 4.4 Best polynomial approximation

The FOC are

$$\begin{aligned} 2\mathbb{E}\left[\left(\mathbb{E}[y|x] - \alpha_0 - \alpha_1x - \dots - \alpha_kx^k\right)(-1)\right] &= 0, \\ 2\mathbb{E}\left[\left(\mathbb{E}[y|x] - \alpha_0 - \alpha_1x - \dots - \alpha_kx^k\right)(-x)\right] &= 0, \\ &\vdots \\ 2\mathbb{E}\left[\left(\mathbb{E}[y|x] - \alpha_0 - \alpha_1x - \dots - \alpha_kx^k\right)(-x^k)\right] &= 0, \end{aligned}$$

or, using the LIME,

$$\begin{aligned} \mathbb{E}\left[y - \alpha_0 - \alpha_1x - \dots - \alpha_kx^k\right] &= 0, \\ \mathbb{E}\left[\left(y - \alpha_0 - \alpha_1x - \dots - \alpha_kx^k\right)x\right] &= 0, \\ &\vdots \\ \mathbb{E}\left[\left(y - \alpha_0 - \alpha_1x - \dots - \alpha_kx^k\right)x^k\right] &= 0, \end{aligned}$$

from where

$$\alpha = (\alpha_0, \alpha_1, \dots, \alpha_k)' = \mathbb{E}[zz']^{-1} \mathbb{E}[zy],$$

where  $z = (1, x, \dots, x^k)'$ . The best  $k^{\text{th}}$  order polynomial approximation is then

$$\mathbb{BPA}_k[y|x] = z'\alpha = z'\mathbb{E}[zz']^{-1} \mathbb{E}[zy].$$

The associated prediction error  $u_k$  has the properties that follow from the FOC:

$$\mathbb{E}[zu_k] = 0,$$

i.e. it has zero mean and is uncorrelated with  $x, x^2, \dots, x^k$ .

## 4.5 Handling conditional expectations

1. By the LIME,  $\mathbb{E}[y|x, z] = \alpha + \beta x + \gamma z$ . Thus we know that in the linear prediction  $y = \alpha + \beta x + \gamma z + e_y$ , the prediction error  $e_y$  is uncorrelated with the predictors, i.e.  $\mathbb{C}[e_y, x] = \mathbb{C}[e_y, z] = 0$ . Consider the linear prediction of  $z$  by  $x$ :  $z = \zeta + \delta x + e_z$ ,  $\mathbb{C}[e_z, x] = 0$ . But since  $\mathbb{C}[z, x] = 0$ , we know that  $\delta = 0$ . Now, if we linearly predict  $y$  only by  $x$ , we will have  $y = \alpha + \beta x + \gamma(\zeta + e_z) + e_y = \alpha + \gamma\zeta + \beta x + \gamma e_z + e_y$ . Here the composite error  $\gamma e_z + e_y$  is uncorrelated with  $x$  and thus is the best linear prediction error. As a result, the OLS estimator of  $\beta$  is consistent.

Checking the properties of the second option is more involved. Notice that the OLS coefficients in the linear prediction of  $y$  by  $x$  and  $w$  converge in probability to

$$\text{plim} \begin{pmatrix} \hat{\beta} \\ \hat{\omega} \end{pmatrix} = \begin{pmatrix} \sigma_x^2 & \sigma_{xw} \\ \sigma_{xw} & \sigma_w^2 \end{pmatrix}^{-1} \begin{pmatrix} \sigma_{xy} \\ \sigma_{wy} \end{pmatrix} = \begin{pmatrix} \sigma_x^2 & \sigma_{xw} \\ \sigma_{xw} & \sigma_w^2 \end{pmatrix}^{-1} \begin{pmatrix} \beta\sigma_x^2 \\ \beta\sigma_{xw} + \gamma\sigma_{wz} \end{pmatrix},$$

so we can see that

$$\text{plim} \hat{\beta} = \beta + \frac{\sigma_{xw}\sigma_{wz}}{\sigma_x^2\sigma_w^2 - \sigma_{xw}^2}\gamma.$$

Thus in general the second option gives an inconsistent estimator.

2. Because  $\hat{E}[x|z] = g(z)$  is a strictly increasing and continuous function,  $g^{-1}(\cdot)$  exists and  $\mathbb{E}[x|z] = \gamma$  is equivalent to  $z = g^{-1}(\gamma)$ . If  $\hat{E}[y|z] = f(z)$ , then  $\hat{E}[y|\mathbb{E}[x|z] = \gamma] = f(g^{-1}(\gamma))$ .

# 5. LINEAR REGRESSION AND OLS

## 5.1 Fixed and random regressors

1. It simplifies a lot. First, we can use simpler versions of LLNs and CLTs; second, we do not need additional conditions apart from existence of some moments. For example, for consistency of the OLS estimator in the linear mean regression model  $y = x\beta + e$ ,  $\mathbb{E}[e|x] = 0$ , only existence of moments is needed, while in the case of fixed regressors we (i) have to use the LLN for heterogeneous sequences, (ii) have to add the condition  $\frac{1}{n} \sum_{i=1}^n x_i^2 \rightarrow M$  as  $n \rightarrow \infty$ .
2. The economist is probably right about treating the regressors as random if he/she has a random sampling experiment. But the reasoning is ridiculous. For a sampled individual, his/her characteristics (whether true or false) are fixed; randomness arises from the fact that this individual is *randomly selected*.
3. The OLS estimator is unbiased conditional on the whole sample of  $x$ -variables, irrespective of how  $x$ 's are generated. The conditional unbiasedness property implies unbiasedness.

## 5.2 Consistency of OLS under serially correlated errors

1. Indeed,

$$\mathbb{E}[u_t] = \mathbb{E}[y_t - \beta y_{t-1}] = \mathbb{E}[y_t] - \beta \mathbb{E}[y_{t-1}] = 0 - \beta \cdot 0 = 0$$

and

$$\mathbb{C}[u_t, y_{t-1}] = \mathbb{C}[y_t - \beta y_{t-1}, y_{t-1}] = \mathbb{C}[y_t, y_{t-1}] - \beta \mathbb{V}[y_{t-1}] = 0.$$

- (ii) Now let us show that  $\hat{\beta}$  is consistent. Because  $\mathbb{E}[y_t] = 0$ , it immediately follows that

$$\hat{\beta} = \frac{T^{-1} \sum_{t=2}^T y_t y_{t-1}}{T^{-1} \sum_{t=2}^T y_{t-1}^2} = \beta + \frac{T^{-1} \sum_{t=2}^T u_t y_{t-1}}{T^{-1} \sum_{t=2}^T y_{t-1}^2} \xrightarrow{p} \beta + \frac{\mathbb{E}[u_t y_{t-1}]}{\mathbb{E}[y_{t-1}^2]} = \beta.$$

- (iii) To show that  $u_t$  is serially correlated, consider

$$\mathbb{C}[u_t, u_{t-1}] = \mathbb{C}[y_t - \beta y_{t-1}, y_{t-1} - \beta y_{t-2}] = \beta (\beta \mathbb{C}[y_t, y_{t-1}] - \mathbb{C}[y_t, y_{t-2}]),$$

which is generally not zero unless  $\beta = 0$  or  $\beta = \frac{\mathbb{C}[y_t, y_{t-2}]}{\mathbb{C}[y_t, y_{t-1}]}$ . As an example of a serially correlated  $u_t$  take the AR(2) process

$$y_t = \alpha y_{t-2} + \varepsilon_t,$$

where  $\varepsilon_t$  is a strict white noise. Then  $\beta = 0$  and thus  $u_t = y_t$ , serially correlated.

- (iv) The OLS estimator is inconsistent if the error term is correlated with the right-hand-side variables. This is not the same as serial correlatedness of the error term.

### 5.3 Estimation of linear combination

1. Consider the class of linear estimators, i.e. one having the form  $\tilde{\theta} = \mathcal{A}\mathcal{Y}$ , where  $\mathcal{A}$  depends only on data  $\mathcal{X} = ((1, x_1, z_1)' \dots (1, x_n, z_n)')'$ . The conditional unbiasedness requirement yields the condition  $\mathcal{A}\mathcal{X} = (1, c_x, c_z)\delta$ , where  $\delta = (\alpha, \beta, \gamma)'$ . The best linear unbiased estimator is  $\hat{\theta} = (1, c_x, c_z)\hat{\delta}$ , where  $\hat{\delta}$  is the OLS estimator. Indeed, this estimator belongs to the class considered, since  $\hat{\theta} = (1, c_x, c_z)(\mathcal{X}'\mathcal{X})^{-1}\mathcal{X}'\mathcal{Y} = \mathcal{A}^*\mathcal{Y}$  for  $\mathcal{A}^* = (1, c_x, c_z)(\mathcal{X}'\mathcal{X})^{-1}\mathcal{X}'$  and  $\mathcal{A}^*\mathcal{X} = (1, c_x, c_z)$ . Besides,

$$\mathbb{V}[\hat{\theta}|\mathcal{X}] = \sigma^2(1, c_x, c_z)(\mathcal{X}'\mathcal{X})^{-1}(1, c_x, c_z)'$$

and is minimal in the class because the key relationship  $(\mathcal{A} - \mathcal{A}^*)\mathcal{A}^* = 0$  holds.

2. Observe that  $\sqrt{n}(\hat{\theta} - \theta) = (1, c_x, c_z)\sqrt{n}(\hat{\delta} - \delta) \xrightarrow{d} \mathcal{N}(0, V_{\hat{\theta}})$ , where

$$V_{\hat{\theta}} = \sigma^2 \left( 1 + \frac{\phi_x^2 + \phi_z^2 - 2\rho\phi_x\phi_z}{1 - \rho^2} \right),$$

$\phi_x = (\mathbb{E}[x] - c_x)/\sqrt{\mathbb{V}[x]}$ ,  $\phi_z = (\mathbb{E}[z] - c_z)/\sqrt{\mathbb{V}[z]}$ , and  $\rho$  is the correlation coefficient between  $x$  and  $z$ .

3. Minimization of  $V_{\hat{\theta}}$  with respect to  $\rho$  yields

$$\rho^{opt} = \begin{cases} \frac{\phi_x}{\phi_z} & \text{if } \left| \frac{\phi_x}{\phi_z} \right| < 1, \\ \frac{\phi_z}{\phi_x} & \text{if } \left| \frac{\phi_x}{\phi_z} \right| \geq 1. \end{cases}$$

4. Multicollinearity between  $x$  and  $z$  means that  $\rho = 1$  and  $\delta$  and  $\theta$  are unidentified. An implication is that the asymptotic variance of  $\hat{\theta}$  is infinite.

### 5.4 Incomplete regression

1. Note that

$$y = x'\beta + z'\gamma + \eta.$$

We know that  $\mathbb{E}[\eta|z] = 0$ , so  $\mathbb{E}[z\eta] = 0$ . However,  $\mathbb{E}[x\eta] \neq 0$  unless  $\gamma = 0$ , because  $0 = \mathbb{E}[xe] = \mathbb{E}[x(z'\gamma + \eta)] = \mathbb{E}[xz']\gamma + \mathbb{E}[x\eta]$ , and we know that  $\mathbb{E}[xz'] \neq 0$ . The regression of  $y$  on  $x$  and  $z$  yields the OLS estimates with the probability limit

$$p \lim \begin{pmatrix} \hat{\beta} \\ \hat{\gamma} \end{pmatrix} = \begin{pmatrix} \beta \\ \gamma \end{pmatrix} + Q^{-1} \begin{pmatrix} \mathbb{E}[x\eta] \\ 0 \end{pmatrix},$$

where

$$Q = \begin{pmatrix} \mathbb{E}[xx'] & \mathbb{E}[xz'] \\ \mathbb{E}[zx'] & \mathbb{E}[zz'] \end{pmatrix}.$$

We can see that the estimators  $\hat{\beta}$  and  $\hat{\gamma}$  are in general inconsistent. To be more precise, the inconsistency of *both*  $\hat{\beta}$  and  $\hat{\gamma}$  is proportional to  $\mathbb{E}[x\eta]$ , so that unless  $\gamma = 0$  (or, more subtly, unless  $\gamma$  lies in the null space of  $\mathbb{E}[xz']$ ), the estimators are inconsistent.

2. The first step yields a consistent OLS estimate  $\hat{\beta}$  of  $\beta$  because the OLS estimator is consistent in a linear mean regression. At the second step, we get the OLS estimate

$$\begin{aligned}\hat{\gamma} &= \left( \sum z_i z_i' \right)^{-1} \sum z_i \hat{e}_i = \left( \sum z_i z_i' \right)^{-1} \left( \sum z_i e_i - \sum z_i x_i' (\hat{\beta} - \beta) \right) = \\ &= \gamma + \left( n^{-1} \sum z_i z_i' \right)^{-1} \left( n^{-1} \sum z_i \eta_i - n^{-1} \sum z_i x_i' (\hat{\beta} - \beta) \right).\end{aligned}$$

Since  $n^{-1} \sum z_i z_i' \xrightarrow{p} \mathbb{E}[zz']$ ,  $n^{-1} \sum z_i x_i' \xrightarrow{p} \mathbb{E}[zx']$ ,  $n^{-1} \sum z_i \eta_i \xrightarrow{p} \mathbb{E}[z\eta] = 0$ ,  $\hat{\beta} - \beta \xrightarrow{p} 0$ , we have that  $\hat{\gamma}$  is consistent for  $\gamma$ .

Therefore, from the point of view of consistency of  $\hat{\beta}$  and  $\hat{\gamma}$ , we recommend the second method. The limiting distribution of  $\sqrt{n}(\hat{\gamma} - \gamma)$  can be deduced by using the Delta-method. Observe that

$$\sqrt{n}(\hat{\gamma} - \gamma) = \left( n^{-1} \sum_i z_i z_i' \right)^{-1} \left( n^{-1/2} \sum_i z_i \eta_i - n^{-1} \sum_i z_i x_i' \left( n^{-1} \sum_i x_i x_i' \right)^{-1} n^{-1/2} \sum_i x_i e_i \right)$$

and

$$\frac{1}{\sqrt{n}} \sum_i \begin{pmatrix} z_i \eta_i \\ x_i e_i \end{pmatrix} \xrightarrow{d} \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \mathbb{E}[zz'\eta^2] & \mathbb{E}[zx'\eta e] \\ \mathbb{E}[xz'\eta e] & \sigma^2 \mathbb{E}[xx'] \end{pmatrix} \right).$$

Having applied the Delta-method and the Continuous Mapping Theorems, we get

$$\sqrt{n}(\hat{\gamma} - \gamma) \xrightarrow{d} \mathcal{N} \left( 0, (\mathbb{E}[zz'])^{-1} V (\mathbb{E}[zz'])^{-1} \right),$$

where

$$\begin{aligned}V &= \mathbb{E}[zz'\eta^2] + \sigma^2 \mathbb{E}[zx'] (\mathbb{E}[xx'])^{-1} \mathbb{E}[xz'] \\ &\quad - \mathbb{E}[zx'] (\mathbb{E}[xx'])^{-1} \mathbb{E}[xz'\eta e] - \mathbb{E}[xz'\eta e] (\mathbb{E}[xx'])^{-1} \mathbb{E}[xz'].\end{aligned}$$

## 5.5 Generated coefficient

Observe that

$$\sqrt{n}(\hat{\beta} - \beta) = \left( \frac{1}{n} \sum_{i=1}^n x_i^2 \right)^{-1} \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i u_i - \sqrt{n}(\hat{\alpha} - \alpha) \cdot \frac{1}{n} \sum_{i=1}^n x_i z_i \right).$$

Now,

$$\frac{1}{n} \sum_{i=1}^n x_i^2 \xrightarrow{p} \gamma_x^2, \quad \frac{1}{n} \sum_{i=1}^n x_i z_i \xrightarrow{p} \gamma_{xz}$$

by the LLN, and

$$\begin{pmatrix} \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i u_i \\ \sqrt{n}(\hat{\alpha} - \alpha) \end{pmatrix} \xrightarrow{d} \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \gamma_x^2 & 0 \\ 0 & 1 \end{pmatrix} \right)$$

by the CLT. We can assert that the convergence here is joint (i.e., as of a *vector* sequence) because of independence of the components. Because of their independence, their joint CDF is just a product of marginal CDFs, and pointwise convergence of these marginal CDFs implies pointwise

convergence of the joint CDF. This is important, since generally weak convergence of components of a vector sequence separately does not imply joint weak convergence.

Now, by the Slutsky theorem,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n x_i u_i - \sqrt{n}(\hat{\alpha} - \alpha) \cdot \frac{1}{n} \sum_{i=1}^n x_i z_i \xrightarrow{d} \mathcal{N}(0, \gamma_x^2 + \gamma_{xz}^2).$$

Applying the Slutsky theorem again, we find:

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} (\gamma_x^2)^{-1} \mathcal{N}(0, \gamma_x^2 + \gamma_{xz}^2) = \mathcal{N}\left(0, \frac{1}{\gamma_x^2} + \frac{\gamma_{xz}^2}{\gamma_x^4}\right).$$

Note how the preliminary estimation step affects the precision in estimation of other parameters: the asymptotic variance is blown up. The implication is that sequential estimation makes “naive” (i.e. which ignore preliminary estimation steps) standard errors invalid.

## 5.6 OLS in nonlinear model

Observe that  $\mathbb{E}[y|x] = \alpha + \beta x$ ,  $\mathbb{V}[y|x] = (\alpha + \beta x)^2$ . Consequently, we can use the usual OLS estimator and White’s standard errors. By the way, the model  $y = (\alpha + \beta x)e$  can be viewed as  $y = \alpha + \beta x + u$ , where  $u = (\alpha + \beta x)(e - 1)$ ,  $\mathbb{E}[u|x] = 0$ ,  $\mathbb{V}[u|x] = (\alpha + \beta x)^2$ .

## 5.7 Long and short regressions

Let us denote this estimator by  $\check{\beta}_1$ . We have

$$\begin{aligned} \check{\beta}_1 &= (\mathcal{X}'_1 \mathcal{X}_1)^{-1} \mathcal{X}'_1 \mathcal{Y} = (\mathcal{X}'_1 \mathcal{X}_1)^{-1} \mathcal{X}'_1 (\mathcal{X}_1 \beta_1 + \mathcal{X}_2 \beta_2 + \mathcal{E}) = \\ &= \beta_1 + \left(\frac{1}{n} \mathcal{X}'_1 \mathcal{X}_1\right)^{-1} \left(\frac{1}{n} \mathcal{X}'_1 \mathcal{X}_2\right) \beta_2 + \left(\frac{1}{n} \mathcal{X}'_1 \mathcal{X}_1\right)^{-1} \left(\frac{1}{n} \mathcal{X}'_1 \mathcal{E}\right). \end{aligned}$$

Since  $\mathbb{E}[ex_1] = 0$ , we have that  $n^{-1} \mathcal{X}'_1 \mathcal{E} \xrightarrow{p} 0$  by the LLN. Also, by the LLN,  $n^{-1} \mathcal{X}'_1 \mathcal{X}_1 \xrightarrow{p} \mathbb{E}[x_{1i} x'_{1i}]$  and  $n^{-1} \mathcal{X}'_1 \mathcal{X}_2 \xrightarrow{p} \mathbb{E}[x_1 x'_2]$ . Therefore,

$$\check{\beta}_1 \xrightarrow{p} \beta_1 + (\mathbb{E}[x_1 x'_1])^{-1} \mathbb{E}[x_1 x'_2] \beta_2.$$

So, in general,  $\check{\beta}_1$  is inconsistent. It will be consistent if  $\beta_2$  lies in the null space of  $\mathbb{E}[x_1 x'_2]$ . Two special cases of this are: (1) when  $\beta_2 = 0$ , i.e. when the true model is  $\mathcal{Y} = \mathcal{X}_1 \beta_1 + e$ ; (2) when  $\mathbb{E}[x_1 x'_2] = 0$ .

## 5.8 Ridge regression

1. There is conditional bias:  $\mathbb{E}[\check{\beta}|\mathcal{X}] = (\mathcal{X}'\mathcal{X} + \lambda I_k)^{-1} \mathcal{X}'\mathbb{E}[\mathcal{Y}|\mathcal{X}] = \beta - (\mathcal{X}'\mathcal{X} + \lambda I_k)^{-1} \lambda \beta$ , unless  $\beta = 0$ . Next,  $\mathbb{E}[\check{\beta}] = \beta - \mathbb{E}[(\mathcal{X}'\mathcal{X} + \lambda I_k)^{-1}] \lambda \beta \neq \beta$  unless  $\beta = 0$ . Therefore, estimator is in general biased.

2. Observe that

$$\begin{aligned}\tilde{\beta} &= (\mathcal{X}'\mathcal{X} + \lambda I_k)^{-1} \mathcal{X}'\mathcal{X}\beta + (\mathcal{X}'\mathcal{X} + \lambda I_k)^{-1} \mathcal{X}'\varepsilon \\ &= \left( \frac{1}{n} \sum_i x_i x_i' + \frac{\lambda}{n} I_k \right)^{-1} \frac{1}{n} \sum_i x_i x_i' \beta + \left( \frac{1}{n} \sum_i x_i x_i' + \frac{\lambda}{n} I_k \right)^{-1} \frac{1}{n} \sum_i x_i \varepsilon_i.\end{aligned}$$

Since  $n^{-1} \sum x_i x_i' \xrightarrow{p} \mathbb{E}[xx']$ ,  $n^{-1} \sum x_i \varepsilon_i \xrightarrow{p} \mathbb{E}[x\varepsilon] = 0$ ,  $\lambda/n \rightarrow 0$ , we have:

$$\tilde{\beta} \xrightarrow{p} (\mathbb{E}[xx'])^{-1} \mathbb{E}[xx'] \beta + (\mathbb{E}[xx'])^{-1} 0 = \beta,$$

that is,  $\tilde{\beta}$  is consistent.

3. The math is straightforward:

$$\begin{aligned}\sqrt{n}(\tilde{\beta} - \beta) &= \underbrace{\left( \frac{1}{n} \sum_i x_i x_i' + \frac{\lambda}{n} I_k \right)^{-1}}_{\downarrow^p} \underbrace{\frac{-\lambda}{\sqrt{n}} \beta}_{\downarrow 0} + \underbrace{\left( \frac{1}{n} \sum_i x_i x_i' + \frac{\lambda}{n} I_k \right)^{-1}}_{\downarrow^p} \underbrace{\frac{1}{\sqrt{n}} \sum_i x_i \varepsilon_i}_{\downarrow^d} \\ &\xrightarrow{d} \mathcal{N}(0, Q_{xx}^{-1} Q_{xx\varepsilon^2} Q_{xx}^{-1}).\end{aligned}$$

4. The conditional mean squared error criterion  $\mathbb{E}[(\tilde{\beta} - \beta)^2 | \mathcal{X}]$  can be used. For the OLS estimator,

$$\mathbb{E}[(\hat{\beta} - \beta)^2 | \mathcal{X}] = \mathbb{V}[\hat{\beta}] = (\mathcal{X}'\mathcal{X})^{-1} \mathcal{X}'\Omega\mathcal{X}(\mathcal{X}'\mathcal{X})^{-1}.$$

For the ridge estimator,

$$\mathbb{E}[(\tilde{\beta} - \beta)^2 | \mathcal{X}] = (\mathcal{X}'\mathcal{X} + \lambda I_k)^{-1} (\mathcal{X}'\Omega\mathcal{X} + \lambda^2 \beta\beta') (\mathcal{X}'\mathcal{X} + \lambda I_k)^{-1}.$$

By the first order approximation, if  $\lambda$  is small,  $(\mathcal{X}'\mathcal{X} + \lambda I_k)^{-1} \approx (\mathcal{X}'\mathcal{X})^{-1} (I_k - \lambda(\mathcal{X}'\mathcal{X})^{-1})$ . Hence,

$$\begin{aligned}\mathbb{E}[(\tilde{\beta} - \beta)^2 | \mathcal{X}] &\approx (\mathcal{X}'\mathcal{X})^{-1} (I - \lambda(\mathcal{X}'\mathcal{X})^{-1}) (\mathcal{X}'\Omega\mathcal{X}) (I - \lambda(\mathcal{X}'\mathcal{X})^{-1}) (\mathcal{X}'\mathcal{X})^{-1} \\ &\approx \mathbb{E}[(\hat{\beta} - \beta)^2] - \lambda(\mathcal{X}'\mathcal{X})^{-1} [\mathcal{X}'\Omega\mathcal{X}(\mathcal{X}'\mathcal{X})^{-1} + (\mathcal{X}'\mathcal{X})^{-1} \mathcal{X}'\Omega\mathcal{X}] (\mathcal{X}'\mathcal{X})^{-1}.\end{aligned}$$

That is  $\mathbb{E}[(\hat{\beta} - \beta)^2 | \mathcal{X}] - \mathbb{E}[(\tilde{\beta} - \beta)^2 | \mathcal{X}] = A$ , where  $A$  is positive definite (exercise: show this). Thus for small  $\lambda$ ,  $\tilde{\beta}$  is a preferable estimator to  $\hat{\beta}$  according to the mean squared error criterion, despite its biasedness.

## 5.9 Inconsistency under alternative

We are interesting in the question whether the  $t$ -statistics can be used to check  $H_0 : \beta = 0$ . In order to answer this question we have to investigate the asymptotic properties of  $\hat{\beta}$ . First of all, under the null

$$\hat{\beta} \xrightarrow{p} \frac{\mathbb{C}[z, y]}{\mathbb{V}[z]} = \beta \frac{\mathbb{V}[x]}{\mathbb{V}[z]} = 0.$$

It is straightforward to show that under the null the conventional standard error correctly estimates (i.e. if correctly normalized, is consistent for) the asymptotic variance of  $\hat{\beta}$ . That is, under the null

$$t_{\beta} \xrightarrow{d} \mathcal{N}(0, 1),$$

which means that we can use the conventional  $t$ -statistics for testing  $H_0$ .

## 5.10 Returns to schooling

1. There is nothing wrong in dividing insignificant estimates by each other. In fact, this is the correct way to get a point estimate for the ratio of parameters, according to the Delta-method. Of course, the Delta-method with

$$g\left(\begin{pmatrix} u_1 \\ u_2 \end{pmatrix}\right) = \frac{u_1}{u_2}$$

should be used in full to get an asymptotic confidence interval for this ratio, if needed. But the confidence interval will necessarily be bounded and centered at 3 in this example.

2. The person from the audience is probably right in his feelings and intentions. But the suggested way of implementing those ideas are full of mistakes. First, to compare coefficients in a separate regression, one needs to combine them into one regression using dummy variables. But what is even more important – there is no place for the sup-Wald test here, because the dummy variables are totally known (in different words, the threshold is known and need not be estimates as in threshold regressions). Second, computing standard errors using the bootstrap will not change the approximation principle because the person will still use asymptotic critical values. So nothing will essentially change, unless full scale bootstrapping is used. To say nothing about the fact that in practice hardly the significance of many coefficients will change after switching from asymptotics to bootstrap.



# 6. HETEROSKEDASTICITY AND GLS

## 6.1 Conditional variance estimation

In the OLS case, the method works not because each  $\sigma^2(x_i)$  is estimated by  $\hat{e}_i^2$ , but because

$$\frac{1}{n} \mathcal{X}' \hat{\Omega} \mathcal{X} = \frac{1}{n} \sum_{i=1}^n x_i x_i' \hat{e}_i^2$$

consistently estimates  $\mathbb{E}[xx'e^2] = \mathbb{E}[xx'\sigma^2(x)]$ . In the GLS case, the same trick does not work:

$$\frac{1}{n} \mathcal{X}' \hat{\Omega}^{-1} \mathcal{X} = \frac{1}{n} \sum_{i=1}^n \frac{x_i x_i'}{\hat{e}_i^2}$$

can potentially consistently estimate  $\mathbb{E}[xx'/e^2]$ , but this is not the same as  $\mathbb{E}[xx'/\sigma^2(x)]$ . Of course,  $\hat{\Omega}$  cannot consistently estimate  $\Omega$ , econometrician B is right about this, but the trick in the OLS case works for a completely different reason.

## 6.2 Exponential heteroskedasticity

1. At the first step, get  $\hat{\beta}$ , a consistent estimate of  $\beta$  (for example, OLS). Then construct  $\hat{\sigma}_i^2 \equiv \exp(x_i' \hat{\beta})$  for all  $i$  (we don't need  $\exp(\alpha)$  since it is a multiplicative scalar that eventually cancels out) and use these weights at the second step to construct a feasible GLS estimator of  $\beta$ :

$$\tilde{\beta} = \left( \frac{1}{n} \sum_i \hat{\sigma}_i^{-2} x_i x_i' \right)^{-1} \frac{1}{n} \sum_i \hat{\sigma}_i^{-2} x_i y_i.$$

2. The feasible GLS estimator is asymptotically efficient, since it is asymptotically equivalent to GLS. It is finite-sample inefficient, since we changed the weights from what GLS presumes.

## 6.3 OLS and GLS are identical

1. Evidently,  $\mathbb{E}[\mathcal{Y}|\mathcal{X}] = \mathcal{X}\beta$  and  $\Sigma = \mathbb{V}[\mathcal{Y}|\mathcal{X}] = \mathcal{X}\Gamma\mathcal{X}' + \sigma^2 I_n$ . Since the latter depends on  $\mathcal{X}$ , we are in the heteroskedastic environment.
2. The OLS estimator is

$$\hat{\beta} = (\mathcal{X}'\mathcal{X})^{-1} \mathcal{X}'\mathcal{Y},$$

and the GLS estimator is

$$\tilde{\beta} = (\mathcal{X}'\Sigma^{-1}\mathcal{X})^{-1} \mathcal{X}'\Sigma^{-1}\mathcal{Y}.$$

First,

$$\mathcal{X}'\widehat{\varepsilon} = \mathcal{X}'(\mathbf{y} - \mathcal{X}(\mathcal{X}'\mathcal{X})^{-1}\mathcal{X}'\mathbf{y}) = \mathcal{X}'\mathbf{y} - \mathcal{X}'\mathcal{X}(\mathcal{X}'\mathcal{X})^{-1}\mathcal{X}'\mathbf{y} = \mathcal{X}'\mathbf{y} - \mathcal{X}'\mathbf{y} = 0.$$

Premultiply this by  $\mathcal{X}\Gamma$ :

$$\mathcal{X}\Gamma\mathcal{X}'\widehat{\varepsilon} = 0.$$

Add  $\sigma^2\widehat{\varepsilon}$  to both sides and combine the terms on the left-hand side:

$$(\mathcal{X}\Gamma\mathcal{X}' + \sigma^2 I_n)\widehat{\varepsilon} \equiv \Sigma\widehat{\varepsilon} = \sigma^2\widehat{\varepsilon}.$$

Now predividing by matrix  $\Sigma$  gives

$$\widehat{\varepsilon} = \sigma^2\Sigma^{-1}\widehat{\varepsilon}.$$

Premultiply once again by  $\mathcal{X}'$  to get

$$0 = \mathcal{X}'\widehat{\varepsilon} = \sigma^2\mathcal{X}'\Sigma^{-1}\widehat{\varepsilon},$$

or just  $\mathcal{X}'\Sigma^{-1}\widehat{\varepsilon} = 0$ . Recall now what  $\widehat{\varepsilon}$  is:

$$\mathcal{X}'\Sigma^{-1}\mathbf{y} = \mathcal{X}'\Sigma^{-1}\mathcal{X}(\mathcal{X}'\mathcal{X})^{-1}\mathcal{X}'\mathbf{y}$$

which implies  $\widehat{\beta} = \widetilde{\beta}$ . The fact that the two estimators are identical implies that all the statistics based on the two will be identical and thus have the same distribution.

3. Evidently, in this model the coincidence of the two estimators gives unambiguous superiority of the OLS estimator. In spite of heteroskedasticity, it is efficient in the class of linear unbiased estimators, since it coincides with GLS. The GLS estimator is worse since its feasible version requires estimation of  $\Sigma$ , while the OLS estimator does not. Additional estimation of  $\Sigma$  adds noise which may spoil finite sample performance of the GLS estimator. But all this is not typical for ranking OLS and GLS estimators and is a result of a special form of matrix  $\Sigma$ .

## 6.4 OLS and GLS are equivalent

1. When  $\Sigma\mathcal{X} = \mathcal{X}\Theta$ , we have  $\mathcal{X}'\Sigma\mathcal{X} = \mathcal{X}'\mathcal{X}\Theta$  and  $\Sigma^{-1}\mathcal{X} = \mathcal{X}\Theta^{-1}$ , so that

$$\mathbb{V}[\widehat{\beta}|\mathcal{X}] = (\mathcal{X}'\mathcal{X})^{-1}\mathcal{X}'\Sigma\mathcal{X}(\mathcal{X}'\mathcal{X})^{-1} = \Theta(\mathcal{X}'\mathcal{X})^{-1}$$

and

$$\mathbb{V}[\widetilde{\beta}|\mathcal{X}] = (\mathcal{X}'\Sigma^{-1}\mathcal{X})^{-1} = (\mathcal{X}'\mathcal{X}\Theta^{-1})^{-1} = \Theta(\mathcal{X}'\mathcal{X})^{-1}.$$

2. In this example,

$$\Sigma = \sigma^2 \begin{bmatrix} 1 & \rho & \dots & \rho \\ \rho & 1 & \dots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \dots & 1 \end{bmatrix},$$

and  $\Sigma\mathcal{X} = \sigma^2(1 + \rho(n-1)) \cdot (1, 1, \dots, 1)'$  =  $\mathcal{X}\Theta$ , where

$$\Theta = \sigma^2(1 + \rho(n-1)).$$

Thus one does not need to use GLS but instead do OLS to achieve the same finite-sample efficiency.

## 6.5 Equicorrelated observations

This is essentially a repetition of the second part of the previous problem, from which it follows that under the circumstances  $\bar{x}_n$  the best linear conditionally (on a constant which is the same as unconditionally) unbiased estimator of  $\theta$  because of coincidence of its variance with that of the GLS estimator. Appealing to the case when  $|\gamma| > 1$  (which is tempting because then the variance of  $\bar{x}_n$  is larger than that of, say,  $x_1$ ) is invalid, because it is ruled out by the Cauchy–Schwartz inequality.

One cannot appeal to the usual LLNs because  $x$  is non-ergodic. The variance of  $\bar{x}_n$  is  $\mathbb{V}[\bar{x}_n] = \frac{1}{n} \cdot 1 + \frac{n-1}{n} \cdot \gamma \rightarrow \gamma$  as  $n \rightarrow \infty$ , so the estimator  $\bar{x}_n$  is in general inconsistent (except in the case when  $\gamma = 0$ ). For an example of inconsistent  $\bar{x}_n$ , assume that  $\gamma > 0$  and consider the following construct:  $u_i = \varepsilon + \varsigma_i$ , where  $\varsigma_i \sim IID(0, 1 - \gamma)$  and  $\varepsilon \sim (0, \gamma)$  independent of  $\varsigma_i$  for all  $i$ . Then the correlation structure is exactly as in the problem, and  $\frac{1}{n} \sum u_i \xrightarrow{P} \varepsilon$ , a random nonzero limit.

## 6.6 Unbiasedness of certain FGLS estimators

(a)  $0 = \mathbb{E}[z - z] = \mathbb{E}[z] + \mathbb{E}[-z] = \mathbb{E}[z] + \mathbb{E}[z] = 2\mathbb{E}[z]$ . It follows that  $\mathbb{E}[z] = 0$ .

(b)  $\mathbb{E}[q(\varepsilon)] = \mathbb{E}[-q(-\varepsilon)] = \mathbb{E}[-q(\varepsilon)] = -\mathbb{E}[q(\varepsilon)]$ . It follows that  $\mathbb{E}[q(\varepsilon)] = 0$ .

Consider

$$\tilde{\beta}_F - \beta = \left( \mathcal{X}' \hat{\Sigma}^{-1} \mathcal{X} \right)^{-1} \mathcal{X}' \hat{\Sigma}^{-1} \mathcal{E}.$$

Let  $\hat{\Sigma}$  be an estimate of  $\Sigma$  which is a function of products of least squares residuals, i.e.

$$\hat{\Sigma} = F(\mathcal{M} \mathcal{E} \mathcal{E}' \mathcal{M}) = H(\mathcal{E} \mathcal{E}')$$

for  $\mathcal{M} = I - \mathcal{X}(\mathcal{X}' \mathcal{X})^{-1} \mathcal{X}'$ . Conditional on  $\mathcal{X}$ , the expression  $\left( \mathcal{X}' \hat{\Sigma}^{-1} \mathcal{X} \right)^{-1} \mathcal{X}' \hat{\Sigma}^{-1} \mathcal{E}$  is odd in  $\mathcal{E}$ , and  $\mathcal{E}$  and  $-\mathcal{E}$  have the same conditional distribution. Hence by (b),

$$\mathbb{E} \left[ \tilde{\beta}_F - \beta \right] = 0.$$



# 7. VARIANCE ESTIMATION

## 7.1 White estimator

1. Yes, one should use the White formula, but not because  $\sigma^2 Q_{xx}^{-1}$  does not make sense. It does make sense, but is poorly related to the OLS asymptotic variance, which in general takes the “sandwich” form. It is not true that  $\sigma^2$  varies from observation to observation, if by  $\sigma^2$  we mean unconditional variance of the error term.
2. The first part of the claim is totally correct. But availability of the  $t$  or Wald statistics is not enough to do inference. We need critical values for these statistics, and they can be obtained only from some distribution theory, asymptotic in particular.

## 7.2 HAC estimation under homoskedasticity

Under conditional homoskedasticity,

$$\begin{aligned}\mathbb{E} [x_t x'_{t-j} e_t e_{t-j}] &= \mathbb{E} [\mathbb{E} [x_t x'_{t-j} e_t e_{t-j} | x_t, x_{t-1}, \dots]] \\ &= \mathbb{E} [x_t x'_{t-j} \mathbb{E} [e_t e_{t-j} | x_t, x_{t-1}, \dots]] \\ &= \gamma_j \mathbb{E} [x_t x'_{t-j}].\end{aligned}$$

Using this, the long-run variance of  $x_t e_t$  will be

$$\sum_{j=-\infty}^{+\infty} \mathbb{E} [x_t x'_{t-j} e_t e_{t-j}] = \sum_{j=-\infty}^{+\infty} \gamma_j \mathbb{E} [x_t x'_{t-j}],$$

and its Newey–West-type HAC estimator will be

$$\sum_{j=-m}^{+m} \left(1 - \frac{|j|}{m+1}\right) \hat{\gamma}_j \widehat{\mathbb{E} [x_t x'_{t-j}]},$$

where

$$\begin{aligned}\hat{\gamma}_j &= \frac{1}{T} \sum_{t=\max(1,1+j)}^{\min(T,T+j)} (y_t - x'_t \hat{\beta}) (y_{t-j} - x'_{t-j} \hat{\beta}), \\ \widehat{\mathbb{E} [x_t x'_{t-j}]} &= \frac{1}{T} \sum_{t=\max(1,1+j)}^{\min(T,T+j)} x_t x'_{t-j}.\end{aligned}$$

It is not clear whether the estimate is positive definite.

### 7.3 Expectations of White and Newey–West estimators in IID setting

The White formula (apart from the factor  $n$ ) is

$$\hat{V}_{\hat{\beta}} = (\mathcal{X}'\mathcal{X})^{-1} \sum_{i=1}^n x_i x_i' \hat{e}_i^2 (\mathcal{X}'\mathcal{X})^{-1}.$$

Because  $\hat{e}_i = e_i - x_i'(\hat{\beta} - \beta)$ , we have  $\hat{e}_i^2 = e_i^2 - 2x_i'(\hat{\beta} - \beta)e_i + x_i'(\hat{\beta} - \beta)(\hat{\beta} - \beta)'x_i$ . Also recall that  $\hat{\beta} - \beta = (\mathcal{X}'\mathcal{X})^{-1} \sum_{j=1}^n x_j e_j$ . Hence,

$$\begin{aligned} \mathbb{E} \left[ \sum_{i=1}^n x_i x_i' \hat{e}_i^2 | \mathcal{X} \right] &= \mathbb{E} \left[ \sum_{i=1}^n x_i x_i' e_i^2 | \mathcal{X} \right] - 2 \mathbb{E} \left[ \sum_{i=1}^n x_i x_i' x_i' (\hat{\beta} - \beta) e_i | \mathcal{X} \right] \\ &\quad + \mathbb{E} \left[ \sum_{i=1}^n x_i x_i' x_i' (\hat{\beta} - \beta) (\hat{\beta} - \beta)' x_i | \mathcal{X} \right] \\ &= \sigma^2 \sum_{i=1}^n x_i x_i' \left( 1 - x_i' (\mathcal{X}'\mathcal{X})^{-1} x_i \right), \end{aligned}$$

because  $\mathbb{E}[(\hat{\beta} - \beta)e_i | \mathcal{X}] = (\mathcal{X}'\mathcal{X})^{-1} \mathcal{X}' \mathbb{E}[e_i \mathcal{E} | \mathcal{X}] = \sigma^2 (\mathcal{X}'\mathcal{X})^{-1} x_i$  and  $\mathbb{E}[(\hat{\beta} - \beta)(\hat{\beta} - \beta)' | \mathcal{X}] = \sigma^2 (\mathcal{X}'\mathcal{X})^{-1}$ . Finally,

$$\begin{aligned} \mathbb{E}[\hat{V}_{\hat{\beta}} | \mathcal{X}] &= (\mathcal{X}'\mathcal{X})^{-1} \mathbb{E} \left[ \sum_{i=1}^n x_i x_i' \hat{e}_i^2 | \mathcal{X} \right] (\mathcal{X}'\mathcal{X})^{-1} \\ &= \sigma^2 (\mathcal{X}'\mathcal{X})^{-1} \sum_{i=1}^n x_i x_i' \left( 1 - x_i' (\mathcal{X}'\mathcal{X})^{-1} x_i \right) (\mathcal{X}'\mathcal{X})^{-1}. \end{aligned}$$

Let  $\omega_j = 1 - |j|/(m+1)$ . The Newey–West estimator of the asymptotic variance matrix of  $\hat{\beta}$  with lag truncation parameter  $m$  is  $\check{V}_{\hat{\beta}} = (\mathcal{X}'\mathcal{X})^{-1} \hat{S} (\mathcal{X}'\mathcal{X})^{-1}$ , where

$$\hat{S} = \sum_{j=-m}^{+m} \omega_j \sum_{i=\max(1,1+j)}^{\min(n,n+j)} x_i x_{i-j}' \left( e_i - x_i'(\hat{\beta} - \beta) \right) \left( e_{i-j} - x_{i-j}'(\hat{\beta} - \beta) \right).$$

Thus

$$\begin{aligned} \mathbb{E}[\hat{S} | \mathcal{X}] &= \sum_{j=-m}^{+m} \omega_j \sum_{i=\max(1,1+j)}^{\min(n,n+j)} x_i x_{i-j}' \mathbb{E} \left[ \left( e_i - x_i'(\hat{\beta} - \beta) \right) \left( e_{i-j} - x_{i-j}'(\hat{\beta} - \beta) \right) | \mathcal{X} \right] \\ &= \sum_{j=-m}^{+m} \omega_j \sum_{i=\max(1,1+j)}^{\min(n,n+j)} x_i x_{i-j}' \begin{pmatrix} \sigma^2 \mathbb{I}_{\{j=0\}} + x_i' \mathbb{E} \left[ (\hat{\beta} - \beta)(\hat{\beta} - \beta)' | \mathcal{X} \right] x_{i-j} \\ -x_i' \mathbb{E} \left[ (\hat{\beta} - \beta) e_{i-j} | \mathcal{X} \right] - x_{i-j}' \mathbb{E} \left[ e_i (\hat{\beta} - \beta) | \mathcal{X} \right] \end{pmatrix} \\ &= \sigma^2 \mathcal{X}'\mathcal{X} - \sigma^2 \sum_{j=-m}^{+m} \omega_j \sum_{i=\max(1,1+j)}^{\min(n,n+j)} \left( x_i' (\mathcal{X}'\mathcal{X})^{-1} x_{i-j} \right) x_i x_{i-j}'. \end{aligned}$$

Finally,

$$\mathbb{E}[\check{V}_{\hat{\beta}} | \mathcal{X}] = \sigma^2 (\mathcal{X}'\mathcal{X})^{-1} - \sigma^2 (\mathcal{X}'\mathcal{X})^{-1} \sum_{j=-m}^{+m} \omega_j \sum_{i=\max(1,1+j)}^{\min(n,n+j)} \left( x_i' (\mathcal{X}'\mathcal{X})^{-1} x_{i-j} \right) x_i x_{i-j}' (\mathcal{X}'\mathcal{X})^{-1}.$$

# 8. NONLINEAR REGRESSION

## 8.1 Local and global identification

The quasiregressor is  $g_\beta = (1, 2\beta_2 x)'$ . The local ID condition that  $\mathbb{E}[g_\beta g_\beta']$  is of full rank is satisfied since it is equivalent to  $\det \mathbb{E}[g_\beta g_\beta'] = \mathbb{V}[2\beta_2 x] \neq 0$  which holds due to  $\beta_2 \neq 0$  and  $\mathbb{V}[x] \neq 0$ . But the global ID condition fails because the sign of  $\beta_2$  is not identified: together with the true pair  $(\beta_1, \beta_2)'$ , another pair  $(\beta_1, -\beta_2)'$  also minimizes the population least squares criterion.

## 8.2 Identification when regressor is nonrandom

In the linear case,  $Q_{xx} = \mathbb{E}[x^2]$ , a scalar. Its rank (i.e. it itself) equals zero if and only if  $\Pr\{x = 0\} = 1$ , i.e. when  $a = 0$ , the identification condition fails. When  $a \neq 0$ , the identification condition is satisfied. Graphically, when all points lie on a vertical line, we can unambiguously draw a line from the origin through them except when all points are lying on the ordinate axis.

In the nonlinear case,  $Q_{gg} = \mathbb{E}[g_\beta(x, \beta) g_\beta(x, \beta)'] = g_\beta(a, \beta) g_\beta(a, \beta)'$ , a  $k \times k$  matrix. This matrix is a square of a vector having rank 1, hence its rank can be only one or zero. Hence, if  $k > 1$  (there are more than one parameter), this matrix cannot be of full rank, and identification fails. Graphically, there are an infinite number of curves passing through a set of points on a vertical line. If  $k = 1$  and  $g_\beta(a, \beta) \neq 0$ , there is identification; if  $k = 1$  and  $g_\beta(a, \beta) = 0$ , there is identification failure (see the linear case). Intuition in the case  $k = 1$ : if marginal changes in  $\beta$  shift the only regression value  $g(a, \beta)$ , it can be identified; if they do not shift it, many values of  $\beta$  are consistent with the same value of  $g(a, \beta)$ .

## 8.3 Cobb–Douglas production function

1. This idea must have come from a temptation to take logs of the production function:

$$\log Q - \log L = \log \alpha + \theta (\log K - \log L) + \log \varepsilon.$$

But the error in this equation,  $e = \log \varepsilon - \mathbb{E}[\log \varepsilon]$ , even after centering, may not have the property  $\mathbb{E}[e|K, L] = 0$ , or even  $\mathbb{E}[(K, L)'e] = 0$ . As a result, OLS will generally give an inconsistent estimate for  $\theta$ .

2. Note that the regression function is  $\mathbb{E}[Q|K, L] = \alpha K^\theta L^{1-\theta} \mathbb{E}[\varepsilon|K, L] = \alpha K^\theta L^{1-\theta}$ . The regression is nonlinear, but using the concentration method leads to the algorithm described in the problem. Hence, this suggestion gives a consistent estimate for  $\theta$ . [Remark: in practice, though, a researcher would probably include also an intercept to be more confident about validity of specification.]

## 8.4 Exponential regression

The local IC is satisfied: the matrix

$$Q_{gg} = \mathbb{E} \left[ \frac{\partial \exp(\alpha + \beta x)}{\partial \begin{pmatrix} \alpha \\ \beta \end{pmatrix}} \frac{\partial \exp(\alpha + \beta x)}{\partial \begin{pmatrix} \alpha \\ \beta \end{pmatrix}'} \right]_{\beta=0} = \exp(\alpha)^2 \mathbb{E} \left[ \begin{pmatrix} 1 & x \\ x & x^2 \end{pmatrix} \right] = \exp(2\alpha) I_2$$

is invertable. The asymptotic distribution is normal with variance matrix

$$V_{NLLS} = \frac{\sigma^2}{\exp(2\alpha)} I_2.$$

The concentration algorithm uses the grid on  $\beta$ . For each  $\beta$  on this grid, we can estimate  $\exp(\alpha(\beta))$  by OLS from the regression of  $y$  on  $\exp(\beta x)$ , so the estimate and sum of squared residuals are

$$\hat{\alpha}(\beta) = \log \frac{\sum_{i=1}^n \exp(\beta x_i) y_i}{\sum_{i=1}^n \exp(2\beta x_i)},$$

$$SSR(\beta) = \sum_{i=1}^n (y_i - \exp(\hat{\alpha}(\beta) + \beta x_i))^2.$$

Choose such  $\hat{\beta}$  that yields minimum value of  $SSR(\beta)$  on the grid. Set  $\hat{\alpha} = \hat{\alpha}(\hat{\beta})$ . The standard errors  $se(\hat{\alpha})$  and  $se(\hat{\beta})$  can be computed as square roots of the diagonal elements of

$$\frac{SSR(\hat{\beta})}{n} \left( \sum_{i=1}^n \exp(2\hat{\alpha} + 2\hat{\beta} x_i) \begin{pmatrix} 1 & x_i \\ x_i & x_i^2 \end{pmatrix} \right)^{-1}.$$

Note that we cannot use the above expression for  $V_{NLLS}$  since in practice we do not know the distribution of  $x$  and that  $\beta = 0$ .

## 8.5 Power regression

Under  $H_0 : \alpha = 0$ , the parameter  $\beta$  is not identified. Therefore, the Wald (or  $t$ ) statistic does not have a usual asymptotic distribution, and we should use the sup-Wald statistic

$$\sup W = \sup_{\beta} \mathcal{W}(\beta),$$

where  $\mathcal{W}(\beta)$  is the Wald statistic for  $\alpha = 0$  when the unidentified parameter is fixed at value  $\beta$ . The asymptotic distribution is non-standard and can be obtained via simulations.

## 8.6 Simple transition regression

1. The marginal influence is

$$\left. \frac{\partial (\beta_1 + \beta_2 / (1 + \beta_3 x))}{\partial x} \right|_{x=0} = \left. \frac{-\beta_2 \beta_3}{(1 + \beta_3 x)^2} \right|_{x=0} = -\beta_2 \beta_3.$$



So the null is  $H_0 : \beta_2\beta_3 + 1 = 0$ . The  $t$ -statistic is

$$t = \frac{\hat{\beta}_2\hat{\beta}_3 + 1}{se(\hat{\beta}_2\hat{\beta}_3)},$$

where  $\hat{\beta}_2$  and  $\hat{\beta}_3$  are elements of the NLLS estimator, and  $se(\hat{\beta}_2\hat{\beta}_3)$  is a standard error for  $\hat{\beta}_2\hat{\beta}_3$  which can be computed from the NLLS asymptotics and Delta-method. The test rejects when  $|t| > q_{1-\alpha/2}^{N(0,1)}$ .

2. The regression function does not depend on  $x$  when, for example,  $H_0 : \beta_2 = 0$ . As under  $H_0$  the parameter  $\hat{\beta}_3$  is not identified, inference is nonstandard. The Wald statistic for a particular value of  $\beta_3$  is

$$\mathcal{W}(\beta_3) = \left( \frac{\hat{\beta}_2}{se(\hat{\beta}_2)} \right)^2,$$

and the test statistic is

$$\sup W = \sup_{\beta_3} \mathcal{W}(\beta_3).$$

The test rejects when  $\sup W > q_{1-\alpha}^{\mathcal{D}}$ , where the limiting distribution  $\mathcal{D}$  is obtained via simulations.

## 8.7 Nonlinear consumption function

1. When  $\delta = 0$ , the parameter  $\gamma$  is not identified, and when  $\gamma = 0$ , the parameters  $\alpha$  and  $\delta$  are not separately identified. A third situation occurs if the support of  $y_t$  lies entirely above or entirely below  $\Delta$ , the parameters  $\alpha$  and  $\beta$  are not separately identified (note: a separation of this into two “different” situations is unfair).
2. Run the concentration method with respect to the parameter  $\gamma$ . The quasi-regressor is

$$(1, \mathbb{I}\{y_{t-1} > \Delta\}, y_{t-1}^\gamma, \delta y_{t-1}^\gamma \log y_{t-1})'$$

When constructing the standard errors, beside using the quasiregressor we have to apply HAC estimation of asymptotic variance, because the regression error does not have a martingale difference property. Indeed, note that the information  $y_{t-1}, y_{t-2}, y_{t-3}, \dots$  does not include  $c_{t-1}, c_{t-2}, c_{t-3}, \dots$  and therefore, for  $j > 0$  and  $e_t = c_t - \mathbb{E}[c_t | y_{t-1}, y_{t-2}, y_{t-3}, \dots]$

$$\mathbb{E}[e_t e_{t-j}] = \mathbb{E}[\mathbb{E}[e_t e_{t-j} | y_{t-1}, y_{t-2}, y_{t-3}, \dots]] \neq \mathbb{E}[\mathbb{E}[e_t | y_{t-1}, y_{t-2}, y_{t-3}, \dots] e_{t-j}] = 0,$$

as  $e_{t-j}$  is not necessarily measurable relative to  $y_{t-1}, y_{t-2}, y_{t-3}, \dots$

- 3.

- (a) Asymptotically,

$$t_{\hat{\gamma}=1} = \frac{\hat{\gamma} - 1}{se(\hat{\gamma})}$$

is distributed as  $N(0, 1)$ . Thus, reject if  $t_{\hat{\gamma}=1} < q_\alpha^{N(0,1)}$ , or, equivalently,  $\hat{\gamma} < 1 + se(\hat{\gamma})q_\alpha^{N(0,1)}$ .

- (b) Bootstrap  $\hat{\gamma} - 1$  whose bootstrap analog is  $\hat{\gamma}^* - \hat{\gamma}$ , get the left  $\alpha$ -quantile  $q_\alpha^*$ , and reject if  $\hat{\gamma} < 1 + q_\alpha^*$ .



# 9. EXTREMUM ESTIMATORS

## 9.1 Regression on constant

For the first estimator use standard LLN and CLT:

$$\hat{\beta}_1 = \frac{1}{n} \sum_{i=1}^n y_i \xrightarrow{p} \mathbb{E}[y] = \beta \text{ (consistency),}$$

$$\sqrt{n}(\hat{\beta}_1 - \beta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n e_i \xrightarrow{d} \mathcal{N}(0, \mathbb{V}[e]) = \mathcal{N}(0, \beta^2) \text{ (asymptotic normality).}$$

Consider

$$\hat{\beta}_2 = \arg \min_b \left\{ \log b^2 + \frac{1}{nb^2} \sum_{i=1}^n (y_i - b)^2 \right\}. \quad (9.1)$$

Denote

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \quad \overline{y^2} = \frac{1}{n} \sum_{i=1}^n y_i^2.$$

The FOC for this problem gives after rearrangement:

$$\hat{\beta}^2 + \hat{\beta}\bar{y} - \overline{y^2} = 0 \Leftrightarrow \hat{\beta}_{\pm} = -\frac{\bar{y}}{2} \pm \frac{\sqrt{\bar{y}^2 + 4\overline{y^2}}}{2}.$$

The two values  $\hat{\beta}_{\pm}$  correspond to the two different solutions of a *local* minimization problem in population:

$$\mathbb{E} \left[ \log b^2 + \frac{1}{b^2} (y - b)^2 \right] \rightarrow \min_b \Leftrightarrow b_{\pm} = -\frac{\mathbb{E}[y]}{2} \pm \frac{\sqrt{\mathbb{E}[y]^2 + 4\mathbb{E}[y^2]}}{2} = \beta, -2\beta. \quad (9.2)$$

Note that the SOC are satisfied at both  $b_{\pm}$ , so both are local minima. Direct comparison yields that  $b_+ = \beta$  is a global minimum. Hence, the extremum estimator  $\hat{\beta}_2 = \hat{\beta}_+$  is consistent for  $\beta$ . The asymptotics of  $\hat{\beta}_2$  can be found using the theory of extremum estimators. For  $f(y, b) = \log b^2 + b^{-2} (y - b)^2$ ,

$$\frac{\partial f(y, b)}{\partial b} = \frac{2}{b} - \frac{2(y - b)^2}{b^3} - \frac{2(y - b)}{b^2} \Rightarrow \mathbb{E} \left[ \left( \frac{\partial f(y, \beta)}{\partial b} \right)^2 \right] = \frac{4\kappa}{\beta^6},$$

$$\frac{\partial^2 f(y, b)}{\partial b^2} = \frac{6(y - b)^2}{b^4} + \frac{8(y - b)}{b^3} \Rightarrow \mathbb{E} \left[ \frac{\partial^2 f(y, \beta)}{\partial b^2} \right] = \frac{6}{\beta^2}.$$

Consequently,

$$\sqrt{n}(\hat{\beta}_2 - \beta) \xrightarrow{d} \mathcal{N} \left( 0, \frac{\kappa}{9\beta^2} \right).$$

Consider now

$$\hat{\beta}_3 = \frac{1}{2} \arg \min_b \sum_{i=1}^n f(y_i, b),$$

where

$$f(y, b) = \left(\frac{y}{b} - 1\right)^2.$$

Note that

$$\frac{\partial f(y, b)}{\partial b} = -\frac{2y^2}{b^3} + \frac{2y}{b^2}, \quad \frac{\partial^2 f(y, b)}{\partial b^2} = \frac{6y^2}{b^4} - \frac{4y}{b^3}.$$

The FOC is

$$\sum_{i=1}^n \frac{\partial f(y_i, \hat{b})}{\partial b} = 0,$$

which implies

$$\hat{b} = \frac{\overline{y^2}}{\bar{y}},$$

and the estimate is

$$\hat{\beta}_3 = \frac{\hat{b}}{2} = \frac{1}{2} \frac{\overline{y^2}}{\bar{y}} \xrightarrow{p} \frac{1}{2} \frac{\mathbb{E}[y^2]}{\mathbb{E}[y]} = \beta.$$

To find the asymptotic variance calculate

$$\mathbb{E} \left[ \left( \frac{\partial f(y, 2\beta)}{\partial b} \right)^2 \right] = \frac{\kappa - \beta^4}{16\beta^6}, \quad \mathbb{E} \left[ \frac{\partial^2 f(y, 2\beta)}{\partial b^2} \right] = \frac{1}{4\beta^2}.$$

The derivatives are taken at point  $b = 2\beta$  because  $2\beta$ , and not  $\beta$ , is the solution of the extremum problem  $\min_b \mathbb{E}[f(y, b)]$ . As follows from our discussion,

$$\sqrt{n}(\hat{b} - 2\beta) \xrightarrow{d} \mathcal{N} \left( 0, \frac{\kappa - \beta^4}{\beta^2} \right) \Rightarrow \sqrt{n}(\hat{\beta}_3 - \beta) \xrightarrow{d} \mathcal{N} \left( 0, \frac{\kappa - \beta^4}{4\beta^2} \right).$$

A safer way to obtain this asymptotics is probably to change variable in the minimization problem from the beginning:

$$\hat{\beta}_3 = \arg \min_b \sum_{i=1}^n \left( \frac{y}{2b} - 1 \right)^2,$$

and proceed as above.

No one of these estimators is *a priori* asymptotically better than the others. The idea behind these estimators is:  $\hat{\beta}_1$  is just the usual OLS estimator,  $\hat{\beta}_2$  is the ML estimator using conditional normality:  $y|x \sim \mathcal{N}(\beta, \beta^2)$ . The third estimator may be thought of as the WNLLS estimator with the conditional variance function  $\sigma^2(x, b) = b^2$ , though it is not exactly that: one should divide by  $\sigma^2(x, \beta)$  in the construction of WNLLS.

## 9.2 Quadratic regression

Note that we have conditional homoskedasticity. The regression function is  $g(x, \beta) = (\beta + x)^2$ . The estimator  $\hat{\beta}$  is NLLS, with  $\frac{\partial g(x, \beta)}{\partial \beta} = 2(\beta + x)$ . Then  $Q_{gg} = \mathbb{E} \left[ (\partial g(x, 0) / \partial \beta)^2 \right] = \frac{28}{3}$ . Therefore,  $\sqrt{n}\hat{\beta} \xrightarrow{d} \mathcal{N}(0, \frac{3}{28}\sigma_0^2)$ .

The estimator  $\tilde{\beta}$  is an extremum one, with

$$h(x, y, \beta) = -\frac{y}{(\beta + x)^2} - \ln(\beta + x)^2.$$

First we check the identification condition. Indeed,

$$\frac{\partial h(x, y, \beta)}{\partial \beta} = \frac{2y}{(\beta + x)^3} - \frac{2}{\beta + x},$$

so the FOC to the population problem is

$$\mathbb{E} \left[ \frac{\partial h(x, y, \beta)}{\partial \beta} \right] = -2\beta \mathbb{E} \left[ \frac{\beta + 2x}{(\beta + x)^3} \right],$$

which equals zero if and only if  $\beta = 0$ . It can be verified that the Hessian is negative on all  $\mathbb{B}$ , hence we have a global maximum. Note that the ID condition would not be satisfied if the true parameter was different from zero. Thus,  $\tilde{\beta}$  works only for  $\beta_0 = 0$ .

Next,

$$\frac{\partial^2 h(x, y, \beta)}{\partial \beta^2} = -\frac{6Y}{(\beta + x)^4} + \frac{2}{(\beta + x)^2}.$$

Then

$$\Sigma = \mathbb{E} \left[ \left( \frac{2y}{x^3} - \frac{2}{x} \right)^2 \right] = \frac{31}{40} \sigma_0^2, \quad \Omega = \mathbb{E} \left[ -\frac{6y}{x^4} + \frac{2}{x^2} \right] = -2.$$

Therefore,  $\sqrt{n}\tilde{\beta} \xrightarrow{d} \mathcal{N}(0, \frac{31}{160}\sigma_0^2)$ .

We can see that  $\hat{\beta}$  asymptotically dominates  $\tilde{\beta}$ .

### 9.3 Nonlinearity at left hand side

1. The FOCs for the NLLS problem are

$$\begin{aligned} 0 &= \frac{\partial \sum_{i=1}^n \left( (y_i + \hat{\alpha})^2 - \hat{\beta}x_i \right)^2}{\partial a} = 4 \sum_{i=1}^n \left( (y_i + \hat{\alpha})^2 - \hat{\beta}x_i \right) (y_i + \hat{\alpha}), \\ 0 &= \frac{\partial \sum_{i=1}^n \left( (y_i + \hat{\alpha})^2 - \hat{\beta}x_i \right)^2}{\partial b} = -2 \sum_{i=1}^n \left( (y_i + \hat{\alpha})^2 - \hat{\beta}x_i \right) x_i. \end{aligned}$$

Consider the first of these. The associated population analog is

$$0 = \mathbb{E} [e(y + \alpha)],$$

and it does not follow from the model structure. The model implies that any function of  $x$  is uncorrelated with the error  $e$ , but  $y + \alpha = \pm\sqrt{\beta x + e}$  is generally correlated with  $e$ . The invalidity of population conditions on which the estimator is based leads to inconsistency. The model differs from a nonlinear regression in that the derivative of  $e$  with respect to parameters is not only a function of  $x$ , the conditioning variable, but also of  $y$ , while in a nonlinear regression it *is* (it equals minus the quasi-regressor).

2. Let us select a just identifying set of instruments which one would use if the left hand side of the equation was not squared. This corresponds to the use the following moment conditions implied by the model:

$$0 = \mathbb{E} \left[ \begin{pmatrix} e \\ ex \end{pmatrix} \right].$$

The corresponding CMM estimator results from applying the analogy principle:

$$0 = \sum_{i=1}^n \begin{pmatrix} (y_i + \tilde{\alpha})^2 - \tilde{\beta}x_i \\ ((y_i + \tilde{\alpha})^2 - \tilde{\beta}x_i) x_i \end{pmatrix}.$$

According to the GMM asymptotic theory,  $(\tilde{\alpha}, \tilde{\beta})'$  is consistent for  $(\alpha, \beta)'$  and asymptotically normal with asymptotic variance

$$V = \sigma^2 \begin{pmatrix} 2(\mathbb{E}[y] + \alpha) & -\mathbb{E}[x] \\ 2(\mathbb{E}[yx] + \alpha\mathbb{E}[x]) & -\mathbb{E}[x^2] \end{pmatrix}^{-1} \begin{pmatrix} 1 & \mathbb{E}[x] \\ \mathbb{E}[x] & \mathbb{E}[x^2] \end{pmatrix} \\ \times \begin{pmatrix} 2(\mathbb{E}[y] + \alpha) & 2(\mathbb{E}[yx] + \alpha\mathbb{E}[x]) \\ -\mathbb{E}[x] & -\mathbb{E}[x^2] \end{pmatrix}^{-1}.$$

## 9.4 Least fourth powers

Consider the population level objective function

$$\begin{aligned} \mathbb{E}[(y - bx)^4] &= \mathbb{E}[(e + (\beta - b)x)^4] \\ &= \mathbb{E}[e^4 + 4e^3(\beta - b)x + 6e^2(\beta - b)^2x^2 + 4e(\beta - b)^3x^3 + (\beta - b)^4x^4] \\ &= \mathbb{E}[e^4] + 6(\beta - b)^2\mathbb{E}[e^2x^2] + (\beta - b)^4\mathbb{E}[x^4], \end{aligned}$$

where some of the terms disappear because of the independence of  $x$  and  $e$  and symmetry of the distribution of  $e$ . The last two terms in the objective function are nonnegative, and are zero if and only if (assuming that  $x$  has a nondegenerate distribution)  $b = \beta$ . Thus the (global) ID condition is satisfied.

The squared “score” and second derivative are

$$\left( \frac{\partial (y - bx)^4}{\partial b} \Big|_{b=\beta} \right)^2 = 16e^6x^2, \quad \frac{\partial^2 (y - bx)^4}{\partial b^2} \Big|_{b=\beta} = 12e^2x^2,$$

with expectations  $16\mathbb{E}[e^6]\mathbb{E}[x^2]$  and  $12\mathbb{E}[e^2]\mathbb{E}[x^2]$ . According to the properties of extremum estimators,  $\hat{\beta}$  is consistent and asymptotically normally distributed with asymptotic variance

$$V_{\hat{\beta}} = (12\mathbb{E}[e^2]\mathbb{E}[x^2])^{-1} \cdot 16\mathbb{E}[e^6]\mathbb{E}[x^2] \cdot (12\mathbb{E}[e^2]\mathbb{E}[x^2])^{-1} = \frac{1}{9} \frac{\mathbb{E}[e^6]}{(\mathbb{E}[e^2])^2} \frac{1}{\mathbb{E}[x^2]}.$$

When  $x$  and  $e$  are normally distributed,

$$V_{\hat{\beta}} = \frac{5}{3} \frac{\sigma_e^2}{\sigma_x^2}.$$

The OLS estimator is also consistent and asymptotically normally distributed with asymptotic variance (note that there is conditional homoskedasticity)

$$V_{OLS} = \frac{\mathbb{E}[e^2]}{\mathbb{E}[x^2]}.$$

When  $x$  and  $e$  are normally distributed,

$$V_{OLS} = \frac{\sigma_e^2}{\sigma_x^2},$$

which is smaller than  $V_{\hat{\beta}}$ .

## 9.5 Asymmetric loss

We first need to make sure that we are consistently estimating the right thing. Assume conveniently that  $\mathbb{E}[e] = 0$  to fix the scale of  $\alpha$ . Let  $F$  and  $f$  denote the CDF and PDF of  $e$ , respectively. Assume that these are continuous. Note that

$$\begin{aligned} (y - a - x'b)^3 &= (e + \alpha - a + x'(\beta - b))^3 \\ &= e^3 + 3e^2(\alpha - a + x'(\beta - b)) \\ &\quad + 3e(\alpha - a + x'(\beta - b))^2 + (\alpha - a + x'(\beta - b))^3. \end{aligned}$$

Now,

$$\begin{aligned} \mathbb{E}[\rho(y - a - x'b)] &= \left( \begin{array}{l} \gamma \mathbb{E}[(y - a - x'b)^3 | y - a - x'b \geq 0] \Pr\{y - a - x'b \geq 0\} \\ -(1 - \gamma) \mathbb{E}[(y - a - x'b)^3 | y - a - x'b < 0] \Pr\{y - a - x'b < 0\} \end{array} \right) \\ &= \gamma \int dF_x \int_{e + \alpha - a + x'(\beta - b) \geq 0} \left( \begin{array}{l} e^3 + 3e^2(\alpha - a + x'(\beta - b)) \\ + 3e(\alpha - a + x'(\beta - b))^2 \\ + (\alpha - a + x'(\beta - b))^3 \end{array} \right) dF_{e|x} \\ &\quad \times (1 - \mathbb{E}[F(-(\alpha - a) - x'(\beta - b))]) \\ &\quad - (1 - \gamma) \int dF_x \int_{e + \alpha - a + x'(\beta - b) < 0} \left( \begin{array}{l} e^3 + 3e^2(\alpha - a + x'(\beta - b)) \\ + 3e(\alpha - a + x'(\beta - b))^2 \\ + (\alpha - a + x'(\beta - b))^3 \end{array} \right) dF_{e|x} \\ &\quad \times \mathbb{E}[F(-(\alpha - a) - x'(\beta - b))]. \end{aligned}$$

Is this minimized at  $\alpha$  and  $\beta$ ? The question about global minimum is very hard to answer. Let us restrict ourselves to the local optimum analysis. Take the derivatives and evaluate them at  $\alpha$  and  $\beta$ :

$$\left. \frac{\partial \mathbb{E}[\rho(y - a - x'b)]}{\partial \begin{pmatrix} a \\ b \end{pmatrix}} \right|_{\alpha, \beta} = 3(-\gamma \mathbb{E}[e^2 | e \geq 0] (1 - F(0)) + (1 - \gamma) \mathbb{E}[e^2 | e < 0] F(0)) \begin{pmatrix} 1 \\ \mathbb{E}[x] \end{pmatrix},$$

where we used that infinitesimal change of  $a$  and  $b$  around  $\alpha$  and  $\beta$  does not change the sign of

$$e + \alpha - a + x'(\beta - b).$$

For consistency, we need these derivatives to be zero. This holds if the expression in round brackets is zero, which is true when

$$\frac{\mathbb{E}[e^2 | e < 0]}{\mathbb{E}[e^2 | e \geq 0]} = \frac{1 - F(0)}{F(0)} \cdot \frac{\gamma}{1 - \gamma}.$$

When there is consistency, the asymptotic normality follows from the theory of extremum estimators. Because

$$\begin{aligned} \partial \rho(u) / \partial u &= 3u^2 \begin{cases} \gamma & \text{from the right,} \\ -(1 - \gamma) & \text{from the left,} \end{cases} \\ \partial^2 \rho(u) / \partial u^2 &= 6u \begin{cases} \gamma & \text{from the right,} \\ -(1 - \gamma) & \text{from the left,} \end{cases} \end{aligned}$$

the expected derivatives of the extremum function are

$$\begin{aligned}
\mathbb{E} [h_{\theta} h'_{\theta}] &= \mathbb{E} \left[ \frac{\partial \rho(y - a - x'b)}{\partial \begin{pmatrix} a \\ b \end{pmatrix}} \frac{\partial \rho(y - a - x'b)}{\partial \begin{pmatrix} a \\ b \end{pmatrix}'} \right]_{\alpha, \beta} \\
&= 9\gamma^2 \mathbb{E} \left[ e^4 \begin{pmatrix} 1 \\ x \end{pmatrix} \begin{pmatrix} 1 \\ x \end{pmatrix}' \mid e \geq 0 \right] \Pr\{e \geq 0\} + 9(1 - \gamma)^2 \mathbb{E} \left[ e^4 \begin{pmatrix} 1 \\ x \end{pmatrix} \begin{pmatrix} 1 \\ x \end{pmatrix}' \mid e < 0 \right] \Pr\{e < 0\} \\
&= 9\mathbb{E} \left[ \begin{pmatrix} 1 \\ x \end{pmatrix} \begin{pmatrix} 1 \\ x \end{pmatrix}' \right] (\gamma^2 \mathbb{E} [e^4 \mid e \geq 0] (1 - F(0)) + (1 - \gamma)^2 \mathbb{E} [e^4 \mid e < 0] F(0)),
\end{aligned}$$

$$\begin{aligned}
\mathbb{E} [h_{\theta\theta}] &= \mathbb{E} \left[ \frac{\partial^2 \rho(y - a - x'b)}{\partial \begin{pmatrix} a \\ b \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix}'} \right]_{\alpha, \beta} \\
&= 6\gamma \mathbb{E} \left[ e \begin{pmatrix} 1 \\ x \end{pmatrix} \begin{pmatrix} 1 \\ x \end{pmatrix}' \mid e \geq 0 \right] \Pr\{e \geq 0\} - 6(1 - \gamma) \mathbb{E} \left[ e \begin{pmatrix} 1 \\ x \end{pmatrix} \begin{pmatrix} 1 \\ x \end{pmatrix}' \mid e < 0 \right] \Pr\{e < 0\} \\
&= 6\mathbb{E} \left[ \begin{pmatrix} 1 \\ x \end{pmatrix} \begin{pmatrix} 1 \\ x \end{pmatrix}' \right] (\gamma \mathbb{E} [e \mid e \geq 0] (1 - F(0)) - (1 - \gamma) \mathbb{E} [e \mid e < 0] F(0)).
\end{aligned}$$

If the last expression in round brackets is non-zero, and no element of  $x$  is deterministically constant, then the local ID condition is satisfied. The formula for the asymptotic variance easily follows.



# 10. MAXIMUM LIKELIHOOD ESTIMATION

## 10.1 Normal distribution

Since  $x_1, \dots, x_n$  are from  $\mathcal{N}(\mu, \mu^2)$ , the loglikelihood function is

$$\ell_n = \text{const} - n \ln |\mu| - \frac{1}{2\mu^2} \sum_{i=1}^n (x_i - \mu)^2 = \text{const} - n \ln |\mu| - \frac{1}{2\mu^2} \left( \sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i + n\mu^2 \right).$$

The equation for the ML estimator is  $\mu^2 + \bar{x}\mu - \bar{x}^2 = 0$ . The equation has two solutions  $\mu_1 > 0$ ,  $\mu_2 < 0$ :

$$\mu_1 = \frac{1}{2} \left( -\bar{x} + \sqrt{\bar{x}^2 + 4\bar{x}^2} \right), \quad \mu_2 = \frac{1}{2} \left( -\bar{x} - \sqrt{\bar{x}^2 + 4\bar{x}^2} \right).$$

Note that  $\ell_n$  is a symmetric function of  $\mu$  except for the term  $\mu^{-1} \sum_{i=1}^n x_i$ . This term determines the solution. If  $\bar{x} > 0$  then the global maximum of  $\ell_n$  will be at  $\mu_1$ , otherwise at  $\mu_2$ . That is, the ML estimator is

$$\hat{\mu}_{ML} = \frac{1}{2} \left( -\bar{x} + \text{sgn}(\bar{x}) \sqrt{\bar{x}^2 + 4\bar{x}^2} \right).$$

It is consistent because, if  $\mu \neq 0$ ,  $\text{sgn}(\bar{x}) \xrightarrow{p} \text{sgn}(\mu)$  and

$$\hat{\mu}_{ML} \xrightarrow{p} \frac{1}{2} \left( -\mathbb{E}[x] + \text{sgn}(\mathbb{E}[x]) \sqrt{(\mathbb{E}[x])^2 + 4\mathbb{E}[x^2]} \right) = \frac{1}{2} \left( -\mu + \text{sgn}(\mu) \sqrt{\mu^2 + 8\mu^2} \right) = \mu.$$

## 10.2 Pareto distribution

The loglikelihood function is

$$\ell_n = n \ln \lambda - (\lambda + 1) \sum_{i=1}^n \ln x_i.$$

(a) The ML estimator  $\hat{\lambda}$  of  $\lambda$  is the solution of  $\partial \ell_n / \partial \lambda = 0$ . That is,  $\hat{\lambda}_{ML} = 1/\overline{\ln x}$ , which is consistent for  $\lambda$ , because  $1/\overline{\ln x} \xrightarrow{p} 1/\mathbb{E}[\ln x] = \lambda$ . The asymptotics is  $\sqrt{n} (\hat{\lambda}_{ML} - \lambda) \xrightarrow{d} \mathcal{N}(0, \mathcal{I}^{-1})$ , where the information matrix is  $\mathcal{I} = -\mathbb{E}[\partial s / \partial \lambda] = -\mathbb{E}[-1/\lambda^2] = 1/\lambda^2$ .

(b) The Wald test for a simple hypothesis is

$$\mathcal{W} = n(\hat{\lambda} - \lambda)' \mathcal{I}(\hat{\lambda})(\hat{\lambda} - \lambda) = n \frac{(\hat{\lambda} - \lambda_0)^2}{\hat{\lambda}^2} \xrightarrow{d} \chi_1^2.$$

The Likelihood Ratio test statistic for a simple hypothesis is

$$\begin{aligned}\mathcal{LR} &= 2 \left( \ell_n(\hat{\lambda}) - \ell_n(\lambda_0) \right) \\ &= 2 \left( n \ln \hat{\lambda} - (\hat{\lambda} + 1) \sum_{i=1}^n \ln x_i - n \ln \lambda_0 + (\lambda_0 + 1) \sum_{i=1}^n \ln x_i \right) \\ &= 2 \left( n \ln \frac{\hat{\lambda}}{\lambda_0} - (\hat{\lambda} - \lambda_0) \sum_{i=1}^n \ln x_i \right) \xrightarrow{d} \chi_1^2.\end{aligned}$$

The Lagrange Multiplier test statistic for a simple hypothesis is

$$\begin{aligned}\mathcal{LM} &= \frac{1}{n} \sum_{i=1}^n s(x_i, \lambda_0)' \mathcal{I}(\lambda_0)^{-1} \sum_{i=1}^n s(x_i, \lambda_0) = \frac{1}{n} \left[ \sum_{i=1}^n \left( \frac{1}{\lambda_0} - \ln x_i \right) \right]^2 \lambda_0^2 \\ &= n \frac{(\hat{\lambda} - \lambda_0)^2}{\hat{\lambda}^2} \xrightarrow{d} \chi_1^2.\end{aligned}$$

Observe that  $\mathcal{W}$  and  $\mathcal{LM}$  are numerically equal.

## 10.3 Comparison of ML tests

Recall that for the ML estimator  $\hat{\lambda}$  and the simple hypothesis  $H_0 : \lambda = \lambda_0$ ,

$$\mathcal{W} = n(\hat{\lambda} - \lambda_0)' \mathcal{I}(\hat{\lambda})(\hat{\lambda} - \lambda_0)$$

and

$$\mathcal{LM} = \frac{1}{n} \sum_i s(x_i, \lambda_0)' \mathcal{I}(\lambda_0)^{-1} \sum_i s(x_i, \lambda_0).$$

1. The density of a Poisson distribution with parameter  $\lambda$  is

$$f(x|\lambda) = \frac{\lambda^x}{x!} e^{-\lambda},$$

so  $\hat{\lambda}_{ML} = \bar{x}$ ,  $\mathcal{I}(\lambda) = 1/\lambda$ . For the simple hypothesis with  $\lambda_0 = 3$  the test statistics are

$$\mathcal{W} = \frac{n(\bar{x} - 3)^2}{\bar{x}}, \quad \mathcal{LM} = \frac{1}{n} \left( \sum_i \frac{x_i}{3} - n \right)^2 3 = \frac{n(\bar{x} - 3)^2}{3},$$

and  $\mathcal{W} \geq \mathcal{LM}$  for  $\bar{x} \leq 3$  and  $\mathcal{W} \leq \mathcal{LM}$  for  $\bar{x} \geq 3$ .

2. The density of an exponential distribution with parameter  $\theta$  is

$$f(x|\theta) = \frac{1}{\theta} e^{-\frac{x}{\theta}},$$

so  $\hat{\theta}_{ML} = \bar{x}$ ,  $\mathcal{I}(\theta) = 1/\theta^2$ . For the simple hypothesis with  $\theta_0 = 3$  the test statistics are

$$\mathcal{W} = \frac{n(\bar{x} - 3)^2}{\bar{x}^2}, \quad \mathcal{LM} = \frac{1}{n} \left( \sum_i \frac{x_i}{3^2} - \frac{n}{3} \right)^2 3^2 = \frac{n(\bar{x} - 3)^2}{9},$$

and  $\mathcal{W} \geq \mathcal{LM}$  for  $0 < \bar{x} \leq 3$  and  $\mathcal{W} \leq \mathcal{LM}$  for  $\bar{x} \geq 3$ .

3. The density of a Bernoulli distribution with parameter  $\theta$  is

$$f(x|\theta) = \theta^x(1 - \theta)^{1-x},$$

so  $\hat{\theta}_{ML} = \bar{x}$ ,  $\mathcal{I}(\theta) = \theta^{-1}(1 - \theta)^{-1}$ . For the simple hypothesis with  $\theta_0 = \frac{1}{2}$  the test statistics are

$$\mathcal{W} = n \frac{(\bar{x} - \frac{1}{2})^2}{\bar{x}(1 - \bar{x})}, \quad \mathcal{LM} = \frac{1}{n} \left( \frac{\sum_i x_i}{\frac{1}{2}} - \frac{n - \sum_i x_i}{\frac{1}{2}} \right)^2 \frac{1}{2} \frac{1}{2} = 4n \left( \bar{x} - \frac{1}{2} \right)^2,$$

and  $\mathcal{W} \geq \mathcal{LM}$  (since  $\bar{x}(1 - \bar{x}) \leq 1/4$ ). For the simple hypothesis with  $\theta_0 = \frac{2}{3}$  the test statistics are

$$\mathcal{W} = n \frac{(\bar{x} - \frac{2}{3})^2}{\bar{x}(1 - \bar{x})}, \quad \mathcal{LM} = \frac{1}{n} \left( \frac{\sum_i x_i}{\frac{2}{3}} - \frac{n - \sum_i x_i}{\frac{1}{3}} \right)^2 \frac{2}{3} \frac{1}{3} = \frac{9}{2} n \left( \bar{x} - \frac{2}{3} \right)^2,$$

therefore  $\mathcal{W} \leq \mathcal{LM}$  when  $2/9 \leq \bar{x}(1 - \bar{x})$  and  $\mathcal{W} \geq \mathcal{LM}$  when  $2/9 \geq \bar{x}(1 - \bar{x})$ . Equivalently,  $\mathcal{W} \leq \mathcal{LM}$  for  $\frac{1}{3} \leq \bar{x} \leq \frac{2}{3}$  and  $\mathcal{W} \geq \mathcal{LM}$  for  $0 < \bar{x} \leq \frac{1}{3}$  or  $\frac{2}{3} \leq \bar{x} \leq 1$ .

## 10.4 Invariance of ML tests to reparameterizations of null

1. Denote by  $\Theta_0$  the set of  $\theta$ 's that satisfy the null. Since  $f$  is one-to-one,  $\Theta_0$  is the same under both parameterizations of the null. Then the restricted and unrestricted ML estimators are invariant to how  $H_0$  is formulated, and so is the  $\mathcal{LR}$  statistic.
2. When  $f$  is linear,  $f(h(\theta)) - f(0) = Fh(\theta) - F0 = Fh(\theta)$ , and the matrix of derivatives of  $h$  translates linearly into the matrix of derivatives of  $g$ :  $G = FH$ , where  $F = \partial f(x)/\partial x'$  does not depend on its argument  $x$ , and thus need not be estimated. Then

$$\begin{aligned} \mathcal{W}_g &= n \left( \frac{1}{n} \sum_{i=1}^n g(\hat{\theta}) \right)' (G \hat{V}_\theta G')^{-1} \left( \frac{1}{n} \sum_{i=1}^n g(\hat{\theta}) \right) \\ &= n \left( \frac{1}{n} \sum_{i=1}^n Fh(\hat{\theta}) \right)' (FH \hat{V}_\theta H' F')^{-1} \left( \frac{1}{n} \sum_{i=1}^n Fh(\hat{\theta}) \right) \\ &= n \left( \frac{1}{n} \sum_{i=1}^n h(\hat{\theta}) \right)' (H \hat{V}_\theta H')^{-1} \left( \frac{1}{n} \sum_{i=1}^n h(\hat{\theta}) \right) = \mathcal{W}_h, \end{aligned}$$

but this sequence of equalities does not work when  $f$  is nonlinear.

3. The  $\mathcal{W}$  statistic for the reparameterized null equals

$$\begin{aligned} \mathcal{W} &= \frac{n \left( \frac{\hat{\theta}_1 - \alpha}{\hat{\theta}_2 - \alpha} - 1 \right)^2}{\left( \begin{array}{c} 1 \\ \frac{\hat{\theta}_2 - \alpha}{\hat{\theta}_1 - \alpha} \\ -\frac{1}{(\hat{\theta}_2 - \alpha)^2} \end{array} \right)' \left( \begin{array}{cc} i_{11} & i_{12} \\ i_{12} & i_{22} \end{array} \right)^{-1} \left( \begin{array}{c} 1 \\ \frac{\hat{\theta}_2 - \alpha}{\hat{\theta}_1 - \alpha} \\ -\frac{1}{(\hat{\theta}_2 - \alpha)^2} \end{array} \right)} \\ &= \frac{n (\hat{\theta}_1 - \hat{\theta}_2)^2}{i^{11} - 2i^{12} \frac{\hat{\theta}_1 - \alpha}{\hat{\theta}_2 - \alpha} + i^{22} \left( \frac{\hat{\theta}_1 - \alpha}{\hat{\theta}_2 - \alpha} \right)^2}, \end{aligned}$$

where

$$\widehat{\mathcal{I}} = \begin{pmatrix} i_{11} & i_{12} \\ i_{12} & i_{22} \end{pmatrix}, \quad \widehat{\mathcal{I}}^{-1} = \begin{pmatrix} i^{11} & i^{12} \\ i^{12} & i^{22} \end{pmatrix}.$$

By choosing  $\alpha$  close to  $\hat{\theta}_2$ , we can make  $\mathcal{W}$  as close to zero as desired. The value of  $\alpha$  equal to  $(\hat{\theta}_1 - \hat{\theta}_2 i^{12}/i^{22})/(1 - i^{12}/i^{22})$  gives the largest possible value to the  $\mathcal{W}$  statistic equal to

$$\frac{n (\hat{\theta}_1 - \hat{\theta}_2)^2}{i^{11} - (i^{12})^2/i^{22}}.$$

## 10.5 Misspecified maximum likelihood

1. *Method 1.* It is straightforward to derive the loglikelihood function and see that the problem of its maximization implies minimization of the sum of squares of deviations of  $y$  from  $g(x, b)$  over  $b$ , i.e. the NLLS problem. But we know that the NLLS estimator is consistent.

*Method 2.* It is straightforward to see that the population analog of the FOC for the ML problem is that the expected product of quasi-regressor and deviation of  $y$  from  $g(x, \beta)$  equals zero, but this system of moment conditions follows from the regression model.

2. By construction, it is an extremum estimator. It will be consistent for the value that solves the analogous extremum problem in population:

$$\hat{\theta} \xrightarrow{P} \theta^* \equiv \arg \max_{q \in \Theta} \mathbb{E}[f(z|q)],$$

provided that this  $\theta^*$  is unique (if it is not unique, no nice asymptotic properties are expected). It is unlikely that this limit will be at true  $\theta$ . As an extremum estimator,  $\hat{\theta}$  will be asymptotically normal, although centered around wrong value of the parameter:

$$\sqrt{n} (\hat{\theta} - \theta^*) \xrightarrow{d} \mathcal{N}(0, V_{\hat{\theta}}).$$

3. The parameter vector is  $q = (b', c)'$ , where  $c$  enters the assumed form of  $\sigma_t^2$ . The conditional density is

$$\log f(y_t | I_{t-1}, b, c) = -\frac{1}{2} \log \sigma_t^2(b, c) - \frac{(y_t - x_t' b)^2}{2\sigma_t^2(b, c)}.$$

The conditional score is

$$s(y_t | I_{t-1}, b, c) = \begin{pmatrix} \frac{(y_t - x_t' b) x_t}{\sigma_t^2(b, c)} + \frac{1}{2\sigma_t^2(b, c)} \left( \frac{(y_t - x_t' b)^2}{\sigma_t^2(b, c)} - 1 \right) \frac{\partial \sigma_t^2(b, c)}{\partial b} \\ \frac{1}{2\sigma_t^2(b, c)} \left( \frac{(y_t - x_t' b)^2}{\sigma_t^2(b, c)} - 1 \right) \frac{\partial \sigma_t^2(b, c)}{\partial c} \end{pmatrix}.$$

The ML estimator will estimate such  $b$  and  $c$  that make  $\mathbb{E}[s(y_t | I_{t-1}, b, c)] = 0$ . Will this  $b$  be equal to true  $\beta$ ? It is easy to see that if we evaluate the conditional score at  $b = \beta$ , its expectation will be

$$\mathbb{E} \left( \begin{pmatrix} \frac{1}{2\sigma_t^2(\beta, c)} \left( \frac{\mathbb{V}[y_t | I_{t-1}]}{\sigma_t^2(\beta, c)} - 1 \right) \frac{\partial \sigma_t^2(\beta, c)}{\partial b} \\ \frac{1}{2\sigma_t^2(\beta, c)} \left( \frac{\mathbb{V}[y_t | I_{t-1}]}{\sigma_t^2(\beta, c)} - 1 \right) \frac{\partial \sigma_t^2(\beta, c)}{\partial c} \end{pmatrix} \right).$$

This may not be zero when  $\sigma_t^2(\beta, c)$  does not coincide with  $\mathbb{V}[y_t|I_{t-1}]$  for any  $c$ , so it does not necessarily follow that the ML estimator of  $\beta$  will be consistent. However, if  $\sigma_t^2(\beta, c)$  does not depend on  $\beta$ , it is clear that the  $\beta$ -part of this will be zero, and the  $c$ -part will be zero too for some  $c$  (this condition will define the probability limit of  $\hat{c}$ , the ML estimator of  $c$ ). In this case, the ML estimator of  $\beta$  will be consistent. In fact, it is easy to see that it will equal to a weighted least squares estimator with weights  $\sigma_t^2(\hat{c})$

$$\hat{\beta} = \left( \sum_t \frac{x_t x_t'}{\sigma_t^2(\hat{c})} \right)^{-1} \sum_t \frac{x_t y_t}{\sigma_t^2(\hat{c})}.$$

## 10.6 Individual effects

The loglikelihood is

$$\ell_n(\mu_1, \dots, \mu_n, \sigma^2) = \text{const} - n \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n \{(x_i - \mu_i)^2 + (y_i - \mu_i)^2\}.$$

FOC give

$$\hat{\mu}_{iML} = \frac{x_i + y_i}{2}, \quad \sigma_{ML}^2 = \frac{1}{2n} \sum_{i=1}^n \{(x_i - \hat{\mu}_{iML})^2 + (y_i - \hat{\mu}_{iML})^2\},$$

so that

$$\hat{\sigma}_{ML}^2 = \frac{1}{4n} \sum_{i=1}^n (x_i - y_i)^2.$$

Because

$$\hat{\sigma}_{ML}^2 = \frac{1}{4n} \sum_{i=1}^n \{(x_i - \mu_i)^2 + (y_i - \mu_i)^2 - 2(x_i - \mu_i)(y_i - \mu_i)\} \xrightarrow{p} \frac{\sigma^2}{4} + \frac{\sigma^2}{4} - 0 = \frac{\sigma^2}{2},$$

the ML estimator is inconsistent. Why? The maximum likelihood method presumes a parameter vector of fixed dimension. In our case the dimension instead increases with the number of observations. The information from new observations goes to estimation of new parameters instead of enhancing precision of the already existing ones. To construct a consistent estimator, just multiply  $\hat{\sigma}_{ML}^2$  by 2. There are also other possibilities.

## 10.7 Irregular confidence interval

The maximum likelihood estimator of  $\theta$  is

$$\hat{\theta}_{ML} = x_{(n)} \equiv \max\{x_1, \dots, x_n\}$$

whose asymptotic distribution is exponential:

$$F_{n(\hat{\theta}_{ML} - \theta)}(t) \rightarrow \exp(t/\theta) \cdot \mathbb{I}_{\{t \leq 0\}} + \mathbb{I}_{\{t > 0\}}.$$

The most elegant way to proceed is by pivoting this distribution first:

$$F_{n(\hat{\theta}_{ML} - \theta)/\theta}(t) \rightarrow \exp(t) \cdot \mathbb{I}_{\{t \leq 0\}} + \mathbb{I}_{\{t > 0\}}.$$

The left 5%-quantile for the limiting distribution is  $\log(.05)$ . Thus, with probability 95%,  $\log(.05) \leq n(\hat{\theta}_{ML} - \theta)/\theta \leq 0$ , so the confidence interval for  $\theta$  is

$$[x_{(n)}, (1 + \log(.05)/n)^{-1}x_{(n)}].$$

## 10.8 Trivial parameter space

Since the parameter space contains only one point, the latter is the optimizer. If  $\theta_1 = \theta_0$ , then the estimator  $\hat{\theta}_{ML} = \theta_1$  is consistent for  $\theta_0$  and has infinite rate of convergence. If  $\theta_1 \neq \theta_0$ , then the ML estimator is inconsistent.

## 10.9 Nuisance parameter in density

We need to apply the Taylor's expansion twice, i.e. for both stages of estimation.

The FOC for the second stage of estimation is

$$\frac{1}{n} \sum_{i=1}^n s_c(y_i, x_i, \tilde{\gamma}, \hat{\delta}_m) = 0,$$

where  $s_c(y, x, \gamma, \delta) \equiv \frac{\partial \log f_c(y|x, \gamma, \delta)}{\partial \gamma}$  is the conditional score. Taylor's expansion with respect to the  $\gamma$ -argument around  $\gamma_0$  yields

$$\frac{1}{n} \sum_{i=1}^n s_c(y_i, x_i, \gamma_0, \hat{\delta}_m) + \frac{1}{n} \sum_{i=1}^n \frac{\partial s_c(y_i, x_i, \gamma^*, \hat{\delta}_m)}{\partial \gamma'} (\tilde{\gamma} - \gamma_0) = 0,$$

where  $\gamma^*$  lies between  $\tilde{\gamma}$  and  $\gamma_0$  componentwise.

Now Taylor-expand the first term around  $\delta_0$ :

$$\frac{1}{n} \sum_{i=1}^n s_c(y_i, x_i, \gamma_0, \hat{\delta}_m) = \frac{1}{n} \sum_{i=1}^n s_c(y_i, x_i, \gamma_0, \delta_0) + \frac{1}{n} \sum_{i=1}^n \frac{\partial s_c(y_i, x_i, \gamma_0, \delta^*)}{\partial \delta'} (\hat{\delta}_m - \delta_0),$$

where  $\delta^*$  lies between  $\hat{\delta}_m$  and  $\delta_0$  componentwise.

Combining the two pieces, we get:

$$\begin{aligned} \sqrt{n}(\tilde{\gamma} - \gamma_0) &= - \left( \frac{1}{n} \sum_{i=1}^n \frac{\partial s_c(y_i, x_i, \gamma^*, \hat{\delta}_m)}{\partial \gamma'} \right)^{-1} \times \\ &\quad \times \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n s_c(y_i, x_i, \gamma_0, \delta_0) + \frac{1}{n} \sum_{i=1}^n \frac{\partial s_c(y_i, x_i, \gamma_0, \delta^*)}{\partial \delta'} \sqrt{n}(\hat{\delta}_m - \delta_0) \right). \end{aligned}$$

Now let  $n \rightarrow \infty$ . Under ULLN for the second derivative of the log of the conditional density, the first factor converges in probability to  $-(\mathcal{I}_c^{\gamma\gamma})^{-1}$ , where  $\mathcal{I}_c^{\gamma\gamma} \equiv -\mathbb{E} \left[ \frac{\partial^2 \log f_c(y|x, \gamma_0, \delta_0)}{\partial \gamma \partial \gamma'} \right]$ . There

are two terms inside the brackets that have nontrivial distributions. We will compute asymptotic variance of each and asymptotic covariance between them. The first term behaves as follows:

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n s_c(y_i, x_i, \gamma_0, \delta_0) \xrightarrow{d} \mathcal{N}(0, \mathcal{I}_c^{\gamma\gamma})$$

due to the CLT (recall that the score has zero expectation and the information matrix equality). Turn to the second term. Under the ULLN,  $\frac{1}{n} \sum_{i=1}^n \frac{\partial s_c(y_i, x_i, \gamma_0, \delta^*)}{\partial \delta'}$  converges to  $-\mathcal{I}_c^{\gamma\delta} = \mathbb{E} \left[ \frac{\partial^2 \log f_c(y|x, \gamma_0, \delta_0)}{\partial \gamma \partial \delta'} \right]$ . Next, we know from the MLE theory that  $\sqrt{n}(\hat{\delta}_m - \delta_0) \xrightarrow{d} \mathcal{N}(0, (\mathcal{I}_m^{\delta\delta})^{-1})$ , where  $\mathcal{I}_m^{\delta\delta} \equiv -\mathbb{E} \left[ \frac{\partial^2 \log f_m(x|\delta_0)}{\partial \delta \partial \delta'} \right]$ . Finally, the asymptotic covariance term is zero because of the “marginal/conditional” relationship between the two terms, the Law of Iterated Expectations and zero expected score.

Collecting the pieces, we find:

$$\sqrt{n}(\hat{\gamma} - \gamma_0) \xrightarrow{d} \mathcal{N} \left( 0, (\mathcal{I}_c^{\gamma\gamma})^{-1} \left( \mathcal{I}_c^{\gamma\gamma} + \mathcal{I}_c^{\gamma\delta} (\mathcal{I}_m^{\delta\delta})^{-1} \mathcal{I}_c^{\delta\gamma} \right) (\mathcal{I}_c^{\gamma\gamma})^{-1} \right).$$

It is easy to see that the asymptotic variance is larger (in matrix sense) than  $(\mathcal{I}_c^{\gamma\gamma})^{-1}$  that would be the asymptotic variance if we knew the nuisance parameter  $\delta_0$ . But it is impossible to compare to the asymptotic variance for  $\hat{\gamma}_c$ , which is *not*  $(\mathcal{I}_c^{\gamma\gamma})^{-1}$ .

## 10.10 MLE versus OLS

1. Observe that  $\hat{\alpha}_{OLS} = \frac{1}{n} \sum_{i=1}^n y_i$ ,  $\mathbb{E}[\hat{\alpha}_{OLS}] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[y] = \alpha$ , so  $\hat{\alpha}_{OLS}$  is unbiased. Next,  $\frac{1}{n} \sum_{i=1}^n y_i \xrightarrow{p} \mathbb{E}[y] = \alpha$ , so  $\hat{\alpha}_{OLS}$  is consistent. Yes, as we know from the theory,  $\hat{\alpha}_{OLS}$  is the best linear unbiased estimator. Note that the members of this class are allowed to be of the form  $\{\mathcal{A}\mathcal{Y} \text{ s.t. } \mathcal{A}\mathcal{X} = I\}$ , where  $\mathcal{A}$  is a *constant* matrix, since there are no regressors beside the constant. There is no heteroskedasticity, since there are no regressors to condition on (more precisely, we should condition on a constant, i.e. the trivial  $\sigma$ -field, which gives just an unconditional variance which is constant by the random sampling assumption). The asymptotic distribution is

$$\sqrt{n}(\hat{\alpha}_{OLS} - \alpha) = \frac{1}{\sqrt{n}} \sum_{i=1}^n e_i \xrightarrow{d} \mathcal{N}(0, \sigma^2 \mathbb{E}[x^2]),$$

since the variance of  $e$  is  $\mathbb{E}[e^2] = \mathbb{E}[\mathbb{E}[e^2|x]] = \sigma^2 \mathbb{E}[x^2]$ .

2. The conditional likelihood function is

$$\mathcal{L}(y_1, \dots, y_n, x_1, \dots, x_n, \alpha, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi x_i^2 \sigma^2}} \exp \left\{ -\frac{(y_i - \alpha)^2}{2x_i^2 \sigma^2} \right\}.$$

The conditional loglikelihood is

$$\ell_n(y_1, \dots, y_n, x_1, \dots, x_n, \alpha, \sigma^2) = \text{const} - \sum_{i=1}^n \frac{(y_i - \alpha)^2}{2x_i^2 \sigma^2} - \frac{n}{2} \log \sigma^2 \rightarrow \max_{\alpha, \sigma^2}.$$

From the FOC

$$\frac{\partial \ell_n}{\partial \alpha} = \sum_{i=1}^n \frac{y_i - \alpha}{x_i^2 \sigma^2} = 0,$$

the ML estimator is

$$\hat{\alpha}_{ML} = \frac{\sum_{i=1}^n y_i/x_i^2}{\sum_{i=1}^n 1/x_i^2}.$$

Note: it is equal to the OLS estimator of  $\alpha$  in

$$\frac{y_i}{x_i} = \alpha \frac{1}{x_i} + \frac{e_i}{x_i}.$$

The asymptotic distribution is

$$\sqrt{n}(\hat{\alpha}_{ML} - \alpha) = \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n e_i/x_i^2}{\frac{1}{n} \sum_{i=1}^n 1/x_i^2} \xrightarrow{d} \left( \mathbb{E} \left[ \frac{1}{x^2} \right] \right)^{-1} \mathcal{N} \left( 0, \sigma^2 \mathbb{E} \left[ \frac{1}{x^2} \right] \right) = \mathcal{N} \left( 0, \sigma^2 \left( \mathbb{E} \left[ \frac{1}{x^2} \right] \right)^{-1} \right).$$

Note that  $\hat{\alpha}_{ML}$  is unbiased and more efficient than  $\hat{\alpha}_{OLS}$  since

$$\left( \mathbb{E} \left[ \frac{1}{x^2} \right] \right)^{-1} < \mathbb{E} [x^2],$$

but it is not in the class of linear unbiased estimators, since the weights in  $\mathcal{A}_{ML}$  depend on extraneous  $x$ 's. The  $\hat{\alpha}_{ML}$  is efficient in a much larger class. Thus there is no contradiction.

## 10.11 MLE versus GLS

The feasible GLS estimator  $\tilde{\beta}$  is constructed by

$$\tilde{\beta} = \left( \sum_{i=1}^n \frac{x_i x_i'}{(x_i' \hat{\beta})^2} \right)^{-1} \sum_{i=1}^n \frac{x_i y_i}{(x_i' \hat{\beta})^2}.$$

The asymptotic variance matrix is

$$V_{\tilde{\beta}} = \sigma^2 \left( \mathbb{E} \left[ \frac{x x'}{(x' \beta)^2} \right] \right)^{-1}.$$

The conditional logdensity is

$$\ell_n(x, y, b, s^2) = \text{const} - \frac{1}{2} \log s^2 - \frac{1}{2} \log(x'b)^2 - \frac{1}{2s^2} \frac{(y - x'b)^2}{(x'b)^2},$$

so the conditional score is

$$\begin{aligned} s_{\beta}(x, y, b, s^2) &= -\frac{x}{x'b} + \frac{y - x'b}{(x'b)^3} \frac{xy}{s^2}, \\ s_{\sigma^2}(x, y, b, s^2) &= -\frac{1}{2s^2} + \frac{1}{2s^4} \frac{(y - x'b)^2}{(x'b)^2}, \end{aligned}$$



whose derivatives are

$$\begin{aligned} s_{\beta\beta}(x, y, b, s^2) &= \frac{xx'}{(x'b)^2} - (3y - 2x'b) \frac{y}{(x'b)^4} \frac{xx'}{s^2}, \\ s_{\beta\sigma^2}(x, y, b, s^2) &= -\frac{y - x'b}{(x'b)^3} \frac{xy}{s^4}, \\ s_{\sigma^2\sigma^2}(x, y, b, s^2) &= \frac{1}{2s^4} - \frac{1}{s^6} \frac{(y - x'b)^2}{(x'b)^2}. \end{aligned}$$

After taking expectations, we find that the information matrix is

$$\mathcal{I}_{\beta\beta} = \frac{2\sigma^2 + 1}{\sigma^2} \mathbb{E} \left[ \frac{xx'}{(x'\beta)^2} \right], \quad \mathcal{I}_{\beta\sigma^2} = \frac{1}{\sigma^2} \mathbb{E} \left[ \frac{x}{x'\beta} \right], \quad \mathcal{I}_{\sigma^2\sigma^2} = \frac{1}{2\sigma^4}.$$

By inverting a partitioned matrix, we find that the asymptotic variance of the ML estimator of  $\beta$  is

$$V_{ML} = \left( \mathcal{I}_{\beta\beta} - \mathcal{I}_{\beta\sigma^2} \mathcal{I}_{\sigma^2\sigma^2}^{-1} \mathcal{I}'_{\beta\sigma^2} \right)^{-1} = \left( \frac{2\sigma^2 + 1}{\sigma^2} \mathbb{E} \left[ \frac{xx'}{(x'\beta)^2} \right] - 2 \mathbb{E} \left[ \frac{x}{x'\beta} \right] \mathbb{E} \left[ \frac{x'}{x'\beta} \right] \right)^{-1}.$$

Now,

$$V_{ML}^{-1} = \frac{1}{\sigma^2} \mathbb{E} \left[ \frac{xx'}{(x'\beta)^2} \right] + 2 \left( \mathbb{E} \left[ \frac{xx'}{(x'\beta)^2} \right] - \mathbb{E} \left[ \frac{x}{x'\beta} \right] \mathbb{E} \left[ \frac{x'}{x'\beta} \right] \right) \geq \frac{1}{\sigma^2} \mathbb{E} \left[ \frac{xx'}{(x'\beta)^2} \right] = V_{\tilde{\beta}}^{-1},$$

where the inequality follows from  $\mathbb{E}[aa'] - \mathbb{E}[a] \mathbb{E}[a'] = \mathbb{E}[(a - \mathbb{E}[a])(a - \mathbb{E}[a])'] \geq 0$ . Therefore,  $V_{\tilde{\beta}} \geq V_{ML}$ , i.e. the GLS estimator is less asymptotically efficient than the ML estimator. This is because  $\beta$  figures both into the conditional mean and conditional variance, but the GLS estimator ignores this information.

## 10.12 MLE in heteroskedastic time series regression

Observe that the joint density factorizes:

$$\begin{aligned} f(y_t, x_t, y_{t-1}, x_{t-1}, y_{t-2}, x_{t-2}, \dots) &= f^c(y_t|x_t, y_{t-1}, x_{t-1}, y_{t-2}, x_{t-2}, \dots) f^m(x_t|y_{t-1}, x_{t-1}, y_{t-2}, x_{t-2}, \dots) \\ &\quad \times f(y_{t-1}, x_{t-1}, y_{t-2}, x_{t-2}, \dots). \end{aligned}$$

Assume that data  $(y_t, x_t)$ ,  $t = 1, 2, \dots, T$ , are stationary and ergodic and generated by

$$y_t = \alpha + \beta x_t + u_t,$$

where  $u_t|x_t, y_{t-1}, x_{t-1}, y_{t-2}, \dots \sim \mathcal{N}(0, \sigma_t^2)$  and  $x_t|y_{t-1}, x_{t-1}, y_{t-2}, x_{t-2}, \dots \sim \mathcal{N}(0, v)$ . Explain, without going into deep math,

Since the parameter  $v$  is not present the conditional density  $f^c$ , it can be efficiently estimated from the marginal density  $f^m$ , which yields

$$\hat{v} = \frac{1}{T} \sum_{t=1}^T x_t^2.$$

The standard error may be constructed via

$$\hat{V} = \frac{1}{T} \sum_{t=1}^T x_t^4 - \hat{v}^2.$$

1. If the entire function  $\sigma_t^2 = \sigma^2(x_t)$  is fully known, the conditional ML estimator of  $\alpha$  and  $\beta$  is the same as the GLS estimator:

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix}_{ML} = \left( \sum_{t=1}^T \frac{1}{\sigma_t^2} \begin{pmatrix} 1 & x_t \\ x_t & x_t^2 \end{pmatrix} \right)^{-1} \sum_{t=1}^T \frac{y_t}{\sigma_t^2} \begin{pmatrix} 1 \\ x_t \end{pmatrix}.$$

The standard errors may be constructed via

$$\hat{V}_{ML} = T \left( \sum_{t=1}^T \frac{1}{\sigma_t^2} \begin{pmatrix} 1 & x_t \\ x_t & x_t^2 \end{pmatrix} \right)^{-1}.$$

2. If the values of  $\sigma_t^2$  at  $t = 1, 2, \dots, T$  are known, we can use the same procedure as in part 1, since it does not use values of  $\sigma^2(x_t)$  other than those at  $x_1, x_2, \dots, x_T$ .
3. If it is known that  $\sigma_t^2 = (\theta + \delta x_t)^2$ , we have in addition parameters  $\theta$  and  $\delta$  to be estimated jointly from the conditional distribution

$$y_t | x_t, y_{t-1}, x_{t-1}, y_{t-2}, x_{t-2}, \dots \sim \mathcal{N}(\alpha + \beta x_t, (\theta + \delta x_t)^2).$$

The loglikelihood function is

$$\ell_n(\alpha, \beta, \theta, \delta) = \text{const} - \frac{1}{2} \sum_{t=1}^T \log(\theta + \delta x_t)^2 - \frac{1}{2} \sum_{t=1}^T \frac{(y_t - \alpha - \beta x_t)^2}{(\theta + \delta x_t)^2},$$

and  $(\hat{\alpha}, \hat{\beta}, \hat{\theta}, \hat{\delta})'_{ML} = \arg \max_{(\alpha, \beta, \theta, \delta)} \ell_n(\alpha, \beta, \theta, \delta)$ . Note that

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix}_{ML} = \left( \sum_{t=1}^T \frac{1}{(\hat{\theta} + \hat{\delta} x_t)^2} \begin{pmatrix} 1 & x_t \\ x_t & x_t^2 \end{pmatrix} \right)^{-1} \sum_{t=1}^T \frac{y_t}{(\hat{\theta} + \hat{\delta} x_t)^2} \begin{pmatrix} 1 \\ x_t \end{pmatrix}.$$

The standard errors may be constructed via

$$\hat{V}_{ML} = T \left( \sum_{t=1}^T \frac{\partial \ell_n(\hat{\alpha}, \hat{\beta}, \hat{\theta}, \hat{\delta})}{\partial(\alpha, \beta, \theta, \delta)'} \frac{\partial \ell_n(\hat{\alpha}, \hat{\beta}, \hat{\theta}, \hat{\delta})}{\partial(\alpha, \beta, \theta, \delta)} \right)^{-1}.$$

4. Similarly to part 3, if it is known that  $\sigma_t^2 = \theta + \delta u_{t-1}^2$ , we have in addition parameters  $\theta$  and  $\delta$  to be estimated jointly from the conditional distribution

$$y_t | x_t, y_{t-1}, x_{t-1}, y_{t-2}, x_{t-2}, \dots \sim \mathcal{N}(\alpha + \beta x_t, \theta + \delta(y_{t-1} - \alpha - \beta x_{t-1})^2).$$

5. If it is only known that  $\sigma_t^2$  is stationary, conditional maximum likelihood function is unavailable, so we have to use subefficient methods, for example, OLS estimation

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix}_{OLS} = \left( \sum_{t=1}^T \begin{pmatrix} 1 & x_t \\ x_t & x_t^2 \end{pmatrix} \right)^{-1} \sum_{t=1}^T y_t \begin{pmatrix} 1 \\ x_t \end{pmatrix}.$$

The standard errors may be constructed via

$$\hat{V}_{OLS} = T \left( \sum_{t=1}^T \begin{pmatrix} 1 & x_t \\ x_t & x_t^2 \end{pmatrix} \right)^{-1} \cdot \sum_{t=1}^T \begin{pmatrix} 1 & x_t \\ x_t & x_t^2 \end{pmatrix} \hat{e}_t^2 \cdot \left( \sum_{t=1}^T \begin{pmatrix} 1 & x_t \\ x_t & x_t^2 \end{pmatrix} \right)^{-1},$$

where  $\hat{e}_t = y_t - \hat{\alpha}_{OLS} - \hat{\beta}_{OLS} x_t$ .

## 10.13 Does the link matter?

Let the  $x$  variable assume two different values  $x^0$  and  $x^1$ ,  $u^a = \alpha + \beta x^a$  and  $n_{ab} = \#\{x_i = x^a, y_i = b\}$ , for  $a, b = 0, 1$  (i.e.,  $n_{a,b}$  is the number of observations for which  $x_i = x^a, y_i = b$ ). The log-likelihood function is

$$\begin{aligned} l(x_1, \dots, x_n, y_1, \dots, y_n; \alpha, \beta) &= \log \left[ \prod_{i=1}^n F(\alpha + \beta x_i)^{y_i} (1 - F(\alpha + \beta x_i))^{1-y_i} \right] = \\ &= n_{01} \log F(u^0) + n_{00} \log(1 - F(u^0)) + n_{11} \log F(u^1) + n_{10} \log(1 - F(u^1)). \end{aligned} \quad (10.1)$$

The FOC for the problem of maximization of  $l(\dots; \alpha, \beta)$  w.r.t.  $\alpha$  and  $\beta$  are:

$$\begin{aligned} \left[ n_{01} \frac{F'(\hat{u}^0)}{F(\hat{u}^0)} - n_{00} \frac{F'(\hat{u}^0)}{1 - F(\hat{u}^0)} \right] + \left[ n_{11} \frac{F'(\hat{u}^1)}{F(\hat{u}^1)} - n_{10} \frac{F'(\hat{u}^1)}{1 - F(\hat{u}^1)} \right] &= 0, \\ x^0 \left[ n_{01} \frac{F'(\hat{u}^0)}{F(\hat{u}^0)} - n_{00} \frac{F'(\hat{u}^0)}{1 - F(\hat{u}^0)} \right] + x^1 \left[ n_{11} \frac{F'(\hat{u}^1)}{F(\hat{u}^1)} - n_{10} \frac{F'(\hat{u}^1)}{1 - F(\hat{u}^1)} \right] &= 0 \end{aligned}$$

As  $x^0 \neq x^1$ , one obtains for  $a = 0, 1$

$$\frac{n_{a1}}{F(\hat{u}^a)} - \frac{n_{a0}}{1 - F(\hat{u}^a)} = 0 \Leftrightarrow F(\hat{u}^a) = \frac{n_{a1}}{n_{a1} + n_{a0}} \Leftrightarrow \hat{u}^a \equiv \hat{\alpha} + \hat{\beta} x^a = F^{-1} \left( \frac{n_{a1}}{n_{a1} + n_{a0}} \right) \quad (10.2)$$

under the assumption that  $F'(\hat{u}^a) \neq 0$ . Comparing (10.1) and (10.2) one sees that  $l(\dots, \hat{\alpha}, \hat{\beta})$  does not depend on the form of the link function  $F(\cdot)$ . The estimates  $\hat{\alpha}$  and  $\hat{\beta}$  can be found from (10.2):

$$\hat{\alpha} = \frac{x^1 F^{-1} \left( \frac{n_{01}}{n_{01} + n_{00}} \right) - x^0 F^{-1} \left( \frac{n_{11}}{n_{11} + n_{10}} \right)}{x^1 - x^0}, \quad \hat{\beta} = \frac{F^{-1} \left( \frac{n_{11}}{n_{11} + n_{10}} \right) - F^{-1} \left( \frac{n_{01}}{n_{01} + n_{00}} \right)}{x^1 - x^0}.$$

## 10.14 Maximum likelihood and binary variables

Since the parameters in the conditional and marginal densities do not overlap, we can separate the problem. The conditional likelihood function is

$$\mathcal{L}(y_1, \dots, y_n, z_1, \dots, z_n, \gamma) = \prod_{i=1}^n \left( \frac{e^{\gamma z_i}}{1 + e^{\gamma z_i}} \right)^{y_i} \left( 1 - \frac{e^{\gamma z_i}}{1 + e^{\gamma z_i}} \right)^{1-y_i},$$

and the conditional loglikelihood –

$$\ell_n(y_1, \dots, y_n, z_1, \dots, z_n, \gamma) = \sum_{i=1}^n \{y_i \gamma z_i - \ln(1 + e^{\gamma z_i})\}.$$

The first order condition

$$\frac{\partial \ell_n}{\partial \gamma} = \sum_{i=1}^n \left\{ y_i z_i - \frac{z_i e^{\gamma z_i}}{1 + e^{\gamma z_i}} \right\} = 0$$

gives the solution  $\hat{\gamma} = \log \frac{n_{11}}{n_{10}}$ , where  $n_{11} = \#\{z_i = 1, y_i = 1\}$ ,  $n_{10} = \#\{z_i = 1, y_i = 0\}$ .

The marginal likelihood function is

$$\mathcal{L}(z_1, \dots, z_n, \alpha) = \prod_{i=1}^n \alpha^{z_i} (1 - \alpha)^{1-z_i},$$

and the marginal loglikelihood –

$$\ell_n(z_1, \dots, z_n, \alpha) = \sum_{i=1}^n \{z_i \ln \alpha + (1 - z_i) \ln(1 - \alpha)\}.$$

The first order condition

$$\frac{\partial \ell_n}{\partial \alpha} = \frac{\sum_{i=1}^n z_i}{\alpha} - \frac{\sum_{i=1}^n (1 - z_i)}{1 - \alpha} = 0$$

gives the solution  $\hat{\alpha} = \frac{1}{n} \sum_{i=1}^n z_i$ . From the asymptotic theory for ML,

$$\sqrt{n} \left( \begin{pmatrix} \hat{\alpha} \\ \hat{\gamma} \end{pmatrix} - \begin{pmatrix} \alpha \\ \gamma \end{pmatrix} \right) \xrightarrow{d} \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \alpha(1 - \alpha) & 0 \\ 0 & \frac{(1 + e^\gamma)^2}{\alpha e^\gamma} \end{pmatrix} \right).$$

## 10.15 Maximum likelihood and binary dependent variable

1. The conditional ML estimator is

$$\begin{aligned} \hat{\gamma}_{ML} &= \arg \max_c \sum_{i=1}^n \left\{ y_i \log \frac{e^{cx_i}}{1 + e^{cx_i}} + (1 - y_i) \log \frac{1}{1 + e^{cx_i}} \right\} \\ &= \arg \max_c \sum_{i=1}^n \{cy_i x_i - \log(1 + e^{cx_i})\}. \end{aligned}$$

The score is

$$s(y, x, \gamma) = \frac{\partial}{\partial \gamma} (\gamma y x - \log(1 + e^{\gamma x})) = \left( y - \frac{e^{\gamma x}}{1 + e^{\gamma x}} \right) x,$$

and the information matrix is

$$\mathcal{I} = -\mathbb{E} \left[ \frac{\partial s(y, x, \gamma)}{\partial \gamma} \right] = \mathbb{E} \left[ \frac{e^{\gamma x}}{(1 + e^{\gamma x})^2} x^2 \right],$$

so the asymptotic distribution of  $\hat{\gamma}_{ML}$  is  $N(0, \mathcal{I}^{-1})$ .

2. The regression is  $\mathbb{E}[y|x] = 1 \cdot \mathbb{P}\{y = 1|x\} + 0 \cdot \mathbb{P}\{y = 0|x\} = \frac{e^{\gamma x}}{1 + e^{\gamma x}}$ . The NLLS estimator is

$$\hat{\gamma}_{NLLS} = \arg \min_c \sum_{i=1}^n \left( y_i - \frac{e^{cx_i}}{1 + e^{cx_i}} \right)^2.$$

The asymptotic distribution of  $\hat{\gamma}_{NLLS}$  is  $\mathcal{N}(0, Q_{gg}^{-1} Q_{gge^2} Q_{gg}^{-1})$ . Now, since  $\mathbb{E}[e^2|x] = \mathbb{V}[y|x] = \frac{e^{\gamma x}}{(1 + e^{\gamma x})^2}$ , we have

$$Q_{gg} = \mathbb{E} \left[ \frac{e^{2\gamma x}}{(1 + e^{\gamma x})^4} x^2 \right], \quad Q_{gge^2} = \mathbb{E} \left[ \frac{e^{2\gamma x}}{(1 + e^{\gamma x})^4} x^2 \mathbb{E}[e^2|x] \right] = \mathbb{E} \left[ \frac{e^{3\gamma x}}{(1 + e^{\gamma x})^6} x^2 \right].$$

3. We know that  $\mathbb{V}[y|x] = \frac{e^{\gamma x}}{(1 + e^{\gamma x})^2}$ , which is a function of  $x$ . The WNLLS estimator of  $\gamma$  is

$$\hat{\gamma}_{WNLLS} = \arg \min_c \sum_{i=1}^n \frac{(1 + e^{\gamma x_i})^2}{e^{\gamma x_i}} \left( y_i - \frac{e^{cx_i}}{1 + e^{cx_i}} \right)^2.$$

Note that there should be the *true*  $\gamma$  in the weighting function (or its consistent estimate in a feasible version), but *not* the parameter of choice  $c$ ! The asymptotic distribution is  $\mathcal{N}\left(0, Q_{gg/\sigma^2}^{-1}\right)$ , where

$$Q_{gg/\sigma^2} = \mathbb{E} \left[ \frac{1}{\mathbb{V}[y|x]} \frac{e^{2\gamma x}}{(1 + e^{\gamma x})^4} x^2 \right] = \left( \frac{e^{\gamma x}}{(1 + e^{\gamma x})^2} x^2 \right).$$

4. For the ML problem, the moment condition is “zero expected score”

$$\mathbb{E} \left[ \left( y - \frac{e^{\gamma x}}{1 + e^{\gamma x}} \right) x \right] = 0.$$

For the NLLS problem, the moment condition is the FOC (or “no correlation between the error and the quasiregressor”)

$$\mathbb{E} \left[ \left( y - \frac{e^{\gamma x}}{1 + e^{\gamma x}} \right) \frac{e^{\gamma x}}{(1 + e^{\gamma x})^2} x \right] = 0.$$

For the WNLLS problem, the moment condition is similar:

$$\mathbb{E} \left[ \left( y - \frac{e^{\gamma x}}{1 + e^{\gamma x}} \right) x \right] = 0,$$

which is magically the same as for the ML problem. No wonder that the two estimators are asymptotically equivalent (see part 5).

5. Of course, from the general theory we have  $V_{MLE} \leq V_{WNLLS} \leq V_{NLLS}$ . We see a strict inequality  $V_{WNLLS} < V_{NLLS}$ , except maybe for special cases of the distribution of  $x$ , and this is not surprising. Surprising may seem the fact that  $V_{MLE} = V_{WNLLS}$ . It may be surprising because usually the MLE uses distributional assumptions, and the NLLSE does not, so usually we have  $V_{MLE} < V_{WNLLS}$ . In this problem, however, the distributional information is used by all estimators, that is, it is *not* an additional assumption made exclusively for ML estimation.

## 10.16 Poisson regression

1. The conditional logdensity is

$$\begin{aligned} \log f(y|x, \alpha, \beta) &= -\lambda(x, \alpha, \beta) + y \log \lambda(x, \alpha, \beta) - \log y! \\ &= \text{const} - \exp(\alpha + \beta x) + y(\alpha + \beta x), \end{aligned}$$

where const does not depend on parameters. Therefore, the unrestricted ML estimator is  $(\hat{\alpha}, \hat{\beta})$  solving

$$(\hat{\alpha}, \hat{\beta}) = \arg \max_{(a,b)} \sum_{i=1}^n \{y_i (a + bx_i) - \exp(a + bx_i)\},$$

or equivalently solving the system

$$\begin{aligned} \sum_{i=1}^n y_i - \sum_{i=1}^n \exp(\hat{\alpha} + \hat{\beta} x_i) &= 0, \\ \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \exp(\hat{\alpha} + \hat{\beta} x_i) &= 0. \end{aligned}$$

The loglikelihood function value is

$$\ell_n = \sum_{i=1}^n \left\{ \text{const} - \exp(\hat{\alpha} + \hat{\beta}x_i) + y_i(\hat{\alpha} + \hat{\beta}x_i) \right\} = n\text{const} + (\hat{\alpha} - 1)n\bar{y} + \hat{\beta}n\bar{x}\bar{y}.$$

The restricted MLE estimator is  $(\hat{\alpha}^R, 0)$  where

$$\hat{\alpha}^R = \arg \max_a \sum_{i=1}^n \{y_i a - \exp(a)\} = \log \bar{y},$$

with the loglikelihood function value

$$\ell_n^R = \sum_{i=1}^n \left\{ \text{const} - \exp(\hat{\alpha}^R) + y_i \hat{\alpha}^R \right\} = n\text{const} + n\bar{y}(\log \bar{y} - 1).$$

On this basis,

$$\mathcal{LR} = 2(\ell_n - \ell_n^R) = 2\left((\hat{\alpha} - \log \bar{y})n\bar{y} + \hat{\beta}n\bar{x}\bar{y}\right).$$

Now, the conditional score is

$$s(y, x, \alpha, \beta) = \begin{pmatrix} \partial \log f(y|x, \alpha, \beta) / \partial \alpha \\ \partial \log f(y|x, \alpha, \beta) / \partial \beta \end{pmatrix} = (y - \exp(\alpha + \beta x)) \begin{pmatrix} 1 \\ x \end{pmatrix},$$

and the true Information matrix is

$$\mathcal{I} = \mathbb{E} \left[ \exp(\alpha + \beta x) \begin{pmatrix} 1 \\ x \end{pmatrix} \begin{pmatrix} 1 \\ x \end{pmatrix}' \right],$$

consistently estimated by (most easily, albeit not quite in line with prescriptions of the theory)

$$\hat{\mathcal{I}} = \bar{y} \begin{pmatrix} 1 & \bar{x} \\ \bar{x} & \bar{x}^2 \end{pmatrix},$$

then

$$\hat{\mathcal{I}}^{-1} = \frac{1}{\bar{y}(\bar{x}^2 - \bar{x}^2)} \begin{pmatrix} \bar{x}^2 & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix},$$

Therefore,

$$\mathcal{W} = n\hat{\beta}^2 \bar{y} (\bar{x}^2 - \bar{x}^2).$$

If we used unconstrained estimated in  $\hat{\mathcal{I}}$  (in line with prescriptions of the theory), this would be more complex. Finally,

$$\begin{aligned} \mathcal{LM} &= \frac{1}{n} \left( \sum_{i=1}^n \begin{pmatrix} 0 \\ x_i(y_i - \bar{y}) \end{pmatrix} \right)' \frac{1}{\bar{y}(\bar{x}^2 - \bar{x}^2)} \begin{pmatrix} \bar{x}^2 & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix} \left( \sum_{i=1}^n \begin{pmatrix} 0 \\ x_i(y_i - \bar{y}) \end{pmatrix} \right) \\ &= n \frac{(\bar{x}\bar{y} - \bar{x}\bar{y})^2}{\bar{y}(\bar{x}^2 - \bar{x}^2)}. \end{aligned}$$

2. The logdensity is

$$\log f(y|x, \alpha, \beta) = \text{const} - \alpha - \exp(\beta x) + y \log(\alpha + \exp(\beta x)).$$

Therefore, the estimator used by the researcher is an extremum one with

$$h(y, x, \alpha, \beta) = -\alpha - \exp(\beta x) + y \log(\alpha + \exp(\beta x)).$$

The (misspecified) conditional score is

$$h_{\theta}(y, x, a, b) = \left( y (a + \exp(\beta x))^{-1} - 1 \right) \begin{pmatrix} 1 \\ \exp(\beta x) x \end{pmatrix}.$$

This extremum estimator will provide consistent estimation of  $(\alpha^*, \beta)$  only if

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} = \mathbb{E}[h_{\theta}(y, x, \alpha^*, \beta)].$$

Here  $\alpha^*$  may be any because we care only about consistency of estimate of  $\beta$ . But, using the law of iterated expectations,

$$\begin{aligned} \mathbb{E}[h_{\theta}(y, x, \alpha^*, \beta)] &= \mathbb{E} \left[ \left( y (\alpha^* + \exp(\beta x))^{-1} - 1 \right) \begin{pmatrix} 1 \\ \exp(\beta x) x \end{pmatrix} \right] \\ &= \mathbb{E} \left[ \left( \frac{(\exp(\alpha) - 1) \exp(\beta x) - \alpha^*}{\alpha^* + \exp(\beta x)} \right) \begin{pmatrix} 1 \\ \exp(\beta x) x \end{pmatrix} \right], \end{aligned}$$

which is not likely to be zero. This illustrates that misspecification of the conditional mean leads to wrong inference.

## 10.17 Bootstrapping ML tests

1. In the bootstrap world, the constraint is  $g(q) = g(\hat{\theta}_{ML})$ , so

$$\mathcal{LR}^* = 2 \left( \max_{q \in \Theta} \ell_n^*(q) - \max_{q \in \Theta, g(q) = g(\hat{\theta}_{ML})} \ell_n^*(q) \right),$$

where  $\ell_n^*$  is the loglikelihood calculated from the bootstrap sample.

2. In the bootstrap world, the constraint is  $g(q) = g(\hat{\theta}_{ML})$ , so

$$\mathcal{LM}^* = n \left( \frac{1}{n} \sum_{i=1}^n s(z_i^*, \hat{\theta}_{ML}^{*R}) \right)' \left( \hat{\mathcal{I}}^* \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n s(z_i^*, \hat{\theta}_{ML}^{*R}) \right),$$

where  $\hat{\theta}_{ML}^{*R}$  is the restricted (subject to  $g(q) = g(\hat{\theta}_{ML})$ ) ML bootstrap estimate and  $\hat{\mathcal{I}}^*$  is the bootstrap estimate of the information matrix, both calculated from the bootstrap sample. No additional recentering is needed because the zero expected score rule is exactly satisfied at the sample.





# 11. INSTRUMENTAL VARIABLES

## 11.1 Invalid 2SLS

1. Since  $\mathbb{E}[u] = 0$ , we have  $\mathbb{E}[y] = \alpha\mathbb{E}[z^2]$ , so  $\alpha$  is identified as long as  $z$  is not deterministic zero. The analog estimator is

$$\hat{\alpha} = \left( \frac{1}{n} \sum_i z_i^2 \right)^{-1} \frac{1}{n} \sum_i y_i.$$

Since  $\mathbb{E}[v] = 0$ , we have  $\mathbb{E}[z] = \pi\mathbb{E}[x]$ , so  $\pi$  is identified as long as  $x$  is not centered around zero. The analog estimator is

$$\hat{\pi} = \left( \frac{1}{n} \sum_i x_i \right)^{-1} \frac{1}{n} \sum_i z_i.$$

Since  $\Sigma$  does not depend on  $x$ , we have  $\Sigma = \mathbb{V} \left[ \begin{pmatrix} u \\ v \end{pmatrix} \right]$ , so  $\Sigma$  is identified. The analog estimator is

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} \hat{u}_i \\ \hat{v}_i \end{pmatrix} \begin{pmatrix} \hat{u}_i \\ \hat{v}_i \end{pmatrix}',$$

where  $\hat{u}_i = y_i - \hat{\alpha}z_i^2$  and  $\hat{v}_i = z_i - \hat{\pi}x_i$ .

2. The estimator satisfies

$$\tilde{\alpha} = \left( \frac{1}{n} \sum_i \hat{z}_i^4 \right)^{-1} \frac{1}{n} \sum_i \hat{z}_i^2 y_i = \left( \hat{\pi}^4 \frac{1}{n} \sum_i x_i^4 \right)^{-1} \hat{\pi}^2 \frac{1}{n} \sum_i x_i^2 y_i.$$

We know that  $\frac{1}{n} \sum_i x_i^4 \xrightarrow{p} \mathbb{E}[x^4]$ ,  $\frac{1}{n} \sum_i x_i^2 y_i = \alpha\pi^2 \frac{1}{n} \sum_i x_i^4 + 2\alpha\pi \frac{1}{n} \sum_i x_i^3 v_i + \alpha \frac{1}{n} \sum_i x_i^2 v_i^2 + \frac{1}{n} \sum_i x_i^2 u_i \xrightarrow{p} \alpha\pi^2 \mathbb{E}[x^4] + \alpha \mathbb{E}[x^2 v^2]$ , and  $\hat{\pi} \xrightarrow{p} \pi$ . Therefore,

$$\tilde{\alpha} \xrightarrow{p} \alpha + \frac{\alpha}{\pi^2} \frac{\mathbb{E}[x^2 v^2]}{\mathbb{E}[x^4]} \neq \alpha.$$

3. Evidently, we should fit the estimate of the square of  $z_i$ , instead of the square of the estimate. To do this, note that the second equation and properties of the model imply

$$\mathbb{E}[z^2|x] = \mathbb{E}[(\pi x + v)^2|x] = \pi^2 x^2 + 2\mathbb{E}[\pi x v|x] + \mathbb{E}[v^2|x] = \pi^2 x^2 + \sigma_v^2.$$

That is, we have a linear mean regression of  $z^2$  on  $x^2$  and a constant. Therefore, in the first stage we should regress  $z^2$  on  $x^2$  and a constant and construct  $\hat{z}_i^2 = \hat{\pi}^2 x_i^2 + \hat{\sigma}_v^2$ , and in the second stage, we should regress  $y$  on  $\hat{z}^2$ . Consistency of this estimator follows from the theory of 2SLS, when we treat  $z^2$  as a right hand side variable, not  $z$ .

## 11.2 Consumption function

The data are generated by

$$C_t = \frac{\alpha}{1-\lambda} + \frac{\lambda}{1-\lambda}A_t + \frac{1}{1-\lambda}e_t, \quad (11.1)$$

$$Y_t = \frac{\alpha}{1-\lambda} + \frac{1}{1-\lambda}A_t + \frac{1}{1-\lambda}e_t, \quad (11.2)$$

where  $A_t = I_t + G_t$  is exogenous and thus uncorrelated with  $e_t$ . Denote  $\sigma_e^2 = \mathbb{V}[e_t]$  and  $\sigma_A^2 = \mathbb{V}[A_t]$ .

1. The probability limit of the OLS estimator of  $\lambda$  is

$$p \lim \hat{\lambda} = \frac{\mathbb{C}[Y_t, C_t]}{\mathbb{V}[Y_t]} = \lambda + \frac{\mathbb{C}[Y_t, e_t]}{\mathbb{V}[Y_t]} = \lambda + \frac{\frac{1}{1-\lambda}\sigma_e^2}{\left(\frac{1}{1-\lambda}\right)^2\sigma_A^2 + \left(\frac{1}{1-\lambda}\right)^2\sigma_e^2} = \lambda + (1-\lambda)\frac{\sigma_e^2}{\sigma_A^2 + \sigma_e^2}.$$

The amount of inconsistency is  $(1-\lambda)\sigma_e^2/(\sigma_A^2 + \sigma_e^2)$ . Since the MPC lies between zero and one, the OLS estimator of  $\lambda$  is biased upward.

2. Econometrician B is correct in one sense, but incorrect in another. Both instrumental vectors will give rise to estimators that have identical *asymptotic properties*. This can be seen by noting that in population the projections of the right hand side variable  $Y_t$  on both instruments  $\Gamma z$ , where  $\Gamma = \mathbb{E}[xz'](\mathbb{E}[zz'])^{-1}$ , are identical. Indeed, because in (11.2) the  $I_t$  and  $G_t$  enter through their sum only, projecting on  $(1, I_t, G_t)'$  and on  $(1, A_t)'$  gives identical fitted values

$$\frac{\alpha}{1-\lambda} + \frac{1}{1-\lambda}A_t.$$

Consequently, the matrix  $Q_{xz}Q_{zz}^{-1}Q'_{xz}$  that figures into the asymptotic variance will be the same since it equals  $\mathbb{E}[\Gamma z(\Gamma z)']$  which is the same across the two instrumental vectors. However, this does not mean that the *numerical values* of the two estimates of  $(\alpha, \lambda)'$  will be the same. Indeed, the in-sample predicted values (that are used as regressors or instruments at the second stage of the 2SLS “procedure”) are  $\hat{x}_i = \hat{\Gamma}z_i = \mathcal{X}'\mathcal{Z}(\mathcal{Z}'\mathcal{Z})^{-1}z_i$ , and these values need not be the same for the “long” and “short” instrumental vectors.<sup>1</sup>

3. Econometrician C estimates the linear projection of  $Y_t$  on 1 and  $C_t$ , so the coefficient at  $C_t$  estimated by  $\hat{\theta}_C$  is

$$p \lim \hat{\theta}_C = \frac{\mathbb{C}[Y_t, C_t]}{\mathbb{V}[C_t]} = \frac{\frac{\lambda}{1-\lambda}\frac{1}{1-\lambda}\sigma_A^2 + \left(\frac{1}{1-\lambda}\right)^2\sigma_e^2}{\left(\frac{\lambda}{1-\lambda}\right)^2\sigma_A^2 + \left(\frac{1}{1-\lambda}\right)^2\sigma_e^2} = \frac{\lambda\sigma_A^2 + \sigma_e^2}{\lambda^2\sigma_A^2 + \sigma_e^2}.$$

Econometrician D estimates the linear projection of  $Y_t$  on 1,  $C_t$ ,  $I_t$ , and  $G_t$ , so the coefficient at  $C_t$  estimated by  $\hat{\phi}_C$  is 1 because of the perfect fit in the equation  $Y_t = 0 \cdot 1 + 1 \cdot C_t + 1 \cdot I_t + 1 \cdot G_t$ . Moreover, because of the perfect fit, the numerical value of  $(\hat{\phi}_0, \hat{\phi}_C, \hat{\phi}_I, \hat{\phi}_G)'$  will be exactly  $(0, 1, 1, 1)'$ .

<sup>1</sup>There exist a special case, however, when the numerical values will be equal (which is not the case in the problem at hand) – when the fit at the first stage is perfect.

## 11.3 Optimal combination of instruments

1. The necessary properties are validity and relevance:  $\mathbb{E}[ze] = \mathbb{E}[\zeta e] = 0$  and  $\mathbb{E}[zx] \neq 0$ ,  $\mathbb{E}[\zeta x] \neq 0$ . The asymptotic distributions of  $\hat{\beta}_z$  and  $\hat{\beta}_\zeta$  are

$$\sqrt{n} \begin{pmatrix} \hat{\beta}_z \\ \hat{\beta}_\zeta \end{pmatrix} - \begin{pmatrix} \beta \\ \beta \end{pmatrix} \xrightarrow{d} \mathcal{N} \left( 0, \begin{pmatrix} \mathbb{E}[zx]^{-2} \mathbb{E}[z^2 e^2] & \mathbb{E}[xz]^{-1} \mathbb{E}[x\zeta]^{-1} \mathbb{E}[z\zeta e^2] \\ \mathbb{E}[xz]^{-1} \mathbb{E}[x\zeta]^{-1} \mathbb{E}[z\zeta e^2] & \mathbb{E}[\zeta x]^{-2} \mathbb{E}[\zeta^2 e^2] \end{pmatrix} \right)$$

(we will need joint distribution in part 3).

2. The optimal instrument can be derived from the FOC for the GMM problem for the moment conditions

$$\mathbb{E} [m(y, x, z, \zeta, \beta)] = \mathbb{E} \left[ \begin{pmatrix} z \\ \zeta \end{pmatrix} (y - \beta x) \right] = 0.$$

Then

$$Q_{mm} = \mathbb{E} \left[ \begin{pmatrix} z \\ \zeta \end{pmatrix} \begin{pmatrix} z \\ \zeta \end{pmatrix}' e^2 \right], \quad Q_{\partial m} = -\mathbb{E} \left[ x \begin{pmatrix} z \\ \zeta \end{pmatrix} \right].$$

From the FOC for the (infeasible) efficient GMM in population, the optimal weighing of moment conditions and thus of instruments is then

$$\begin{aligned} Q'_{\partial m} Q_{mm}^{-1} &\propto \mathbb{E} \left[ x \begin{pmatrix} z \\ \zeta \end{pmatrix}' \right] \mathbb{E} \left[ \begin{pmatrix} z \\ \zeta \end{pmatrix} \begin{pmatrix} z \\ \zeta \end{pmatrix}' e^2 \right]^{-1} \\ &\propto \mathbb{E} \left[ x \begin{pmatrix} z \\ \zeta \end{pmatrix}' \right] \mathbb{E} \left[ \begin{pmatrix} \zeta \\ -z \end{pmatrix} \begin{pmatrix} \zeta \\ -z \end{pmatrix}' e^2 \right] \\ &\propto \begin{pmatrix} \mathbb{E}[xz] \mathbb{E}[\zeta^2 e^2] - \mathbb{E}[x\zeta] \mathbb{E}[z\zeta e^2] \\ \mathbb{E}[x\zeta] \mathbb{E}[z^2 e^2] - \mathbb{E}[xz] \mathbb{E}[z\zeta e^2] \end{pmatrix}'. \end{aligned}$$

That is, the optimal instrument is

$$(\mathbb{E}[xz] \mathbb{E}[\zeta^2 e^2] - \mathbb{E}[x\zeta] \mathbb{E}[z\zeta e^2]) z + (\mathbb{E}[x\zeta] \mathbb{E}[z^2 e^2] - \mathbb{E}[xz] \mathbb{E}[z\zeta e^2]) \zeta \equiv \gamma_z z + \gamma_\zeta \zeta.$$

This means that the optimally combined moment conditions imply

$$\mathbb{E} [(\gamma_z z + \gamma_\zeta \zeta) (y - \beta x)] = 0 \quad \Leftrightarrow$$

$$\begin{aligned} \beta &= \mathbb{E} [(\gamma_z z + \gamma_\zeta \zeta) x]^{-1} \mathbb{E} [(\gamma_z z + \gamma_\zeta \zeta) y] \\ &= \mathbb{E} [(\gamma_z z + \gamma_\zeta \zeta) x]^{-1} (\gamma_z \mathbb{E}[zx] \beta_z + \gamma_\zeta \mathbb{E}[\zeta x] \beta_\zeta), \end{aligned}$$

where  $\beta_z$  and  $\beta_\zeta$  are determined from the instruments separately. Thus the optimal IV estimator is the following linear combination of  $\hat{\beta}_z$  and  $\hat{\beta}_\zeta$ :

$$\begin{aligned} &\mathbb{E}[xz] \frac{\mathbb{E}[xz] \mathbb{E}[\zeta^2 e^2] - \mathbb{E}[x\zeta] \mathbb{E}[z\zeta e^2]}{\mathbb{E}[xz]^2 \mathbb{E}[\zeta^2 e^2] - 2\mathbb{E}[xz] \mathbb{E}[x\zeta] \mathbb{E}[z\zeta e^2] + \mathbb{E}[x\zeta]^2 \mathbb{E}[z^2 e^2]} \hat{\beta}_z \\ &+ \mathbb{E}[x\zeta] \frac{\mathbb{E}[x\zeta] \mathbb{E}[z^2 e^2] - \mathbb{E}[xz] \mathbb{E}[z\zeta e^2]}{\mathbb{E}[xz]^2 \mathbb{E}[\zeta^2 e^2] - 2\mathbb{E}[xz] \mathbb{E}[x\zeta] \mathbb{E}[z\zeta e^2] + \mathbb{E}[x\zeta]^2 \mathbb{E}[z^2 e^2]} \hat{\beta}_\zeta. \end{aligned}$$

3. Because of the joint convergence in part 1, the  $t$ -type test statistic can be constructed as

$$T = \frac{\sqrt{n} (\hat{\beta}_z - \hat{\beta}_\zeta)}{\sqrt{\hat{\mathbb{E}} [xz]^{-2} \hat{\mathbb{E}} [\zeta^2 e^2] - 2\hat{\mathbb{E}} [xz]^{-1} \hat{\mathbb{E}} [x\zeta]^{-1} \hat{\mathbb{E}} [z\zeta e^2] + \hat{\mathbb{E}} [x\zeta]^{-2} \hat{\mathbb{E}} [z^2 e^2]}}$$

where  $\hat{\mathbb{E}}$  denoted a sample analog of an expectation. The test rejects if  $|T|$  exceeds an appropriate quantile of the standard normal distribution. If the test rejects, one or both of  $z$  and  $\zeta$  may not be valid.

## 11.4 Trade and growth

1. The economic rationale for uncorrelatedness is that the variables  $P$  and  $S$  are exogenous and are unaffected by what's going on in the economy, and on the other hand, hardly can they affect the income in other ways than through the trade. To estimate (11.3), we can use just-identifying IV estimation, where the vector of right-hand-side variables is  $x = (1, T, W)'$  and the instrument vector is  $z = (1, P, S)'$ .
2. When data on within-country trade are not available, none of the coefficients in (11.3) is identifiable without further assumptions. In general, neither of the available variables can serve as instruments for  $T$  in (11.3) where the composite error term is  $\gamma W_i + \varepsilon_i$ .
3. We can exploit the assumption that  $P$  is uncorrelated with the error term in (11.5). Substitute (11.5) into (11.3) to get

$$\log Y = (\alpha + \gamma\eta) + \beta T + \gamma\lambda S + (\gamma\nu + \varepsilon).$$

Now we see that  $S$  and  $P$  are uncorrelated with the composite error term  $\gamma\nu + \varepsilon$  due to their exogeneity and due to their uncorrelatedness with  $\nu$  which follows from the additional assumption and because  $\nu$  is the best linear prediction error in (11.5). As for the coefficients of (11.3), only  $\beta$  will be consistently estimated, but not  $\alpha$  or  $\gamma$ .

4. In general, for this model the OLS is inconsistent, and the IV method is consistent. Thus, the discrepancy may be due to the different probability limits of the two estimators. Let  $\theta_{IV} \xrightarrow{p} \theta$  and  $\theta_{OLS} \xrightarrow{p} \theta + a$ ,  $a < 0$ . Then for large samples,  $\theta_{IV} \approx \theta$  and  $\theta_{OLS} \approx \theta + a$ . The difference is  $a$  which is  $(\mathbb{E}[xx'])^{-1} \mathbb{E}[xe]$ . Since  $(\mathbb{E}[xx'])^{-1}$  is positive definite,  $a < 0$  means that the regressors tend to be negatively correlated with the error term. In the present context this means that the trade variables are negatively correlated with other influences on income.

# 12. GENERALIZED METHOD OF MOMENTS

## 12.1 Nonlinear simultaneous equations

1. Since  $\mathbb{E}[u] = \mathbb{E}[v] = 0$ ,  $m(w, \theta) = \begin{pmatrix} y - \beta x \\ x - \gamma y^2 \end{pmatrix}$ , where  $w = \begin{pmatrix} x \\ y \end{pmatrix}$ ,  $\theta = \begin{pmatrix} \beta \\ \gamma \end{pmatrix}$ , can be used as a moment function. The true  $\beta$  and  $\gamma$  solve  $\mathbb{E}[m(w, \theta)] = 0$ , therefore  $\mathbb{E}[y] = \beta\mathbb{E}[x]$  and  $\mathbb{E}[x] = \gamma\mathbb{E}[y^2]$ , and they are identified as long as  $\mathbb{E}[x] \neq 0$  and  $\mathbb{E}[y^2] \neq 0$ . The analog of the population mean is the sample mean, so the analog estimators are

$$\hat{\beta} = \frac{\sum_i y_i}{\sum_i x_i}, \quad \hat{\gamma} = \frac{\sum_i x_i}{\sum_i y_i^2}.$$

2. If we add  $\mathbb{E}[uv] = 0$ , the moment function is

$$m(w, \theta) = \begin{pmatrix} y - \beta x \\ x - \gamma y^2 \\ (y - \beta x)(x - \gamma y^2) \end{pmatrix}$$

and GMM can be used. The feasible efficient GMM estimator is

$$\hat{\theta}_{GMM} = \arg \min_{q \in \Theta} \left( \frac{1}{n} \sum_{i=1}^n m(w_i, q) \right)' \hat{Q}_{mm}^{-1} \left( \frac{1}{n} \sum_{i=1}^n m(w_i, q) \right),$$

where  $\hat{Q}_{mm} = n^{-1} \sum_{i=1}^n m(w_i, \hat{\theta})m(w_i, \hat{\theta})'$  and  $\hat{\theta}$  is consistent estimator of  $\theta$  that can be taken from part 1. The asymptotic distribution of this estimator is

$$\sqrt{n}(\hat{\theta}_{GMM} - \theta) \xrightarrow{d} \mathcal{N}(0, V_{GMM}),$$

where  $V_{GMM} = (Q'_{\partial m} Q_{mm}^{-1} Q_{\partial m})^{-1}$ . The complete answer presumes expressing this matrix in terms of moments of observable variables.

## 12.2 Improved GMM

The first moment restriction gives GMM estimator  $\hat{\theta} = \bar{x}$  with asymptotic variance  $V_{CMM} = \mathbb{V}[x]$ . The GMM estimation of the full set of moment conditions gives estimator  $\hat{\theta}_{GMM}$  with asymptotic variance  $V_{GMM} = (Q'_{\partial m} Q_{mm}^{-1} Q_{\partial m})^{-1}$ , where  $Q_{\partial m} = \mathbb{E}[\partial m(x, y, \theta)/\partial q] = (-1, 0)'$  and

$$Q_{mm} = \mathbb{E} [m(x, y, \theta)m(x, y, \theta)'] = \begin{pmatrix} \mathbb{V}(x) & \mathbb{C}(x, y) \\ \mathbb{C}(x, y) & \mathbb{V}(y) \end{pmatrix}.$$

Hence,

$$V_{GMM} = \mathbb{V}[x] - \frac{(\mathbb{C}[x, y])^2}{\mathbb{V}[y]}$$

and thus efficient GMM estimation reduces the asymptotic variance when  $\mathbb{C}[x, y] \neq 0$ .

## 12.3 Minimum Distance estimation

1. Since the equation  $\theta_0 - s(\gamma_0) = 0$  can be uniquely solved for  $\gamma_0$ , we have

$$\gamma_0 = \arg \min_{\gamma \in \Gamma} (\theta_0 - s(\gamma))' W (\theta_0 - s(\gamma)).$$

For large  $n$ ,  $\hat{\theta}$  is concentrated around  $\theta_0$ , and  $\hat{W}$  is concentrated around  $W$ . Therefore,  $\hat{\gamma}_{MD}$  will be concentrated around  $\gamma_0$ . To derive the asymptotic distribution of  $\hat{\gamma}_{MD}$ , let us take the first order Taylor expansion of the last factor in the normalized sample FOC

$$0 = S(\hat{\gamma}_{MD})' \hat{W} \sqrt{n} (\hat{\theta} - s(\hat{\gamma}_{MD}))$$

around  $\gamma_0$ :

$$0 = S(\hat{\gamma}_{MD})' \hat{W} \sqrt{n} (\hat{\theta} - \theta_0) - S(\hat{\gamma}_{MD})' \hat{W} S(\bar{\gamma}) \sqrt{n} (\hat{\gamma}_{MD} - \gamma_0),$$

where  $\bar{\gamma}$  lies between  $\hat{\gamma}_{MD}$  and  $\gamma_0$  componentwise, hence  $\bar{\gamma} \xrightarrow{p} \gamma_0$ . Then

$$\begin{aligned} \sqrt{n} (\hat{\gamma}_{MD} - \gamma_0) &= \left( S(\hat{\gamma}_{MD})' \hat{W} S(\bar{\gamma}) \right)^{-1} S(\hat{\gamma}_{MD})' \hat{W} \sqrt{n} (\hat{\theta} - \theta_0) \\ &\stackrel{d}{\rightarrow} \left( S(\gamma_0)' W S(\gamma_0) \right)^{-1} S(\gamma_0)' W \mathcal{N}(0, V_{\hat{\theta}}) \\ &= \mathcal{N} \left( 0, \left( S(\gamma_0)' W S(\gamma_0) \right)^{-1} S(\gamma_0)' W V_{\hat{\theta}} W S(\gamma_0) \left( S(\gamma_0)' W S(\gamma_0) \right)^{-1} \right). \end{aligned}$$

2. By analogy with efficient GMM estimation, the optimal choice for the weight matrix  $W$  is  $V_{\hat{\theta}}^{-1}$ . Then

$$\sqrt{n} (\hat{\gamma}_{MD} - \gamma_0) \stackrel{d}{\rightarrow} \mathcal{N} \left( 0, \left( S(\gamma_0)' V_{\hat{\theta}}^{-1} S(\gamma_0) \right)^{-1} \right).$$

The obvious consistent estimator is  $\hat{V}_{\hat{\theta}}^{-1}$ . Note that it may be freely renormalized by a constant and this will not affect the result numerically.

3. Under  $H_0$ , the sample objective function is close to zero for large  $n$ , while under the alternative, it is far from zero. Let us take the first order Taylor expansion of the “root” of the optimal (i.e., when  $W = V_{\hat{\theta}}^{-1}$ ) sample objective function normalized by  $n$  around  $\gamma_0$ :

$$n \left( \hat{\theta} - s(\hat{\gamma}_{MD}) \right)' \hat{W} \left( \hat{\theta} - s(\hat{\gamma}_{MD}) \right) = \xi' \xi, \quad \xi \equiv \sqrt{n} \hat{W}^{1/2} \left( \hat{\theta} - s(\hat{\gamma}_{MD}) \right),$$

$$\begin{aligned} \xi &= \sqrt{n} \hat{W}^{1/2} (\hat{\theta} - \theta_0) - \sqrt{n} \hat{W}^{1/2} S(\bar{\gamma}) (\hat{\gamma}_{MD} - \gamma_0) \\ &\stackrel{A}{=} \left( I_{\ell} - V_{\hat{\theta}}^{-1/2} S(\gamma_0) \left( S(\gamma_0)' V_{\hat{\theta}}^{-1} S(\gamma_0) \right)^{-1} S(\gamma_0)' V_{\hat{\theta}}^{-1/2} \right) V_{\hat{\theta}}^{-1/2} \sqrt{n} (\hat{\theta} - \theta_0) \\ &\stackrel{A}{=} \left( I_{\ell} - V_{\hat{\theta}}^{-1/2} S(\gamma_0) \left( S(\gamma_0)' V_{\hat{\theta}}^{-1} S(\gamma_0) \right)^{-1} S(\gamma_0)' V_{\hat{\theta}}^{-1/2} \right) \mathcal{N}(0, I_{\ell}). \end{aligned}$$

Thus under  $H_0$

$$n \left( \hat{\theta} - s(\hat{\gamma}_{MD}) \right)' \hat{W} \left( \hat{\theta} - s(\hat{\gamma}_{MD}) \right) \stackrel{d}{\rightarrow} \chi_{\ell-k}^2.$$

4. The parameter of interest  $\rho$  is implicitly defined by the system

$$\begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} = \begin{pmatrix} 2\rho \\ -\rho^2 \end{pmatrix} \equiv s(\rho).$$

The matrix of derivatives is

$$S(\rho) \equiv \frac{\partial s(\rho)}{\partial \rho} = 2 \begin{pmatrix} 1 \\ -\rho \end{pmatrix}.$$

The OLS estimator of  $(\theta_1, \theta_2)'$  is consistent and asymptotically normal with asymptotic variance matrix

$$V_{\hat{\theta}} = \sigma^2 \begin{bmatrix} \mathbb{E}[y_t^2] & \mathbb{E}[y_t y_{t-1}] \\ \mathbb{E}[y_t y_{t-1}] & \mathbb{E}[y_t^2] \end{bmatrix}^{-1} = \frac{1 - \rho^4}{1 + \rho^2} \begin{bmatrix} 1 + \rho^2 & -2\rho \\ -2\rho & 1 + \rho^2 \end{bmatrix},$$

because

$$\mathbb{E}[y_t^2] = \sigma^2 \frac{1 + \rho^2}{(1 - \rho^2)^3}, \quad \frac{\mathbb{E}[y_t y_{t-1}]}{\mathbb{E}[y_t^2]} = \frac{2\rho}{1 + \rho^2}.$$

An optimal MD estimator of  $\rho$  is

$$\hat{\rho}_{MD} = \arg \min_{\rho: |\rho| < 1} \left( \begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{pmatrix} - \begin{pmatrix} 2\rho \\ -\rho^2 \end{pmatrix} \right)' \cdot \sum_{t=2}^n \begin{bmatrix} y_t^2 & y_t y_{t-1} \\ y_t y_{t-1} & y_t^2 \end{bmatrix} \cdot \left( \begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{pmatrix} - \begin{pmatrix} 2\rho \\ -\rho^2 \end{pmatrix} \right)$$

and is consistent and asymptotically normal with asymptotic variance

$$V_{\hat{\rho}_{MD}} = \left( 2 \begin{pmatrix} 1 \\ -\rho \end{pmatrix}' \frac{1 + \rho^2}{(1 - \rho^2)^3} \frac{1}{1 + \rho^2} \begin{bmatrix} 1 & 2\rho \\ 2\rho & 1 \end{bmatrix} 2 \begin{pmatrix} 1 \\ -\rho \end{pmatrix} \right)^{-1} = \frac{1 - \rho^2}{4}.$$

To verify that both autoregressive roots are indeed equal, we can test the hypothesis of correct specification. Let  $\hat{\sigma}^2$  be the estimated residual variance. The test statistic is

$$\frac{1}{\hat{\sigma}^2} \left( \begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{pmatrix} - \begin{pmatrix} 2\hat{\rho}_{MD} \\ -\hat{\rho}_{MD}^2 \end{pmatrix} \right)' \cdot \sum_{t=2}^n \begin{bmatrix} y_t^2 & y_t y_{t-1} \\ y_t y_{t-1} & y_t^2 \end{bmatrix} \cdot \left( \begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{pmatrix} - \begin{pmatrix} 2\hat{\rho}_{MD} \\ -\hat{\rho}_{MD}^2 \end{pmatrix} \right)$$

and is asymptotically distributed as  $\chi_1^2$ .

## 12.4 Formation of moment conditions

1. We know that  $\mathbb{E}[w] = \mu$  and  $\mathbb{E}[(w - \mu)^4] = 3 \left( \mathbb{E}[(w - \mu)^2] \right)^2$ . It is trivial to take care of the former. To take care of the latter, introduce a constant  $\sigma^2 = \mathbb{E}[(w - \mu)^2]$ , then we have  $\mathbb{E}[(w - \mu)^4] = 3(\sigma^2)^2$ . Together, the system of moment conditions is

$$\mathbb{E} \left[ \begin{pmatrix} w - \mu \\ (w - \mu)^2 - \sigma^2 \\ (w - \mu)^4 - 3(\sigma^2)^2 \end{pmatrix} \right] = 0.$$

2. To have overidentification, we need to arrive at least at 4 moment conditions (which must be informative about parameters!). Let, for example,  $z_t = (1, y_{t-1}, y_{t-1}^2, y_{t-2}, y_{t-2}^2)'$ . The moment conditions are

$$\mathbb{E} \left[ z_t \otimes \begin{pmatrix} y_t - \rho y_{t-1} \\ (y_t - \rho y_{t-1})^2 - \omega - \gamma (y_{t-1} - \rho y_{t-2})^2 \end{pmatrix} \right] = 0.$$

A much better way is to think and to use different instruments for the mean and variance equations. For example, let  $z_{1t} = (y_{t-1}, y_{t-2})'$ ,  $z_{2t} = (1, y_{t-1}^2, y_{t-2}^2)'$  so that the moment conditions are

$$\mathbb{E} \left[ \begin{pmatrix} z_{1t} (y_t - \rho y_{t-1}) \\ z_{2t} \left( (y_t - \rho y_{t-1})^2 - \omega - \gamma (y_{t-1} - \rho y_{t-2})^2 \right) \end{pmatrix} \right] = 0.$$

## 12.5 What CMM estimates

The population analog of

$$\sum_{i=1}^n g(z, q) = 0$$

is  $\mathbb{E}[g(z, q)] = 0$ . Thus, if the latter equation has a unique solution  $\theta$ , it is a probability limit of  $\hat{\theta}$ . The asymptotic distribution of  $\hat{\theta}$  is that of the CMM estimator of  $\theta$  based on the moment condition  $\mathbb{E}[g(z, \theta)] = 0$ :

$$\sqrt{n} (\hat{\theta} - \theta) \xrightarrow{d} \mathcal{N} \left( 0, \left( \mathbb{E} [\partial g(z, \theta) / \partial \theta'] \right)^{-1} \mathbb{E} [g(z, \theta) g(z, \theta)'] \left( \mathbb{E} [\partial g(z, \theta) / \partial \theta'] \right)^{-1} \right).$$

## 12.6 Trinity for GMM

The Wald test is the same up to a change in the variance matrix:

$$\mathcal{W} = nh(\hat{\theta}_{GMM})' \left[ H(\hat{\theta}_{GMM}) (\hat{Q}'_{\partial m} \hat{Q}_{mm}^{-1} \hat{Q}_{\partial m})^{-1} H(\hat{\theta}_{GMM})' \right]^{-1} h(\hat{\theta}_{GMM}) \xrightarrow{d} \chi_q^2,$$

where  $\hat{\theta}_{GMM}$  is the unrestricted GMM estimator,  $\hat{Q}_{\partial m}$  and  $\hat{Q}_{mm}$  are consistent estimators of  $Q_{\partial m}$  and  $Q_{mm}$ , relatively, and  $H(\theta) = \partial h(\theta) / \partial \theta'$ .

The Distance Difference test is similar to  $\mathcal{LR}$ , but without factor 2, as  $\partial^2 \hat{Q}_n / \partial \theta \partial \theta' \xrightarrow{p} 2Q'_{\partial m} Q_{mm}^{-1} Q_{\partial m}$ :

$$DD = n \left[ Q_n(\hat{\theta}_{GMM}^R) - Q_n(\hat{\theta}_{GMM}) \right] \xrightarrow{d} \chi_q^2.$$

The LM test is a little bit harder, since the analog of the average score is

$$\lambda(\theta) = 2 \left( \frac{1}{n} \sum_{i=1}^n \frac{\partial m(z_i, \theta)}{\partial \theta'} \right)' \hat{\Sigma}^{-1} \left( \frac{1}{n} \sum_{i=1}^n m(z_i, \theta) \right).$$

It is straightforward to find that

$$\mathcal{LM} = \frac{n}{4} \lambda(\hat{\theta}_{GMM}^R)' (\hat{Q}'_{\partial m} \hat{Q}_{mm}^{-1} \hat{Q}_{\partial m})^{-1} \lambda(\hat{\theta}_{GMM}^R) \xrightarrow{d} \chi_q^2.$$

In the middle one may use either restricted or unrestricted estimators of  $Q_{\partial m}$  and  $Q_{mm}$ .

A more detailed derivation can be found, for example, in section 7.4 of “Econometrics” by Fumio Hayashi (2000, Princeton University Press).



## 12.7 All about J

1. For the system of moment conditions  $\mathbb{E}[m(z, \theta)] = 0$ , the J-test statistic is

$$J = n\bar{m}(z, \hat{\theta})' \hat{Q}_{mm}^{-1} \bar{m}(z, \hat{\theta}).$$

Observe that

$$\sqrt{n}\bar{m}(z, \hat{\theta}) = \sqrt{n} \left( \bar{m}(z, \theta^*) + \frac{\partial \bar{m}(z, \tilde{\theta})}{\partial q'} (\hat{\theta} - \theta^*) \right),$$

where  $\theta^* = p \lim \hat{\theta}$ . By the CLT,  $\sqrt{n}(\bar{m}(z, \theta^*) - \mathbb{E}[m(z, \theta^*)])$  converges to the normal distribution. However, this means that

$$\sqrt{n}\bar{m}(z, \theta^*) = \sqrt{n}(\bar{m}(z, \theta^*) - \mathbb{E}[m(z, \theta^*)]) + \sqrt{n}\mathbb{E}[m(z, \theta^*)]$$

converges to infinity because of the second term which does not equal to zero (the system of moment conditions is misspecified). Hence,  $J$  is a quadratic form of something diverging to infinity, and thus diverges to infinity too.

2. Let the moment function be just  $x$  so that  $\ell = 1$  and  $k = 0$ . Then the J-test statistic corresponding to weight “matrix”  $\hat{w}$  equals  $J = n\hat{w}\bar{x}^2$ . Of course, when  $\hat{w} \xrightarrow{p} \mathbb{V}[x]^{-1}$  corresponds to efficient weighting, we have  $J \xrightarrow{d} \chi_1^2$ , but when  $\hat{w} \xrightarrow{p} w \neq \mathbb{V}[x]^{-1}$ , we have

$$J \xrightarrow{d} w\mathbb{V}[x] \chi_1^2 \neq \chi_1^2.$$

3. The argument would be valid if the model for the conditional mean was known to be correctly specified. Then one could blame instruments for a high value of the  $\mathcal{J}$ -statistic. But in our time series regression of the type  $\mathbb{E}_t[y_{t+1}] = g(x_t)$ , if this regression was correctly specified, then the variables from time  $t$  information set *must* be valid instruments! The failure of the model may be associated with incorrect functional form of  $g(\cdot)$ , or with specification of conditional information. Lastly, asymptotic theory may give a poor approximation to exact distribution of the  $\mathcal{J}$ -statistic.

## 12.8 Interest rates and future inflation

1. The conventional econometric model that tests the hypothesis of conditional unbiasedness of interest rates as predictors of inflation, is

$$\pi_t^k = \alpha_k + \beta_k i_t^k + \eta_t^k, \quad \mathbb{E}_t[\eta_t^k] = 0.$$

Under the null,  $\alpha_k = 0$ ,  $\beta_k = 1$ . Setting  $k = m$  in one case,  $k = n$  in the other case, and subtracting one equation from another, we can get

$$\pi_t^m - \pi_t^n = \alpha_m - \alpha_n + \beta_m i_t^m - \beta_n i_t^n + \eta_t^m - \eta_t^n, \quad \mathbb{E}_t[\eta_t^n - \eta_t^m] = 0.$$

Under the null  $\alpha_m = \alpha_n = 0$ ,  $\beta_m = \beta_n = 1$ , this specification coincides with Mishkin's under the null  $\alpha_{m,n} = 0$ ,  $\beta_{m,n} = 1$ . The restriction  $\beta_{m,n} = 0$  implies that the term structure provides no information about future shifts in inflation. The prediction error  $\eta_t^{m,n}$  is serially correlated of the order that is the farthest prediction horizon, i.e.,  $\max(m, n)$ .

2. Selection of instruments: there is a variety of choices, for instance,

$$\left(1, i_t^m - i_t^n, i_{t-1}^m - i_{t-1}^n, i_{t-2}^m - i_{t-2}^n, \pi_{t-\max(m,n)}^m - \pi_{t-\max(m,n)}^n\right)',$$

or

$$\left(1, i_t^m, i_t^n, i_{t-1}^m, i_{t-1}^n, \pi_{t-\max(m,n)}^m, \pi_{t-\max(m,n)}^n\right)',$$

etc. Construction of the optimal weighting matrix demands a HAC procedure, and so does estimation of asymptotic variance. The rest is more or less standard.

3. Most interesting are the results of the test  $\beta_{m,n} = 0$  which tell us that there is no information in the term structure about future path of inflation. Testing  $\beta_{m,n} = 1$  then seems excessive. This hypothesis would correspond to the conditional bias containing only a systematic component (i.e. a constant unpredictable by the term structure). It also looks like there is no systematic component in inflation (the hypothesis  $\alpha_{m,n} = 0$  can be accepted).

## 12.9 Spot and forward exchange rates

1. This is not the only way to proceed, but it is straightforward. The OLS estimator uses the instrument  $z_t^{OLS} = (1 \ x_t)'$ , where  $x_t = f_t - s_t$ . The additional moment condition adds  $f_{t-1} - s_{t-1}$  to the list of instruments:  $z_t = (1 \ x_t \ x_{t-1})'$ . Let us look at the optimal instrument. If it is proportional to  $z_t^{OLS}$ , then the instrument  $x_{t-1}$ , and hence the additional moment condition, is redundant. The optimal instrument takes the form  $\zeta_t = Q'_{\partial m} Q_{mm}^{-1} z_t$ . But

$$Q_{\partial m} = - \begin{pmatrix} 1 & \mathbb{E}[x_t] \\ \mathbb{E}[x_t] & \mathbb{E}[x_t^2] \\ \mathbb{E}[x_{t-1}] & \mathbb{E}[x_t x_{t-1}] \end{pmatrix}, \quad Q_{mm} = \sigma^2 \begin{pmatrix} 1 & \mathbb{E}[x_t] & \mathbb{E}[x_{t-1}] \\ \mathbb{E}[x_t] & \mathbb{E}[x_t^2] & \mathbb{E}[x_t x_{t-1}] \\ \mathbb{E}[x_{t-1}] & \mathbb{E}[x_t x_{t-1}] & \mathbb{E}[x_{t-1}^2] \end{pmatrix}.$$

It is easy to see that

$$Q'_{\partial m} Q_{mm}^{-1} = \sigma^{-2} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix},$$

which can be verified by postmultiplying this equation by  $Q_{mm}$ . Hence,  $\zeta_t = \sigma^{-2} z_t^{OLS}$ . But the most elegant way to solve this problem goes as follows. Under conditional homoskedasticity, the GMM estimator is asymptotically equivalent to the 2SLS estimator, if both use the same vector of instruments. But if the instrumental vector includes the regressors ( $z_t$  does include  $z_t^{OLS}$ ), the 2SLS estimator is identical to the OLS estimator. In total, GMM is asymptotically equivalent to OLS and thus the additional moment condition is redundant.

2. We can expect asymptotic equivalence of the OLS and efficient GMM estimators when the additional moment function is uncorrelated with the main moment function. Indeed, let us compare the  $2 \times 2$  northwestern block of  $V_{GMM} = (Q'_{\partial m} Q_{mm}^{-1} Q_{\partial m})^{-1}$  with asymptotic variance of the OLS estimator

$$V_{OLS} = \sigma^2 \begin{pmatrix} 1 & \mathbb{E}[x_t] \\ \mathbb{E}[x_t] & \mathbb{E}[x_t^2] \end{pmatrix}^{-1}.$$

Denote  $\Delta f_{t+1} = f_{t+1} - f_t$ . For the full set of moment conditions,

$$Q_{\partial m} = - \begin{pmatrix} 1 & \mathbb{E}[x_t] \\ \mathbb{E}[x_t] & \mathbb{E}[x_t^2] \\ 0 & 0 \end{pmatrix}, \quad Q_{mm} = \begin{pmatrix} \sigma^2 & \sigma^2 \mathbb{E}[x_t] & \mathbb{E}[x_t e_{t+1} \Delta f_{t+1}] \\ \sigma^2 \mathbb{E}[x_t] & \sigma^2 \mathbb{E}[x_t^2] & \mathbb{E}[x_t^2 e_{t+1} \Delta f_{t+1}] \\ \mathbb{E}[x_t e_{t+1} \Delta f_{t+1}] & \mathbb{E}[x_t^2 e_{t+1} \Delta f_{t+1}] & \mathbb{E}[x_t^2 (\Delta f_{t+1})^2] \end{pmatrix}.$$

It is easy to see that when  $\mathbb{E}[x_t e_{t+1} \Delta f_{t+1}] = \mathbb{E}[x_t^2 e_{t+1} \Delta f_{t+1}] = 0$ ,  $Q_{mm}$  is block-diagonal and the  $2 \times 2$  northwest block of  $V_{GMM}$  is the same as  $V_{OLS}$ . A sufficient condition for these two equalities is  $\mathbb{E}[e_{t+1} \Delta f_{t+1} | I_t] = 0$ , i. e. that conditionally on the past, unexpected movements in spot rates are uncorrelated with unexpected movements in forward rates. This is hardly satisfied in practice.

## 12.10 Returns from financial market

Let us denote  $y_t = r_{t+1} - r_t$ ,  $\beta = (\beta_0, \beta_1, \beta_2, \beta_3)'$  and  $x_t = (1, r_t, r_t^2, r_t^{-1})'$ . Then the model may be rewritten as

$$y_t = x_t' \beta + \varepsilon_{t+1}.$$

1. The OLS moment condition is  $\mathbb{E}[(y_t - x_t' \beta) x_t] = 0$ , with the OLS asymptotic distribution

$$\mathcal{N}\left(0, \sigma^2 \mathbb{E}[x_t x_t']^{-1} \mathbb{E}[x_t x_t' r_t^{2\gamma}] \mathbb{E}[x_t x_t']^{-1}\right).$$

The GLS moment condition is  $\mathbb{E}[(y_t - x_t' \beta) x_t r_t^{-2\gamma}] = 0$ , with the GLS asymptotic distribution

$$\mathcal{N}\left(0, \sigma^2 \mathbb{E}[x_t x_t' r_t^{-2\gamma}]^{-1}\right).$$

2. The parameter vector is  $\theta = (\beta', \sigma^2, \gamma)'$ . The conditional logdensity is

$$\log f(y_t | I_t, \beta, \sigma^2, \gamma) = -\frac{1}{2} \log \sigma^2 - \frac{\gamma}{2} \log r_t^2 - \frac{(y_t - x_t' \beta)^2}{2\sigma^2 r_t^{2\gamma}}$$

(note that it is illegitimate to take a log of  $r_t$ ) The conditional score is

$$s(y_t | I_t, \beta, \sigma^2, \gamma) = \begin{pmatrix} \frac{(y_t - x_t' \beta) x_t}{\sigma^2 r_t^{2\gamma}} \\ \frac{1}{2\sigma^2} \left( \frac{(y_t - x_t' \beta)^2}{\sigma^2 r_t^{2\gamma}} - 1 \right) \\ \frac{1}{2} \left( \frac{(y_t - x_t' \beta)^2}{\sigma^2 r_t^{2\gamma}} - 1 \right) \log r_t^2 \end{pmatrix}.$$

Hence, the ML moment condition is

$$\mathbb{E} \begin{pmatrix} \frac{(y_t - x_t' \beta) x_t}{r_t^{2\gamma}} \\ \frac{(y_t - x_t' \beta)^2}{\sigma^2 r_t^{2\gamma}} - 1 \\ \left( \frac{(y_t - x_t' \beta)^2}{\sigma^2 r_t^{2\gamma}} - 1 \right) \log r_t^2 \end{pmatrix} = 0.$$

Note that its first part is the GLS moment condition. The information matrix is

$$\mathcal{I} = \begin{pmatrix} \frac{1}{\sigma^2} \mathbb{E} \left[ \frac{x_t x_t'}{r_t^{2\gamma}} \right] & 0 & 0 \\ 0 & \frac{1}{2\sigma^4} & \frac{1}{2\sigma^2} \mathbb{E} [\log r_t^2] \\ 0 & \frac{1}{2\sigma^2} \mathbb{E} [\log r_t^2] & \frac{1}{2} \mathbb{E} [\log^2 r_t^2] \end{pmatrix}.$$

The asymptotic distribution for the ML estimates of the regression parameters is

$$\mathcal{N}\left(0, \sigma^2 \mathbb{E}[x_t x_t' r_t^{-2\gamma}]^{-1}\right),$$

and for the ML estimates of  $\sigma^2$  and  $\gamma$  it is independent of the former, and is

$$\mathcal{N}\left(0, \frac{2\sigma^2}{\mathbb{V}[\log r_t^2]} \begin{pmatrix} \sigma^2 \mathbb{E}[\log^2 r_t^2] & -\mathbb{E}[\log r_t^2] \\ -\mathbb{E}[\log r_t^2] & \sigma^{-2} \end{pmatrix}\right).$$

3. The method of moments will use the moment conditions

$$\mathbb{E} \begin{pmatrix} (y_t - x_t' \beta) x_t \\ \left( (y_t - x_t' \beta)^2 - \sigma^2 r_t^{2\gamma} \right) r_t^{2\gamma} \\ \left( (y_t - x_t' \beta)^2 - \sigma^2 r_t^{2\gamma} \right) \sigma^2 r_t^{2\gamma} \log r_t^2 \end{pmatrix} = 0,$$

where the former one is from Part 1, while the second and the third are the expected products of the error and quasi-regressors in the skedastic regression. Note that the resulting CMM estimator will have the usual OLS estimator for the  $\beta$ -part because of exact identification. Therefore, the  $\beta$ -part will have the same asymptotic distribution:

$$\mathcal{N}\left(0, \sigma^2 \mathbb{E}[x_t x_t']^{-1} \mathbb{E}[x_t x_t' r_t^{2\gamma}] \mathbb{E}[x_t x_t']^{-1}\right).$$

As for the other two parts, it is more messy. Indeed,

$$Q_{\partial m} = - \begin{pmatrix} \mathbb{E}[x_t x_t'] & 0 & 0 \\ 0 & \mathbb{E}[r_t^{4\gamma}] & \sigma^2 \mathbb{E}[r_t^{4\gamma} \log r_t^2] \\ 0 & \sigma^2 \mathbb{E}[r_t^{4\gamma} \log r_t^2] & \sigma^4 \mathbb{E}[r_t^{4\gamma} \log^2 r_t^2] \end{pmatrix},$$

$$Q_{mm} = \begin{pmatrix} \sigma^2 \mathbb{E}[x_t x_t' r_t^{2\gamma}] & 0 & 0 \\ 0 & 2\sigma^4 \mathbb{E}[r_t^{8\gamma}] & 2\sigma^6 \mathbb{E}[r_t^{8\gamma} \log r_t^2] \\ 0 & 2\sigma^6 \mathbb{E}[r_t^{8\gamma} \log r_t^2] & 2\sigma^8 \mathbb{E}[r_t^{8\gamma} \log^2 r_t^2] \end{pmatrix}.$$

The asymptotic variances follow from  $Q_{\partial m}^{-1} Q_{mm} Q_{\partial m}'^{-1}$ . The sandwich does not collapse, but one can see that the pair  $(\hat{\sigma}_{CMM}^2, \hat{\gamma}_{CMM})'$  is asymptotically independent of  $\hat{\beta}_{CMM}$ .

4. We have already established that the OLS estimator is identical to the  $\beta$ -part of the CMM estimator. From general principles, the GLS estimator is asymptotically more efficient, and is identical to the  $\beta$ -part of the ML estimator. As for the other two parameters, the ML estimator is asymptotically efficient, and hence is asymptotically more efficient than the appropriate part of the CMM estimator which, by construction, is the same as implied by NLLS applied to the skedastic function. To summarize, for  $\beta$  we have  $OLS = CMM \prec GLS = ML$ , while for  $\sigma^2$  and  $\gamma$  we have  $CMM \prec ML$ .

## 12.11 Instrumental variables in ARMA models

1. The instrument  $x_{t-j}$  is scalar, the parameter is scalar, so there is exact identification. The instrument is obviously valid. The asymptotic variance of the just identifying IV estimator of a scalar parameter under homoskedasticity is  $V_{x_{t-j}} = \sigma^2 Q_{xz}^{-2} Q_{zz}$ . Let us calculate all pieces:

$$Q_{zz} = \mathbb{E}[x_{t-j}^2] = \mathbb{V}[x_t] = \sigma^2 (1 - \rho^2)^{-1},$$

$$Q_{xz} = \mathbb{E}[x_{t-1} x_{t-j}] = \mathbb{C}[x_{t-1}, x_{t-j}] = \rho^{j-1} \mathbb{V}[x_t] = \sigma^2 \rho^{j-1} (1 - \rho^2)^{-1}.$$

Thus,  $V_{x_{t-j}} = \rho^{2-2j} (1 - \rho^2)$ . It is monotonically declining in  $j$ , so this suggests that the optimal instrument must be  $x_{t-1}$ . Although this is not a proof of the fact, the optimal instrument is indeed  $x_{t-1}$ . The result makes sense, since the last observation is most informative and embeds all information in all the other instruments.

2. It is possible to use as instruments lags of  $y_t$  starting from  $y_{t-2}$  back to the past. The regressor  $y_{t-1}$  will not do as it is correlated with the error term through  $e_{t-1}$ . Among  $y_{t-2}, y_{t-3}, \dots$  the first one deserves more attention, since, intuitively, it contains more information than older values of  $y_t$ .

## 12.12 Hausman may not work

1. Because the instruments include the regressors,  $\hat{\beta}_{2SLS}$  will be identical to  $\hat{\beta}_{OLS}$ , because the instrument vector contains  $x$ , the “fitted values” from the “first-stage projection” will be equal to  $x$ . Hence,  $\hat{\beta}_{OLS} - \hat{\beta}_{2SLS}$  will be zero irrespective of validity of  $z$ , and so will be the Hausman test statistic.
2. Under the null of conditional homoskedasticity both estimators are asymptotically equivalent and hence equally asymptotically efficient. Hence, the difference in asymptotic variances is identically zero. Note that here the estimators are not identical.

## 12.13 Testing moment conditions

Consider the unrestricted ( $\hat{\beta}_u$ ) and restricted ( $\hat{\beta}_r$ ) estimates of parameter  $\beta \in R^k$ . The first is the CMM estimate:

$$\sum_{i=1}^n x_i (y_i - x_i' \hat{\beta}_u) = 0 \quad \Rightarrow \quad \hat{\beta}_u = \left( \frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \frac{1}{n} \sum_{i=1}^n x_i y_i$$

The second is a feasible efficient GMM estimate:

$$\hat{\beta}_r = \arg \min_b \left( \frac{1}{n} \sum_{i=1}^n m_i(b) \right)' \hat{Q}_{mm}^{-1} \left( \frac{1}{n} \sum_{i=1}^n m_i(b) \right), \quad (12.1)$$

where  $m(b) = \begin{pmatrix} xu(b) \\ xu(b)^3 \end{pmatrix}$ ,  $u(b) = y - xb$ ,  $u \equiv u(\beta)$ , and  $\hat{Q}_{mm}^{-1}$  is a consistent estimator of

$$Q_{mm} = \mathbb{E} [m(\beta) m'(\beta)] = \mathbb{E} \left[ \begin{pmatrix} xx'u^2 & xx'u^4 \\ xx'u^4 & xx'u^6 \end{pmatrix} \right].$$

Denote also  $Q_{\partial m} = \mathbb{E} \left[ \frac{\partial m(\beta)}{\partial b'} \right] = -\mathbb{E} \left[ \begin{pmatrix} xx' \\ 3xx'u^2 \end{pmatrix} \right]$ . Writing out the FOC for (12.1) and expanding  $m(\hat{\beta}_r)$  around  $\beta$  gives after rearrangement

$$\sqrt{n}(\hat{\beta}_r - \beta) \stackrel{A}{=} - (Q'_{\partial m} Q_{mm}^{-1} Q_{\partial m})^{-1} Q'_{\partial m} Q_{mm}^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n m_i(\beta).$$

Here  $\stackrel{A}{\equiv}$  means that we substitute the probability limits for their sample analogues. The last equation holds under the null hypothesis  $H_0 : \mathbb{E} [xu^3] = 0$ .

Note that the unrestricted estimate can be rewritten as

$$\sqrt{n}(\hat{\beta}_u - \beta) \stackrel{A}{\equiv} (\mathbb{E} [xx'])^{-1} ( I_k \quad O_k ) \frac{1}{\sqrt{n}} \sum_{i=1}^n m_i(\beta).$$

Therefore,

$$\sqrt{n}(\hat{\beta}_u - \beta_r) \stackrel{A}{\equiv} \left[ (\mathbb{E} [xx'])^{-1} ( I_k \quad O_k ) + (Q'_{\partial m} Q_{mm}^{-1} Q_{\partial m})^{-1} Q'_{\partial m} Q_{mm}^{-1} \right] \frac{1}{\sqrt{n}} \sum_{i=1}^n m_i(\beta) \xrightarrow{d} \mathcal{N}(0, V),$$

where (after some algebra)

$$V = (\mathbb{E} [xx'])^{-1} \mathbb{E} [xx'u^2] (\mathbb{E} [xx'])^{-1} - (Q'_{\partial m} Q_{mm}^{-1} Q_{\partial m})^{-1}.$$

Note that  $V$  is  $k \times k$  matrix. It can be shown that this matrix is non-degenerate (and thus has a full rank  $k$ ). Let  $\hat{V}$  be a consistent estimate of  $V$ . By the Slutsky and Mann–Wald theorems,

$$\mathcal{W} \equiv n(\hat{\beta}_u - \hat{\beta}_r)' \hat{V}^{-1} (\hat{\beta}_u - \hat{\beta}_r) \xrightarrow{d} \chi_k^2.$$

The test may be implemented as follows. First find the (consistent) estimate  $\hat{\beta}_u$ . Then compute  $\hat{Q}_{mm} = n^{-1} \sum_{i=1}^n m_i(\hat{\beta}_u) m_i(\hat{\beta}_u)'$ , use it to carry out feasible GMM and obtain  $\hat{\beta}_r$ . Use  $\hat{\beta}_u$  or  $\hat{\beta}_r$  to compute  $\hat{V}$ , the sample analog of  $V$ . Finally, compute the Wald statistic  $\mathcal{W}$ , compare it with 95% quantile of  $\chi_k^2$  distribution  $q_{0.95}$ , and reject the null hypothesis if  $\mathcal{W} > q_{0.95}$ , or accept it otherwise.

## 12.14 Bootstrapping OLS

Indeed, we are supposed to recenter, but only when there is overidentification. When the parameter is just identified, as in the case of the OLS estimator, the moment conditions hold exactly in the sample, so the “center” is zero anyway.

## 12.15 Bootstrapping DD

Let  $\hat{\theta}$  denote the GMM estimator. Then the bootstrap  $\mathcal{DD}$  test statistic is

$$\mathcal{DD}^* = n \left[ \min_{q: h(q) = h(\hat{\theta})} \mathcal{Q}_n^*(q) - \min_q \mathcal{Q}_n^*(q) \right],$$

where  $\mathcal{Q}_n^*(q)$  is the bootstrap GMM objective function

$$\mathcal{Q}_n^*(q) = \left( \frac{1}{n} \sum_{i=1}^n m(z_i^*, q) - \frac{1}{n} \sum_{i=1}^n m(z_i, \hat{\theta}) \right)' \hat{Q}_{mm}^{*-1} \left( \frac{1}{n} \sum_{i=1}^n m(z_i^*, q) - \frac{1}{n} \sum_{i=1}^n m(z_i, \hat{\theta}) \right),$$

where  $\hat{Q}_{mm}^*$  uses the formula for  $\hat{Q}_{mm}$  and the bootstrap sample. Note two instances of recentering.

# 13. PANEL DATA

## 13.1 Alternating individual effects

It is convenient to use three indices instead of two in indexing the data. Namely, let

$$t = 2(s - 1) + q, \text{ where } q \in \{1, 2\}, s \in \{1, \dots, T\}.$$

Then  $q = 1$  corresponds to odd periods, while  $q = 2$  corresponds to even periods. The dummy variables will have the form of the Kronecker product of three matrices, which is defined recursively as  $A \otimes B \otimes C = A \otimes (B \otimes C)$ .

Part 1. (a) In this case we rearrange the data column as follows:

$$y_{isq} = y_{it}, \quad y_{is} = \begin{pmatrix} y_{is1} \\ y_{is2} \end{pmatrix}, \quad \mathcal{Y}_i = \begin{pmatrix} y_{i1} \\ \dots \\ y_{iT} \end{pmatrix}, \quad \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_1 \\ \dots \\ \mathcal{Y}_n \end{pmatrix},$$

and  $\mu = (\mu_1^O \ \mu_1^E \ \dots \ \mu_n^O \ \mu_n^E)'$ . The regressors and errors are rearranged in the same manner as  $y$ 's. Then the regression can be rewritten as

$$\mathcal{Y} = D\mu + \mathcal{X}\beta + v, \tag{13.1}$$

where  $D = I_n \otimes i_T \otimes I_2$ , and  $i_T = (1 \ \dots \ 1)'$  ( $T \times 1$  vector). Clearly,

$$D'D = I_n \otimes i_T' i_T \otimes I_2 = T \cdot I_{2n},$$

$$D(D'D)^{-1}D' = \frac{1}{T}I_n \otimes i_T i_T' \otimes I_2 = \frac{1}{T}I_n \otimes J_T \otimes I_2,$$

where  $J_T = i_T i_T'$ . In other words,  $D(D'D)^{-1}D'$  is block-diagonal with  $n$  blocks of size  $2T \times 2T$  of the form:

$$\begin{pmatrix} \frac{1}{T} & 0 & \dots & \frac{1}{T} & 0 \\ 0 & \frac{1}{T} & \dots & 0 & \frac{1}{T} \\ \dots & \dots & \dots & \dots & \dots \\ \frac{1}{T} & 0 & \dots & \frac{1}{T} & 0 \\ 0 & \frac{1}{T} & \dots & 0 & \frac{1}{T} \end{pmatrix}.$$

The  $Q$ -matrix is then  $Q = I_{2nT} - \frac{1}{T}I_n \otimes J_T \otimes I_2$ . Note that  $Q$  is an orthogonal projection and  $QD = 0$ . Thus we have from (13.1)

$$Q\mathcal{Y} = Q\mathcal{X}\beta + Qv. \tag{13.2}$$

Note that  $\frac{1}{T}J_T$  is the operator of taking the mean over the  $s$ -index (i.e. over odd or even periods depending on the value of  $q$ ). Therefore, the transformed regression is:

$$y_{isq} - \bar{y}_{iq} = (x_{isq} - \bar{x}_{iq})'\beta + v^*, \tag{13.3}$$

where  $\bar{y}_{iq} = \sum_{s=1}^T y_{isq}$ .

(b) This time the data are rearranged in the following manner:

$$y_{qis} = y_{it}; \quad \mathcal{Y}_{qi} = \begin{pmatrix} y_{i1} \\ \dots \\ y_{iT} \end{pmatrix}; \quad \mathcal{Y}_q = \begin{pmatrix} \mathcal{Y}_{q1} \\ \dots \\ \mathcal{Y}_{qn} \end{pmatrix}; \quad \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_1 \\ \mathcal{Y}_2 \end{pmatrix};$$

$\mu = (\mu_1^O \dots \mu_n^O \mu_1^E \dots \mu_n^E)'$ . In matrix form the regression is again (13.1) with  $D = I_2 \otimes I_n \otimes i_T$ , and

$$D(D'D)^{-1}D' = \frac{1}{T}I_2 \otimes I_n \otimes i_T i_T' = \frac{1}{T}I_2 \otimes I_n \otimes J_T.$$

This matrix consists of  $2n$  blocks on the main diagonal, each of them being  $\frac{1}{T}J_T$ . The  $Q$ -matrix is  $Q = I_{2n \cdot T} - \frac{1}{T}I_{2n} \otimes J_T$ . The rest is as in part 1(b) with the transformed regression

$$y_{qis} - \bar{y}_{qi} = (x_{qis} - \bar{x}_{qi})'\beta + v^*, \quad (13.4)$$

with  $\bar{y}_{qi} = \sum_{s=1}^T y_{qis}$ , which is essentially the same as (13.3).

Part 2. Take the  $Q$ -matrix as in Part 1(b). The Within estimator is the OLS estimator in (13.4), i.e.  $\hat{\beta} = (\mathcal{X}'Q\mathcal{X})^{-1}\mathcal{X}'Q\mathcal{Y}$ , or

$$\hat{\beta} = \left( \sum_{q,i,s} (x_{qis} - \bar{x}_{qi})(x_{qis} - \bar{x}_{qi})' \right)^{-1} \sum_{q,i,s} (x_{qis} - \bar{x}_{qi})(y_{qis} - \bar{y}_{qi}).$$

Clearly,  $\mathbb{E}[\hat{\beta}] = \beta$ ,  $\hat{\beta} \xrightarrow{p} \beta$  and  $\hat{\beta}$  is asymptotically normal as  $n \rightarrow \infty$ ,  $T$  fixed. For normally distributed errors  $v_{qis}$  the standard F-test for hypothesis

$$H_0 : \mu_1^O = \mu_2^O = \dots = \mu_n^O \text{ and } \mu_1^E = \mu_2^E = \dots = \mu_n^E$$

is

$$F = \frac{(RSS^R - RSS^U)/(2n-2)^{H_0}}{RSS^U/(2nT-2n-k)} \overset{H_0}{\sim} F(2n-2, 2nT-2n-k)$$

(we have  $2n-2$  restrictions in the hypothesis), where  $RSS^U = \sum_{isq} (y_{qis} - \bar{y}_{qi} - (x_{qis} - \bar{x}_{qi})'\beta)^2$ , and  $RSS^R$  is the sum of squared residuals in the restricted regression.

Part 3. Here we start with

$$y_{qis} = x'_{qis}\beta + u_{qis}, \quad u_{qis} := \mu_{qi} + v_{qis}, \quad (13.5)$$

where  $\mu_{1i} = \mu_i^O$  and  $\mu_{2i} = \mu_i^E$ ;  $\mathbb{E}[\mu_{qi}] = 0$ . Let  $\sigma_1^2 = \sigma_O^2$ ,  $\sigma_2^2 = \sigma_E^2$ . We have

$$\mathbb{E}[u_{qis}u_{q'i's'}] = \mathbb{E}[(\mu_{qi} + v_{qis})(\mu_{q'i'} + v_{q'i's'})] = \sigma_q^2 \delta_{qq'} \delta_{ii'} 1_{ss'} + \sigma_v^2 \delta_{qq'} \delta_{ii'} \delta_{ss'},$$

where  $\delta_{aa'} = \{1 \text{ if } a = a', \text{ and } 0 \text{ if } a \neq a'\}$ ,  $1_{ss'} = 1$  for all  $s, s'$ . Consequently,

$$\begin{aligned} \Omega &= \mathbb{E}[uu'] = \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix} \otimes I_n \otimes J_T + \sigma_v^2 I_{2nT} = (T\sigma_1^2 + \sigma_v^2) \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \otimes I_n \otimes \frac{1}{T}J_T + \\ &+ (T\sigma_2^2 + \sigma_v^2) \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \otimes I_n \otimes \frac{1}{T}J_T + \sigma_v^2 I_2 \otimes I_n \otimes (I_T - \frac{1}{T}J_T). \end{aligned}$$

The last expression is the spectral decomposition of  $\Omega$  since all operators in it are idempotent symmetric matrices (orthogonal projections), which are orthogonal to each other and give identity in sum. Therefore,

$$\begin{aligned} \Omega^{-1/2} &= (T\sigma_1^2 + \sigma_v^2)^{-1/2} \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \otimes I_n \otimes \frac{1}{T}J_T + (T\sigma_2^2 + \sigma_v^2)^{-1/2} \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \otimes I_n \otimes \frac{1}{T}J_T + \\ &+ \sigma_v^{-1} I_2 \otimes I_n \otimes (I_T - \frac{1}{T}J_T). \end{aligned}$$

The GLS estimator of  $\beta$  is

$$\hat{\beta} = (\mathcal{X}'\Omega^{-1}\mathcal{X})^{-1}\mathcal{X}'\Omega^{-1}\mathcal{Y}.$$



To put it differently,  $\hat{\beta}$  is the OLS estimator in the transformed regression

$$\sigma_v \Omega^{-1/2} \mathcal{Y} = \sigma_v \Omega^{-1/2} \mathcal{X} \beta + u^*.$$

The latter may be rewritten as

$$y_{qis} - (1 - \sqrt{\theta_q}) \bar{y}_{qi} = (x_{qis} - (1 - \sqrt{\theta_q}) \bar{x}_{qi})' \beta + u^*,$$

where  $\theta_q = \sigma_v^2 / (\sigma_v^2 + T \sigma_q^2)$ .

To make  $\hat{\beta}$  feasible, we should consistently estimate parameter  $\theta_q$ . In the case  $\sigma_1^2 = \sigma_2^2$  we may apply the result obtained in class (we have  $2n$  different objects and  $T$  observations for each of them – see part 1(b)):

$$\hat{\theta} = \frac{2n - k}{2n(T - 1) - k + 1} \frac{\hat{u}' Q \hat{u}}{\hat{u}' P \hat{u}},$$

where  $\hat{u}$  are OLS-residuals for (13.4), and  $Q = I_{2n \cdot T} - \frac{1}{T} I_{2n} \otimes J_T$ ,  $P = I_{2n \cdot T} - Q$ . Suppose now that  $\sigma_1^2 \neq \sigma_2^2$ . Using equations

$$\mathbb{E}[u_{qis}] = \sigma_v^2 + \sigma_q^2; \quad \mathbb{E}[\bar{u}_{is}] = \frac{1}{T} \sigma_v^2 + \sigma_q^2,$$

and repeating what was done in class, we have

$$\hat{\theta}_q = \frac{n - k}{n(T - 1) - k + 1} \frac{\hat{u}' Q_q \hat{u}}{\hat{u}' P_q \hat{u}},$$

with  $Q_1 = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \otimes I_n \otimes (I_T - \frac{1}{T} J_T)$ ,  $Q_2 = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \otimes I_n \otimes (I_T - \frac{1}{T} J_T)$ ,  $P_1 = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \otimes I_n \otimes \frac{1}{T} J_T$ ,  
 $P_2 = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \otimes I_n \otimes \frac{1}{T} J_T$ .

## 13.2 Time invariant regressors

1. (a) Under fixed effects, the  $z_i$  variable is collinear with the dummy for  $\mu_i$ . Thus,  $\gamma$  is unidentifiable.. The Within transformation wipes out the term  $z_i \gamma$  together with individual effects  $\mu_i$ , so the transformed equation looks exactly like it looks if no term  $z_i \gamma$  is present in the model. Under usual assumptions about independence of  $v_{it}$  and  $\mathcal{X}$ , the Within estimator of  $\beta$  is efficient.  
 (b) Under random effects and mutual independence of  $\mu_i$  and  $v_{it}$ , as well as their independence of  $\mathcal{X}$  and  $\mathcal{Z}$ , the GLS estimator is efficient, and the feasible GLS estimator is asymptotically efficient as  $n \rightarrow \infty$ .
2. Recall that the first-step  $\hat{\beta}$  is consistent but  $\hat{\pi}_i$ 's are inconsistent as  $T$  stays fixed and  $n \rightarrow \infty$ . However, the estimator of  $\gamma$  so constructed is consistent under assumptions of random effects (see part 1(b)). Observe that  $\hat{\pi}_i = \bar{y}_i - \bar{x}_i' \hat{\beta}$ . If we regress  $\hat{\pi}_i$  on  $z_i$ , we get the OLS coefficient

$$\begin{aligned} \hat{\gamma} &= \frac{\sum_{i=1}^n z_i \hat{\pi}_i}{\sum_{i=1}^n z_i^2} = \frac{\sum_{i=1}^n z_i (\bar{y}_i - \bar{x}_i' \hat{\beta})}{\sum_{i=1}^n z_i^2} = \frac{\sum_{i=1}^n z_i (\bar{x}_i' \beta + z_i \gamma + \mu_i + \bar{v}_i - \bar{x}_i' \hat{\beta})}{\sum_{i=1}^n z_i^2} \\ &= \gamma + \frac{\frac{1}{n} \sum_{i=1}^n z_i \mu_i}{\frac{1}{n} \sum_{i=1}^n z_i^2} + \frac{\frac{1}{n} \sum_{i=1}^n z_i \bar{v}_i}{\frac{1}{n} \sum_{i=1}^n z_i^2} + \frac{\frac{1}{n} \sum_{i=1}^n z_i \bar{x}_i'}{\frac{1}{n} \sum_{i=1}^n z_i^2} (\beta - \hat{\beta}). \end{aligned}$$

Now, as  $n \rightarrow \infty$ ,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n z_i^2 &\xrightarrow{p} \mathbb{E}[z_i^2] \neq 0, & \frac{1}{n} \sum_{i=1}^n z_i \mu_i &\xrightarrow{p} \mathbb{E}[z_i \mu_i] = \mathbb{E}[z_i] \mathbb{E}[\mu_i] = 0, \\ \frac{1}{n} \sum_{i=1}^n z_i \bar{v}_i &\xrightarrow{p} \mathbb{E}[z_i \bar{v}_i] = \mathbb{E}[z_i] \mathbb{E}[\bar{v}_i] = 0, & \frac{1}{n} \sum_{i=1}^n z_i \bar{x}'_i &\xrightarrow{p} \mathbb{E}[z_i \bar{x}'_i], & \beta - \hat{\beta} &\xrightarrow{p} 0. \end{aligned}$$

In total,  $\hat{\gamma} \xrightarrow{p} \gamma$ . However, so constructed estimator of  $\gamma$  is asymptotically inefficient. A better estimator is the feasible GLS estimator of part 1(b).

### 13.3 Within and Between

Recall that  $\hat{\beta}_W - \beta = (\mathcal{X}'Q\mathcal{X})^{-1} \mathcal{X}'Q\mathcal{U}$  and  $\hat{\beta}_B - \beta = (\mathcal{X}'P\mathcal{X})^{-1} \mathcal{X}'P\mathcal{U}$ , where  $\mathcal{U}$  is a vector of two-component errors, and that  $\mathbb{V}[\mathcal{U}|\mathcal{X}] = \Omega = \alpha_Q Q + \alpha_P P$  for certain weights  $\alpha_Q$  and  $\alpha_P$  that depend on variance components. Then

$$\begin{aligned} \mathbb{C}[\hat{\beta}_W, \hat{\beta}_B|\mathcal{X}] &= (\mathcal{X}'Q\mathcal{X})^{-1} \mathcal{X}'Q\mathbb{V}[\mathcal{U}|\mathcal{X}]P\mathcal{X} (\mathcal{X}'P\mathcal{X})^{-1} \\ &= (\mathcal{X}'Q\mathcal{X})^{-1} \mathcal{X}'Q(\alpha_Q Q + \alpha_P P)P\mathcal{X} (\mathcal{X}'P\mathcal{X})^{-1} \\ &= \alpha_Q (\mathcal{X}'Q\mathcal{X})^{-1} \mathcal{X}'QQP\mathcal{X} (\mathcal{X}'P\mathcal{X})^{-1} + \alpha_P (\mathcal{X}'Q\mathcal{X})^{-1} \mathcal{X}'QPP\mathcal{X} (\mathcal{X}'P\mathcal{X})^{-1} \\ &= 0, \end{aligned}$$

because  $QP = PQ = 0$ . Using this result,

$$\begin{aligned} \mathbb{V}[\hat{\beta}_W - \hat{\beta}_B|\mathcal{X}] &= \mathbb{V}[\hat{\beta}_W|\mathcal{X}] + \mathbb{V}[\hat{\beta}_B|\mathcal{X}] \\ &= (\mathcal{X}'Q\mathcal{X})^{-1} \mathcal{X}'Q\Omega Q\mathcal{X} (\mathcal{X}'Q\mathcal{X})^{-1} + (\mathcal{X}'P\mathcal{X})^{-1} \mathcal{X}'P\Omega P\mathcal{X} (\mathcal{X}'P\mathcal{X})^{-1} \\ &= \alpha_Q (\mathcal{X}'Q\mathcal{X})^{-1} + \alpha_P (\mathcal{X}'P\mathcal{X})^{-1}. \end{aligned}$$

We know that under random effects both  $\hat{\beta}_W$  and  $\hat{\beta}_B$  are consistent and asymptotically normal, hence the test statistic

$$\mathcal{R} = (\hat{\beta}_W - \hat{\beta}_B)' \left( \hat{\alpha}_Q (\mathcal{X}'Q\mathcal{X})^{-1} + \hat{\alpha}_P (\mathcal{X}'P\mathcal{X})^{-1} \right)^{-1} (\hat{\beta}_W - \hat{\beta}_B)$$

is asymptotically chi-squared with  $k = \dim(\beta)$  degrees of freedom under random effects. Here, as usual,  $\hat{\alpha}_Q = RSS_W / (nT - n - k)$  and  $\hat{\alpha}_P = RSS_B / (n - k)$ . When random effects are inappropriate,  $\hat{\beta}_W$  is still consistent but  $\hat{\beta}_B$  is inconsistent, so  $\hat{\beta}_W - \hat{\beta}_B$  converges to a non-zero limit, while  $\mathcal{X}'Q\mathcal{X}$  and  $\mathcal{X}'P\mathcal{X}$  diverge to infinity, making  $\mathcal{R}$  diverge to infinity.

### 13.4 Panels and instruments

Stack the regressors into matrix  $X$ , the instruments – into matrix  $Z$ , the dependent variables – into vector  $\mathcal{Y}$ . The Between transformation is implied by the transformation matrix  $Q = I_n \otimes T^{-1}J_T$ , where  $J_T$  is a  $T \times T$  matrix of ones, the Within transformation is implied by the transformation

matrix  $Q = I_{nT} - P$ , the GLS transformation is implied by the transformation matrix  $\Omega^{-1/2} = \alpha_Q Q + \alpha_P P$  for certain weights  $\alpha_Q$  and  $\alpha_P$  that depend on variance components. Recall that  $P$  and  $Q$  are mutually orthogonal symmetric idempotent matrices.

1. When in the Within regression  $Q\mathcal{Y} = QX\beta + Q\mathcal{U}$  the original instruments  $\mathcal{Z}$  are used, the IV estimator is  $(\mathcal{Z}'Q\mathcal{X})^{-1} \mathcal{Z}'Q\mathcal{Y}$ . When the Within-transformed instruments  $Q\mathcal{Z}$  are used, the IV estimator is  $((Q\mathcal{Z})'Q\mathcal{X})^{-1} (Q\mathcal{Z})'Q\mathcal{Y} = (\mathcal{Z}'Q\mathcal{X})^{-1} \mathcal{Z}'Q\mathcal{Y}$ . These two are identical.
2. The same as in part 1, with  $P$  in place of  $Q$  everywhere.
3. When in the Between regression  $P\mathcal{Y} = PX\beta + P\mathcal{U}$  the Within-transformed instruments  $Q\mathcal{Z}$  are used,  $(Q\mathcal{Z})'PX = \mathcal{Z}'Q'PX = \mathcal{Z}'0X = 0$ , a zero matrix. Such IV estimator does not exist. The same result holds when one uses the Between-transformed instruments in the Within regression.
4. When in the GLS-transformed regression  $\Omega^{-1/2}\mathcal{Y} = \Omega^{-1/2}X\beta + \Omega^{-1/2}\mathcal{U}$  the original instruments  $\mathcal{Z}$  are used, the IV estimator is

$$\left(\mathcal{Z}'\Omega^{-1/2}\mathcal{X}\right)^{-1} \mathcal{Z}'\Omega^{-1/2}\mathcal{Y}.$$

When the GLS-transformed instruments  $\Omega^{-1/2}\mathcal{Z}$  are used, the IV estimator is

$$\left(\left(\Omega^{-1/2}\mathcal{Z}\right)' \Omega^{-1/2}\mathcal{X}\right)^{-1} \left(\Omega^{-1/2}\mathcal{Z}\right)' \Omega^{-1/2}\mathcal{Y} = (\mathcal{Z}'\Omega^{-1}\mathcal{X})^{-1} \mathcal{Z}'\Omega^{-1}\mathcal{Y}.$$

These two are different. The second one is more natural to do: when instruments coincide with regressors, the efficient GLS estimator results, while the former estimator is inefficient and weird.

5. When in the GLS-transformed regression the Within-transformed instruments are used, the IV estimator is

$$\begin{aligned} \left((Q\mathcal{Z})' \Omega^{-1/2}\mathcal{X}\right)^{-1} (Q\mathcal{Z})' \Omega^{-1/2}\mathcal{Y} &= (\mathcal{Z}'Q'(\alpha_Q Q + \alpha_P P)\mathcal{X})^{-1} \mathcal{Z}'Q'(\alpha_Q Q + \alpha_P P)\mathcal{Y} \\ &= (\mathcal{Z}'(\alpha_Q Q)\mathcal{X})^{-1} \mathcal{Z}'(\alpha_Q Q)\mathcal{Y} \\ &= (\mathcal{Z}'Q\mathcal{X})^{-1} \mathcal{Z}'Q\mathcal{Y}, \end{aligned}$$

the estimator from part 1. Similarly, when in the GLS-transformed regression one uses the Between-transformed instruments, the IV estimator is that of part 2.

## 13.5 Differencing transformations

1. OLS on FD-transformed equations is unbiased and consistent as  $n \rightarrow \infty$  since the differenced error has mean zero conditional on the matrix of differenced regressors under the standard FE assumptions. However, OLS is inefficient as the conditional variance matrix is not diagonal. The efficient estimator of structural parameters is the LSDV estimator, which is the OLS estimator on Within-transformed equations.
2. The proposal leads to a consistent, but not very efficient, GMM estimator. The resulting error term  $v_{i,t} - v_{i,2}$  is uncorrelated only with  $y_{i,1}$  among all  $y_{i,1}, \dots, y_{i,T}$  so that for all equations we can find much fewer instruments than in the FD approach, and the same is true for the regular regressors if they are predetermined, but not strictly exogenous. As a result, we lose efficiency but get nothing in return.

## 13.6 Nonlinear panel data model

The nonlinear one-way ECM with random effects is

$$(y_{it} + \alpha)^2 = \beta x_{it} + \mu_i + v_{it}, \quad \mu_i \sim IID(0, \sigma_\mu^2), \quad v_{it} \sim IID(0, \sigma_v^2),$$

where individual effects  $\mu_i$  and idiosyncratic shocks  $v_{it}$  are mutually independent and independent of  $x_{it}$ . The latter assumption is unnecessarily strong and may be relaxed. The estimator of Problem 9.3 obtained from the pooled sample is inefficient since it ignores nondiagonality of the variance matrix of the error vector. We have to construct an analog of the GLS estimator in a linear one-way ECM with random effects. To get preliminary consistent estimates of variance components, we run analogs of Within and Between regressions (we call the resulting estimators Within-CMM and Between-CMM):

$$\begin{aligned} (y_{it} + \alpha)^2 - \frac{1}{T} \sum_{t=1}^T (y_{it} + \alpha)^2 &= \beta \left( x_{it} - \frac{1}{T} \sum_{t=1}^T x_{it} \right) + v_{it} - \frac{1}{T} \sum_{t=1}^T v_{it} \\ \frac{1}{T} \sum_{t=1}^T (y_{it} + \alpha)^2 &= \beta \frac{1}{T} \sum_{t=1}^T x_{it} + \mu_i + \frac{1}{T} \sum_{t=1}^T v_{it} \end{aligned}$$

Numerically estimates can be obtained by concentration as described in Problem 9.3. The estimated variance components and the “GLS-CMM parameter” can be found from

$$\hat{\sigma}_v^2 = \frac{RSS_W}{Tn - n - 2}, \quad \hat{\sigma}_\mu^2 + \frac{1}{T} \hat{\sigma}_v^2 = \frac{RSS_B}{T(n - 2)} \quad \Rightarrow \quad \hat{\theta} = \frac{RSS_W}{RSS_B} \frac{n - 2}{Tn - n - 2}.$$

Note that  $RSS_W$  and  $RSS_B$  are sums of squared residuals in the Within-CMM and Between-CMM systems, *not* the values of CMM objective functions. Then we consider the FGLS-transformed system where the variance matrix of the error vector is (asymptotically) diagonalized:

$$(y_{it} + \alpha)^2 - (1 - \sqrt{\hat{\theta}}) \frac{1}{T} \sum_{t=1}^T (y_{it} + \alpha)^2 = \beta \left( x_{it} - (1 - \sqrt{\hat{\theta}}) \frac{1}{T} \sum_{t=1}^T x_{it} \right) + \text{error term}.$$

## 13.7 Durbin–Watson statistic and panel data

1. In both regressions, the residuals consistently estimate corresponding regression errors. Therefore, to find a probability limit of the Durbin–Watson statistic, it suffices to compute the variance and first-order autocovariance of the errors across the stacked equations:

$$\text{p lim}_{n \rightarrow \infty} \mathcal{DW} = 2 \left( 1 - \frac{\varrho_1}{\varrho_0} \right),$$

where

$$\varrho_0 = \text{p lim}_{n \rightarrow \infty} \frac{1}{nT} \sum_{t=1}^T \sum_{i=1}^n u_{it}^2, \quad \varrho_1 = \text{p lim}_{n \rightarrow \infty} \frac{1}{nT} \sum_{t=2}^T \sum_{i=1}^n u_{it} u_{i,t-1},$$

and  $u_{it}$ 's denote regression errors. Note that the errors are uncorrelated where the index  $i$  switches between individuals, hence summation from  $t = 2$  in  $\varrho_1$ . Consider the original regression

$$y_{it} = x'_{it} \beta + u_{it}, \quad i = 1, \dots, n, \quad t = 1, \dots, T.$$

where  $u_{it} = \mu_i + v_{it}$ . Then  $\varrho_0 = \sigma_v^2 + \sigma_\mu^2$  and

$$\varrho_1 = \frac{1}{T} \sum_{t=2}^T \text{p lim}_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n (\mu_i + v_{it}) (\mu_i + v_{i,t-1}) = \frac{T-1}{T} \sigma_\mu^2.$$

Thus

$$\text{p lim}_{n \rightarrow \infty} \mathcal{DW}_{OLS} = 2 \left( 1 - \frac{T-1}{T} \frac{\sigma_\mu^2}{\sigma_v^2 + \sigma_\mu^2} \right) = 2 \frac{T\sigma_v^2 + \sigma_\mu^2}{T(\sigma_v^2 + \sigma_\mu^2)}.$$

The GLS-transformation orthogonalizes the errors, therefore

$$\text{p lim}_{n \rightarrow \infty} \mathcal{DW}_{GLS} = 2.$$

2. Since all computed probability limits except that for  $\mathcal{DW}_{OLS}$  do not depend on the variance components, the only way to construct an asymptotic test of  $H_0 : \sigma_\mu^2 = 0$  vs.  $H_A : \sigma_\mu^2 > 0$  is by using  $\mathcal{DW}_{OLS}$ . Under  $H_0$ ,  $\sqrt{nT}(\mathcal{DW}_{OLS} - 2) \xrightarrow{d} \mathcal{N}(0, 4)$  as  $n \rightarrow \infty$ . Under  $H_A$ ,  $\text{p lim}_{n \rightarrow \infty} \mathcal{DW}_{OLS} < 2$ . Hence a one-sided asymptotic test for  $\sigma_\mu^2 = 0$  for a given level  $\alpha$  is:

$$\text{Reject if } \mathcal{DW}_{OLS} < 2 \left( 1 + \frac{z_\alpha}{\sqrt{nT}} \right),$$

where  $z_\alpha$  is the  $\alpha$ -quantile of the standard normal distribution.

## 13.8 Higher-order dynamic panel

The original system is

$$y_{it} = \delta_1 y_{i,t-1} + \delta_2 y_{i,t-2} + \beta x_{it} + \mu_i + v_{it}, \quad i = 1, \dots, n, \quad t = 3, \dots, T,$$

where  $\mu_i \sim IID(0, \sigma_\mu^2)$  and  $v_i \sim IID(0, \sigma_v^2)$  are independent of each other and of  $x_i$ 's. The FD-transformed system is

$$\begin{aligned} y_{it} - y_{i,t-1} &= \delta_1 (y_{i,t-1} - y_{i,t-2}) + \delta_2 (y_{i,t-2} - y_{i,t-3}) + \beta (x_{it} - x_{i,t-1}) + v_{it} - v_{i,t-1}, \\ i &= 1, \dots, n, \quad t = 4, \dots, T, \end{aligned}$$

or in the matrix form,

$$\Delta \mathcal{Y} = \delta_1 \Delta \mathcal{Y}_{-1} + \delta_2 \Delta \mathcal{Y}_{-2} + \beta \Delta \mathcal{X} + \Delta \mathcal{V},$$

where the vertical dimension is  $n(T-3)$ . The GMM estimator of the 3-dimensional parameter vector  $\theta = (\delta_1, \delta_2, \beta)'$  is

$$\begin{aligned} \hat{\theta}_{GMM} &= \left( (\Delta \mathcal{Y}_{-1} \ \Delta \mathcal{Y}_{-2} \ \Delta \mathcal{X})' \mathcal{W} (\mathcal{W}' (I_n \otimes G) \mathcal{W})^{-1} \mathcal{W}' (\Delta \mathcal{Y}_{-1} \ \Delta \mathcal{Y}_{-2} \ \Delta \mathcal{X}) \right)^{-1} \\ &\quad \times (\Delta \mathcal{Y}_{-1} \ \Delta \mathcal{Y}_{-2} \ \Delta \mathcal{X})' \mathcal{W} (\mathcal{W}' (I_n \otimes G) \mathcal{W})^{-1} \mathcal{W}' \Delta \mathcal{Y}, \end{aligned}$$

where the  $(T-3) \times ((T+1)(T-3))$  matrix of instruments for the  $i^{th}$  individual is

$$\mathcal{W}_i = \text{diag} \begin{bmatrix} y_{i,1} & y_{i,2} & x_{i,1} & x_{i,2} & x_{i,3} \\ y_{i,1} & y_{i,2} & y_{i,3} & x_{i,1} & x_{i,2} & x_{i,3} & x_{i,4} \\ \vdots & & & & & & \\ y_{i,1} & \dots & y_{i,T-2} & x_{i,1} & \dots & x_{i,T-1} \end{bmatrix}$$

The standard errors can be computed using

$$\widehat{AV}(\hat{\theta}_{GMM}) = \hat{\sigma}_v^2 \left( (\Delta\mathcal{Y}_{-1} \ \Delta\mathcal{Y}_{-2} \ \Delta\mathcal{X})' \mathcal{W} (\mathcal{W}' (I_n \otimes G) \mathcal{W})^{-1} \mathcal{W}' (\Delta\mathcal{Y}_{-1} \ \Delta\mathcal{Y}_{-2} \ \Delta\mathcal{X}) \right)^{-1},$$

where, for example,

$$\hat{\sigma}_v^2 = \frac{1}{2n(T-3)-3} (\widehat{\Delta\mathcal{V}})' \widehat{\Delta\mathcal{V}},$$

and  $\widehat{\Delta\mathcal{V}}$  are residuals from the FD-transformed system.

# 14. NONPARAMETRIC ESTIMATION

## 14.1 Nonparametric regression with discrete regressor

Fix  $a_{(j)}$ ,  $j = 1, \dots, k$ . Observe that

$$g(a_{(j)}) = \mathbb{E}[y_i | x_i = a_{(j)}] = \frac{\mathbb{E}(y_i \mathbb{I}[x_i = a_{(j)}])}{\mathbb{E}(\mathbb{I}[x_i = a_{(j)}])}$$

because of the following equalities:

$$\begin{aligned} \mathbb{E}[\mathbb{I}[x_i = a_{(j)}]] &= 1 \cdot \mathbb{P}\{x_i = a_{(j)}\} + 0 \cdot \mathbb{P}\{x_i \neq a_{(j)}\} = \mathbb{P}\{x_i = a_{(j)}\}, \\ \mathbb{E}[y_i \mathbb{I}[x_i = a_{(j)}]] &= \mathbb{E}[y_i \mathbb{I}[x_i = a_{(j)}] | x_i = a_{(j)}] \cdot \mathbb{P}\{x_i = a_{(j)}\} = \mathbb{E}[y_i | x_i = a_{(j)}] \cdot \mathbb{P}\{x_i = a_{(j)}\}. \end{aligned}$$

According to the analogy principle we can construct  $\hat{g}(a_{(j)})$  as

$$\hat{g}(a_{(j)}) = \frac{\sum_{i=1}^n y_i \mathbb{I}[x_i = a_{(j)}]}{\sum_{i=1}^n \mathbb{I}[x_i = a_{(j)}]}.$$

Now let us find its properties. First, according to the LLN,

$$\hat{g}(a_{(j)}) = \frac{\sum_{i=1}^n y_i \mathbb{I}[x_i = a_{(j)}]}{\sum_{i=1}^n \mathbb{I}[x_i = a_{(j)}]} \xrightarrow{p} \frac{\mathbb{E}[y_i \mathbb{I}[x_i = a_{(j)}]]}{\mathbb{E}[\mathbb{I}[x_i = a_{(j)}]]} = g(a_{(j)}).$$

Second,

$$\sqrt{n}(\hat{g}(a_{(j)}) - g(a_{(j)})) = \sqrt{n} \frac{\sum_{i=1}^n (y_i - \mathbb{E}[y_i | x_i = a_{(j)}]) \mathbb{I}[x_i = a_{(j)}]}{\sum_{i=1}^n \mathbb{I}[x_i = a_{(j)}]}.$$

According to the CLT,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (y_i - \mathbb{E}[y_i | x_i = a_{(j)}]) \mathbb{I}[x_i = a_{(j)}] \xrightarrow{d} \mathcal{N}(0, \omega),$$

where

$$\begin{aligned} \omega &= \mathbb{V}[(y_i - \mathbb{E}[y_i | x_i = a_{(j)}]) \mathbb{I}[x_i = a_{(j)}]] = \mathbb{E}[(y_i - \mathbb{E}[y_i | x_i = a_{(j)}])^2 | x_i = a_{(j)}] \mathbb{P}\{x_i = a_{(j)}\} \\ &= \mathbb{V}[y_i | x_i = a_{(j)}] \mathbb{P}\{x_i = a_{(j)}\}. \end{aligned}$$

Thus

$$\sqrt{n}(\hat{g}(a_{(j)}) - g(a_{(j)})) \xrightarrow{d} \mathcal{N}\left(0, \frac{\mathbb{V}[y_i | x_i = a_{(j)}]}{\mathbb{P}\{x_i = a_{(j)}\}}\right).$$

## 14.2 Nonparametric density estimation

(a) Use the hint that  $\mathbb{E}[\mathbb{I}[x_i \leq x]] = F(x)$  to prove the unbiasedness of estimator:

$$\mathbb{E}[\hat{F}(x)] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \mathbb{I}[x_i \leq x]\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbb{I}[x_i \leq x]] = \frac{1}{n} \sum_{i=1}^n F(x) = F(x).$$

- (b) Use the Taylor expansion  $F(x+h) = F(x) + hf(x) + \frac{1}{2}h^2f'(x) + o(h^2)$  to see that the bias of  $\hat{f}_1(x)$  is

$$\begin{aligned}\mathbb{E}[\hat{f}_1(x)] - f(x) &= h^{-1}(F(x+h) - F(x)) - f(x) \\ &= \frac{1}{h}(F(x) + hf(x) + \frac{1}{2}h^2f'(x) + o(h^2) - F(x)) - f(x) \\ &= \frac{1}{2}hf'(x) + o(h).\end{aligned}$$

Therefore,  $a = 1$ .

- (c) Use the Taylor expansions  $F(x + \frac{h}{2}) = F(x) + \frac{h}{2}f(x) + \frac{1}{2}(\frac{h}{2})^2f'(x) + \frac{1}{6}(\frac{h}{2})^3f''(x) + o(h^3)$  and  $F(x - \frac{h}{2}) = F(x) - \frac{h}{2}f(x) + \frac{1}{2}(\frac{h}{2})^2f'(x) - \frac{1}{6}(\frac{h}{2})^3f''(x) + o(h^3)$  to see that the bias of  $\hat{f}_2(x)$  is

$$\mathbb{E}[\hat{f}_2(x)] - f(x) = h^{-1}(F(x+h/2) - F(x-h/2)) - f(x) = \frac{1}{24}h^2f''(x) + o(h^2).$$

Therefore,  $b = 2$ .

## 14.3 Nadaraya–Watson density estimator

Recall that

$$\sqrt{n}(\hat{f}(x) - f(x)) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (K_h(x_i - x) - f(x)).$$

It is easy to derive that

$$\mathbb{E}[K_h(x_i - x) - f(x)] = O(h^2),$$

and it can similarly be shown that

$$\begin{aligned}\mathbb{V}[K_h(x_i - x) - f(x)] &= \mathbb{E}[K_h(x_i - x)^2] + f(x)^2 - 2f(x)\mathbb{E}[K_h(x_i - x)] - \mathbb{E}[K_h(x_i - x) - f(x)]^2 \\ &= h^{-1} \int K(u)^2 (f(x) + O(h)) du + O(1) - O(h^2)^2 \\ &= h^{-1}f(x)R_K + O(1).\end{aligned}$$

To get non-trivial asymptotic bias, we should work on the bias term further:

$$\begin{aligned}\mathbb{E}[K_h(x_i - x) - f(x)] &= \int K(u) \left( f(x) + hu f'(x) + \frac{1}{2}(hu)^2 f''(x) + O(h^3) \right) du - f(x) \\ &= h^2 \frac{1}{2} f''(x) \sigma_K^2 + O(h^3).\end{aligned}$$

Summarizing, by (some variant of) CLT we have, as  $n \rightarrow \infty$  and  $h \rightarrow 0$ ,  $nh \rightarrow \infty$ , that

$$\sqrt{nh}(\hat{f}(x) - f(x)) \xrightarrow{d} \mathcal{N}\left(\frac{1}{2}\lambda f''(x)\sigma_K^2, f(x)R_K\right),$$

provided that  $\lambda \equiv \lim_{n \rightarrow \infty} \sqrt{nh^5}$  exists and is finite.



The asymptotic bias is proportional to  $f''(x)$ , the amount of which indicates how the density in the neighborhood of  $x$  differs from the density at  $x$  which is estimated. Note that the asymptotic bias does not depend on  $f(x)$ , i.e. how often observations fall into this region, and on  $f'(x)$ , i.e. how density to the left and that to the right of  $x$  differ – indeed, both are equally irrelevant to estimation of the density (unlike in the regression). The asymptotic variance is proportional to  $f(x)$ , the density at  $x$ , which may seem counterintuitive (the higher frequency of observations, the poorer estimation quality). Recall, however, that we are estimating  $f(x)$ , so its higher value implies more dispersion of its estimate around that value, and this effect prevails (the frequency effect gives  $\propto f(x)^{-1}$ , the scale effect gives  $\propto f(x)^2$ ).

## 14.4 First difference transformation and nonparametric regression

1. Let us consider the following average that can be decomposed into three terms:

$$\begin{aligned} \frac{1}{n-1} \sum_{i=2}^n (y_i - y_{i-1})^2 &= \frac{1}{n-1} \sum_{i=2}^n (g(z_i) - g(z_{i-1}))^2 + \frac{1}{n-1} \sum_{i=2}^n (e_i - e_{i-1})^2 \\ &\quad + \frac{2}{n-1} \sum_{i=2}^n (g(z_i) - g(z_{i-1}))(e_i - e_{i-1}). \end{aligned}$$

Since  $z_i$  compose a uniform grid and are increasing in order, i.e.  $z_i - z_{i-1} = \frac{1}{n-1}$ , we can find the limit of the first term using the Lipschitz condition:

$$\left| \frac{1}{n-1} \sum_{i=2}^n (g(z_i) - g(z_{i-1}))^2 \right| \leq \frac{G^2}{n-1} \sum_{i=2}^n (z_i - z_{i-1})^2 = \frac{G^2}{(n-1)^2} \xrightarrow{n \rightarrow \infty} 0$$

Using the Lipschitz condition again we can find the probability limit of the third term:

$$\begin{aligned} \left| \frac{2}{n-1} \sum_{i=2}^n (g(z_i) - g(z_{i-1}))(e_i - e_{i-1}) \right| &\leq \frac{2G}{(n-1)^2} \sum_{i=2}^n |e_i - e_{i-1}| \\ &\leq \frac{2G}{n-1} \frac{1}{n-1} \sum_{i=2}^n (|e_i| + |e_{i-1}|) \xrightarrow[n \rightarrow \infty]{p} 0 \end{aligned}$$

since  $\frac{2G}{n-1} \xrightarrow{n \rightarrow \infty} 0$  and  $\frac{1}{n-1} \sum_{i=2}^n (|e_i| + |e_{i-1}|) \xrightarrow[n \rightarrow \infty]{p} 2\mathbb{E}|e_i| < \infty$ . The second term has the following probability limit:

$$\frac{1}{n-1} \sum_{i=2}^n (e_i - e_{i-1})^2 = \frac{1}{n-1} \sum_{i=2}^n (e_i^2 - 2e_i e_{i-1} + e_{i-1}^2) \xrightarrow[n \rightarrow \infty]{p} 2\mathbb{E}[e_i^2] = 2\sigma^2.$$

Thus the estimator for  $\sigma^2$  whose consistency is proved by previous manipulations is

$$\hat{\sigma}^2 = \frac{1}{2} \frac{1}{n-1} \sum_{i=2}^n (y_i - y_{i-1})^2.$$

2. At the first step estimate  $\beta$  from the FD-regression. The FD-transformed regression is

$$y_i - y_{i-1} = (x_i - x_{i-1})' \beta + g(z_i) - g(z_{i-1}) + e_i - e_{i-1},$$

which can be rewritten as

$$\Delta y_i = \Delta x_i' \beta + \Delta g(z_i) + \Delta e_i.$$

The consistency of the following estimator for  $\beta$

$$\hat{\beta} = \left( \sum_{i=2}^n \Delta x_i \Delta x_i' \right)^{-1} \left( \sum_{i=2}^n \Delta x_i \Delta y_i \right)$$

can be proved in the standard way:

$$\hat{\beta} - \beta = \left( \frac{1}{n-1} \sum_{i=2}^n \Delta x_i \Delta x_i' \right)^{-1} \left( \frac{1}{n-1} \sum_{i=2}^n \Delta x_i (\Delta g(z_i) + \Delta e_i) \right)$$

Here  $\frac{1}{n-1} \sum_{i=2}^n \Delta x_i \Delta x_i'$  has some non-zero probability limit,  $\frac{1}{n-1} \sum_{i=2}^n \Delta x_i \Delta e_i \xrightarrow{p} 0$  since  $\mathbb{E}[e_i | x_i, z_i] = 0$ , and  $\left| \frac{1}{n-1} \sum_{i=2}^n \Delta x_i \Delta g(z_i) \right| \leq \frac{G}{n-1} \frac{1}{n-1} \sum_{i=2}^n |\Delta x_i| \xrightarrow{p} 0$ . Now we can use standard nonparametric tools for the “regression”

$$y_i - x_i' \hat{\beta} = g(z_i) + e_i^*,$$

where  $e_i^* = e_i + x_i'(\beta - \hat{\beta})$ . Consider the following estimator (we use the uniform kernel for algebraic simplicity):

$$\widehat{g}(z) = \frac{\sum_{i=1}^n (y_i - x_i' \hat{\beta}) \mathbb{I}[|z_i - z| \leq h]}{\sum_{i=1}^n \mathbb{I}[|z_i - z| \leq h]}.$$

It can be decomposed into three terms:

$$\widehat{g}(z) = \frac{\sum_{i=1}^n \left( g(z_i) + x_i'(\beta - \hat{\beta}) + e_i \right) \mathbb{I}[|z_i - z| \leq h]}{\sum_{i=1}^n \mathbb{I}[|z_i - z| \leq h]}$$

The first term gives  $g(z)$  in the limit. To show this, use Lipschitz condition:

$$\left| \frac{\sum_{i=1}^n (g(z_i) - g(z)) \mathbb{I}[|z_i - z| \leq h]}{\sum_{i=1}^n \mathbb{I}[|z_i - z| \leq h]} \right| \leq Gh,$$

and introduce the asymptotics for the smoothing parameter:  $h \rightarrow 0$ . Then

$$\begin{aligned} \frac{\sum_{i=1}^n g(z_i) \mathbb{I}[|z_i - z| \leq h]}{\sum_{i=1}^n \mathbb{I}[|z_i - z| \leq h]} &= \frac{\sum_{i=1}^n (g(z) + g(z_i) - g(z)) \mathbb{I}[|z_i - z| \leq h]}{\sum_{i=1}^n \mathbb{I}[|z_i - z| \leq h]} = \\ &= g(z) + \frac{\sum_{i=1}^n (g(z_i) - g(z)) \mathbb{I}[|z_i - z| \leq h]}{\sum_{i=1}^n \mathbb{I}[|z_i - z| \leq h]} \xrightarrow{n \rightarrow \infty} g(z). \end{aligned}$$

The second and the third terms have zero probability limit if the condition  $nh \rightarrow \infty$  is satisfied

$$\underbrace{\frac{\sum_{i=1}^n x_i' \mathbb{I}[|z_i - z| \leq h]}{\sum_{i=1}^n \mathbb{I}[|z_i - z| \leq h]}}_{\downarrow^p \mathbb{E}[x_i']} \underbrace{(\beta - \hat{\beta})}_{\downarrow^p 0} \xrightarrow{n \rightarrow \infty} 0$$

and

$$\frac{\sum_{i=1}^n e_i \mathbb{I}[|z_i - z| \leq h]}{\sum_{i=1}^n \mathbb{I}[|z_i - z| \leq h]} \xrightarrow{n \rightarrow \infty} \mathbb{E}[e_i] = 0.$$

Therefore,  $\widehat{g}(z)$  is consistent when  $n \rightarrow \infty$ ,  $nh \rightarrow \infty$ ,  $h \rightarrow 0$ .

## 14.5 Unbiasedness of kernel estimates

Recall that

$$\hat{g}(x) = \frac{\sum_{i=1}^n y_i K_h(x_i - x)}{\sum_{i=1}^n K_h(x_i - x)},$$

so

$$\begin{aligned} \mathbb{E}[\hat{g}(x)] &= \mathbb{E}\left[\mathbb{E}\left[\frac{\sum_{i=1}^n y_i K_h(x_i - x)}{\sum_{i=1}^n K_h(x_i - x)} \mid x_1, \dots, x_n\right]\right] \\ &= \mathbb{E}\left[\frac{\sum_{i=1}^n \mathbb{E}[y_i | x_i] K_h(x_i - x)}{\sum_{i=1}^n K_h(x_i - x)}\right] = \mathbb{E}\left[\frac{\sum_{i=1}^n c K_h(x_i - x)}{\sum_{i=1}^n K_h(x_i - x)}\right] = c, \end{aligned}$$

i.e.  $\hat{g}(x)$  is unbiased for  $c = g(x)$ . The reason is simple: all points in the sample are equally relevant in estimation of this trivial conditional mean, so bias is not induced when points far from  $x$  are used in estimation.

The local linear estimator will be unbiased if  $g(x) = a + bx$ . Then all points in the sample are equally relevant in estimation since it is a linear regression, albeit locally, is run. Indeed,

$$\hat{g}_1(x) = \bar{y} + \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) K_h(x_i - x)}{\sum_{i=1}^n (x_i - \bar{x})^2 K_h(x_i - x)} (x - \bar{x}),$$

so

$$\begin{aligned} \mathbb{E}[\hat{g}_1(x)] &= \mathbb{E}[\mathbb{E}[\bar{y} | x_1, \dots, x_n]] \\ &\quad + \mathbb{E}\left[\mathbb{E}\left[\frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) K_h(x_i - x)}{\sum_{i=1}^n (x_i - \bar{x})^2 K_h(x_i - x)} \mid x_1, \dots, x_n\right] (x - \bar{x})\right] \\ &= \mathbb{E}[a + b\bar{x}] \\ &\quad + \mathbb{E}\left[\mathbb{E}\left[\frac{\sum_{i=1}^n (a + bx_i - a - b\bar{x})(x_i - \bar{x}) K_h(x_i - x)}{\sum_{i=1}^n (x_i - \bar{x})^2 K_h(x_i - x)} \mid x_1, \dots, x_n\right] (x - \bar{x})\right] \\ &= \mathbb{E}[a + b\bar{x} + b(x - \bar{x})] = a + bx. \end{aligned}$$

As far as the density is concerned, unbiasedness is unlikely. Indeed, recall that

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n K_h(x_i - x),$$

so

$$\mathbb{E}[\hat{f}(x)] = \mathbb{E}[K_h(x_i - x)] = \frac{1}{h} \int K\left(\frac{x_i - x}{h}\right) f(x) dx.$$

This expectation heavily depends on the bandwidth and kernel function, and barely will it equal  $f(x)$ , except under special circumstances (e.g., uniform  $f(x)$ ,  $x$  far from boundaries, etc.).

## 14.6 Shape restriction

The CRS technology has the property that

$$f(l, k) = kf\left(\frac{l}{k}, 1\right).$$

The regression in terms of the rescaled variables is

$$\frac{y_i}{k_i} = f\left(\frac{l_i}{k_i}, 1\right) + \frac{\varepsilon_i}{k_i}.$$

Therefore, we can construct the (one dimensional!) kernel estimate of  $f(l, k)$  as

$$\hat{f}(l, k) = k \times \frac{\sum_{i=1}^n \frac{y_i}{k_i} K_h\left(\frac{l_i}{k_i} - \frac{l}{k}\right)}{\sum_{i=1}^n K_h\left(\frac{l_i}{k_i} - \frac{l}{k}\right)}.$$

In effect, we are using the sample points giving higher weight to those that are close to the ray  $l/k$ .

## 14.7 Nonparametric hazard rate

(i) A simple nonparametric estimator for  $F(t) \equiv \Pr\{z \leq t\}$  is the sample frequency

$$\hat{F}(t) = \frac{1}{n} \sum_{j=1}^n \mathbb{I}[z_j \leq t].$$

By the law of large numbers, it is consistent for  $F(t)$ . By the central limit theorem, its rate of convergence is  $\sqrt{n}$ . This will be helpful later.

(ii) We derived in class that

$$\mathbb{E}\left[\frac{1}{h_n} k \left(\frac{z_j - t}{h_n}\right) - f(t)\right] = O(h_n^2)$$

and

$$\mathbb{V}\left[\frac{1}{h_n} k \left(\frac{z_j - t}{h_n}\right)\right] = \frac{1}{h_n} R_k f(t) + O(1),$$

where  $R_k = \int k(u)^2 du$ . By the central limit theorem applied to

$$\sqrt{nh_n} (\hat{f}(t) - f(t)) = \sqrt{h_n} \frac{1}{\sqrt{n}} \sum_{j=1}^n \left(\frac{1}{h_n} k \left(\frac{z_j - t}{h_n}\right) - f(t)\right),$$

we get

$$\sqrt{nh_n} (\hat{f}(t) - f(t)) \xrightarrow{d} \mathcal{N}(0, R_k f(t)).$$

(iii) By the analogy principle,

$$\hat{H}(t) = \frac{\hat{f}(t)}{1 - \hat{F}(t)}.$$

It is consistent for  $H(t)$  and

$$\begin{aligned}
\sqrt{nh_n} \left( \hat{H}(t) - H(t) \right) &= \sqrt{nh_n} \left( \frac{\hat{f}(t)}{1 - \hat{F}(t)} - \frac{f(t)}{1 - F(t)} \right) \\
&= \sqrt{nh_n} \left( \frac{\hat{f}(t)(1 - F(t)) - f(t)(1 - \hat{F}(t))}{(1 - \hat{F}(t))(1 - F(t))} \right) \\
&= \frac{\sqrt{nh_n}(\hat{f}(t) - f(t))}{1 - \hat{F}(t)} + \sqrt{h_n}f(t) \frac{\sqrt{n}(\hat{F}(t) - F(t))}{(1 - \hat{F}(t))(1 - F(t))} \\
&\xrightarrow{d} \frac{1}{1 - F(t)} \mathcal{N}(0, R_k f(t)) + 0 \\
&= \mathcal{N} \left( 0, R_k \frac{f(t)}{(1 - F(t))^2} \right).
\end{aligned}$$

The reason of the fact that uncertainty in  $\hat{F}(t)$  does not affect the asymptotic distribution of  $\hat{H}(t)$  is that  $\hat{F}(t)$  converges with faster rate than  $\hat{f}(t)$  does.

## 14.8 Nonparametrics and perfect fit

When the variance of the error is zero,

$$\hat{g}(x) - g(x) = \frac{(nh)^{-1} \sum_{i=1}^n (g(x_i) - g(x)) K \left( \frac{x_i - x}{h} \right)}{(nh)^{-1} \sum_{i=1}^n K \left( \frac{x_i - x}{h} \right)}.$$

There is no usual source of variance (regression errors), so the variance should come from the variance of  $x_i$ 's.

The denominator converges to  $f(x)$  when  $n \rightarrow \infty$ ,  $h \rightarrow 0$ . Consider the numerator, which we denote by  $\hat{q}(x)$ , and which is an average of IID random variables, say  $\zeta_i$ . We derived in class that

$$\mathbb{E}[\hat{q}(x)] = h^2 B(x) f(x) + o(h^2), \quad \mathbb{V}[\hat{q}(x)] = o((nh)^{-1}).$$

Now we need to look closer at the variance. Now,

$$\begin{aligned}
\mathbb{E}[\zeta_i^2] &= h^{-2} \int (g(x_i) - g(x))^2 K \left( \frac{x_i - x}{h} \right)^2 f(x_i) dx_i \\
&= h^{-1} \int (g(x + hu) - g(x))^2 K(u)^2 f(x + hu) du \\
&= h^{-1} \int (g'(x)hu + o(h))^2 K(u)^2 (f(x) + o(h)) du \\
&= hg'(x)^2 f(x) \int u^2 K(u)^2 du + o(h),
\end{aligned}$$

so

$$\mathbb{V}[\zeta_i] = \mathbb{E}[\zeta_i^2] - \mathbb{E}[\zeta_i]^2 = hg'(x)^2 f(x) \Psi_K^2 + o(h),$$

where  $\Psi_K^2 = \int u^2 K(u)^2 du$ . Hence, by some CLT,

$$\begin{aligned} & \sqrt{nh^{-1}} \left( (nh)^{-1} \sum_{i=1}^n (g(x_i) - g(x)) K\left(\frac{x_i - x}{h}\right) - h^2 B(x) f(x) + o(h^2) \right) \\ & \xrightarrow{d} \mathcal{N}\left(0, g'(x)^2 f(x) \Psi_K^2\right). \end{aligned}$$

Let  $\lambda = \lim_{n \rightarrow \infty, h \rightarrow 0} \sqrt{nh^3}$ , assuming  $\lambda < \infty$ . Then,

$$\sqrt{nh^{-1}} (\hat{g}(x) - g(x)) \xrightarrow{d} \mathcal{N}\left(\lambda B(x), \frac{g'(x)^2}{f(x)} \Psi_K^2\right).$$

## 14.9 Nonparametrics and extreme observations

- (a) When  $x_1$  is very big, we will have a big influence of this observation when estimating the regression in the right region of the support. With bounded kernels, we may find that no observations fall into the window. With unbounded kernels, the estimated regression line will go almost through  $(x_1, y_1)$ .
- (b) When  $y_1$  is very big, the estimated regression line will have a hump in the region centered at  $x_1$ .

# 15. CONDITIONAL MOMENT RESTRICTIONS

## 15.1 Usefulness of skedastic function

Denote  $\theta = \begin{pmatrix} \beta \\ \pi \end{pmatrix}$  and  $e = y - x'\beta$ . The moment function is

$$m(y, x, \theta) = \begin{pmatrix} m_1 \\ m_2 \end{pmatrix} = \begin{pmatrix} y - x'\beta \\ (y - x'\beta)^2 - h(x, \beta, \pi) \end{pmatrix}$$

The general theory for the conditional moment restriction  $\mathbb{E}[m(w, \theta)|x] = 0$  gives the optimal restriction  $\mathbb{E}[D(x)'\Omega(x)^{-1}m(w, \theta)] = 0$ , where  $D(x) = \mathbb{E}\left[\frac{\partial m}{\partial \theta'}|x\right]$  and  $\Omega(x) = \mathbb{E}[mm'|x]$ . The variance of the optimal estimator is  $V = (\mathbb{E}[D(x)'\Omega(x)^{-1}D(x)])^{-1}$ . For the problem at hand,

$$D(x) = \mathbb{E}\left[\frac{\partial m}{\partial \theta'}|x\right] = -\mathbb{E}\left[\begin{pmatrix} x' & 0 \\ 2ex' + h'_\beta & h'_\pi \end{pmatrix}|x\right] = -\begin{pmatrix} x' & 0 \\ h'_\beta & h'_\pi \end{pmatrix},$$

$$\Omega(x) = \mathbb{E}[mm'|x] = \mathbb{E}\left[\begin{pmatrix} e^2 & e(e^2 - h) \\ e(e^2 - h) & (e^2 - h)^2 \end{pmatrix}|x\right] = \mathbb{E}\left[\begin{pmatrix} e^2 & e^3 \\ e^3 & (e^2 - h)^2 \end{pmatrix}|x\right],$$

since  $\mathbb{E}[ex|x] = 0$  and  $\mathbb{E}[eh|x] = 0$ .

Let  $\Delta(x) \equiv \det \Omega(x) = \mathbb{E}[e^2|x]\mathbb{E}[(e^2 - h)^2|x] - (\mathbb{E}[e^3|x])^2$ . The inverse of  $\Omega$  is

$$\Omega(x)^{-1} = \frac{1}{\Delta(x)} \mathbb{E}\left[\begin{pmatrix} (e^2 - h)^2 & -e^3 \\ -e^3 & e^2 \end{pmatrix}|x\right],$$

and the asymptotic variance of the efficient GMM estimator is

$$V^{-1} = \mathbb{E}[D(x)'\Omega(x)^{-1}D(x)] = \begin{pmatrix} A & B' \\ B & C \end{pmatrix},$$

where

$$A = \mathbb{E}\left[\frac{(e^2 - h)^2 xx' - e^3(xh'_\beta + h_\beta x') + e^2 h_\beta h'_\beta}{\Delta(x)}\right],$$

$$B = \mathbb{E}\left[\frac{-e^3 h_\pi x' + e^2 h_\pi h'_\beta}{\Delta(x)}\right], \quad C = \mathbb{E}\left[\frac{e^2 h_\pi h'_\pi}{\Delta(x)}\right].$$

Using the formula for inversion of the partitioned matrices, find that

$$V = \begin{pmatrix} (A - B'C^{-1}B)^{-1} & * \\ * & * \end{pmatrix},$$

where \*'s denote submatrices which are not of interest.

To answer the problem we need to compare  $V_{11} = (A - B'C^{-1}B)^{-1}$  with  $V_0 = \left(\mathbb{E}\left[\frac{xx'}{h}\right]\right)^{-1}$ , the variance of the optimal GMM estimator constructed with the use of  $m_1$  only. We need to show that  $V_{11} \leq V_0$ , or, alternatively,  $V_{11}^{-1} \geq V_0^{-1}$ . Note that

$$V_{11}^{-1} - V_0^{-1} = \tilde{A} - B'C^{-1}B,$$

where  $\tilde{A} = A - V_0^{-1}$  can be simplified to

$$\tilde{A} = \mathbb{E} \left[ \frac{1}{\Delta(x)} \left( \frac{xx' (\mathbb{E}[e^3|x])^2}{\mathbb{E}[e^2|x]} - e^3(xh'_\beta + h_\beta x') + e^2 h_\beta h'_\beta \right) \right].$$

Next, we can use the following representation:

$$\tilde{A} - B' C^{-1} B = \mathbb{E}[ww'],$$

where

$$w = \frac{\mathbb{E}[e^3|x] x - \mathbb{E}[e^2|x] h_\beta}{\sqrt{\mathbb{E}[e^2|x]} \sqrt{\Delta(x)}} + B' C^{-1} h_\pi \sqrt{\frac{\mathbb{E}[e^2|x]}{\Delta(x)}}.$$

This representation concludes that  $V_{11}^{-1} \geq V_0^{-1}$  and gives the condition under which  $V_{11} = V_0$ . This condition is  $w(x) = 0$  almost surely. It can be written as

$$\frac{\mathbb{E}[e^3|x]}{\mathbb{E}[e^2|x]} x = h_\beta - B' C^{-1} h_\pi \text{ almost surely.}$$

Consider the special cases.

1.  $h_\beta = 0$ . Then the condition modifies to

$$\frac{\mathbb{E}[e^3|x]}{\mathbb{E}[e^2|x]} x = -\mathbb{E} \left[ \frac{e^3 h_\pi x'}{\Delta(x)} \right] \mathbb{E} \left[ \frac{e^2 h_\pi h'_\pi}{\Delta(x)} \right]^{-1} h_\pi \text{ almost surely.}$$

2.  $h_\beta = 0$  and the distribution of  $e$  conditional on  $x$  is symmetric. Then the previous condition is satisfied automatically since  $\mathbb{E}[e^3|x] = 0$ .

## 15.2 Symmetric regression error

Part 1. The maintained hypothesis is  $\mathbb{E}[e|x] = 0$ . We can use the null hypothesis  $H_0 : \mathbb{E}[e^3|x] = 0$  to test for the conditional symmetry. We could in addition use more conditional moment restrictions (e.g., involving higher odd powers) to increase the power of the test, but in finite samples that would probably lead to more distorted test sizes. The alternative hypothesis is  $H_1 : \mathbb{E}[e^3|x] \neq 0$ .

An estimator that is consistent under both  $H_0$  and  $H_1$  is, for example, the OLS estimator  $\hat{\alpha}_{OLS}$ . The estimator that is consistent and asymptotically efficient (in the same class where  $\hat{\alpha}_{OLS}$  belongs) under  $H_0$  and (hopefully) inconsistent under  $H_1$  is the instrumental variables (GMM) estimator  $\hat{\alpha}_{IV}$  that uses the optimal instrument for the system  $\mathbb{E}[e|x] = 0$ ,  $\mathbb{E}[e^3|x] = 0$ . We derived in class that the optimal unconditional moment restriction is

$$\mathbb{E} \left[ a_1(x) (y - \alpha x) + a_2(x) (y - \alpha x)^3 \right] = 0,$$

where

$$\begin{pmatrix} a_1(x) \\ a_2(x) \end{pmatrix} = \frac{x}{\mu_2(x)\mu_6(x) - \mu_4(x)^2} \begin{pmatrix} \mu_6(x) - 3\mu_2(x)\mu_4(x) \\ 3\mu_2(x)^2 - \mu_4(x) \end{pmatrix}$$

and  $\mu_r(x) = \mathbb{E}[(y - \alpha x)^r | x]$ ,  $r = 2, 4, 6$ . To construct a feasible  $\hat{\alpha}_{IV}$ , one needs to first estimate  $\mu_r(x)$  at the points  $x_i$  of the sample. This may be done nonparametrically using nearest neighbor,



series expansion or other approaches. Denote the resulting estimates by  $\hat{\mu}_r(x_i)$ ,  $i = 1, \dots, n$ ,  $r = 2, 4, 6$  and compute  $\hat{a}_1(x_i)$  and  $\hat{a}_2(x_i)$ ,  $i = 1, \dots, n$ . Then  $\hat{\alpha}_{OIV}$  is a solution of the equation

$$\frac{1}{n} \sum_{i=1}^n \left( \hat{a}_1(x_i) (y_i - \hat{\alpha}_{OIV} x_i) + \hat{a}_2(x_i) (y_i - \hat{\alpha}_{OIV} x_i)^3 \right) = 0,$$

which can be turned into an optimization problem, if convenient.

The Hausman test statistic is then

$$H = n \frac{(\hat{\alpha}_{OLS} - \hat{\alpha}_{OIV})^2}{\hat{V}_{OLS} - \hat{V}_{OIV}} \xrightarrow{d} \chi_1^2,$$

where  $\hat{V}_{OLS} = n \left( \sum_{i=1}^n x_i^2 \right)^{-2} \sum_{i=1}^n x_i^2 (y_i - \hat{\alpha}_{OLS} x_i)^2$  and  $\hat{V}_{OIV}$  is a consistent estimate of the efficiency bound

$$V_{OIV} = \left( \mathbb{E} \left[ \frac{x^2 (\mu_6(x) - 6\mu_2(x)\mu_4(x) + 9\mu_2^3(x))}{\mu_2(x)\mu_6(x) - \mu_4^2(x)} \right] \right)^{-1}.$$

Note that the constructed Hausman test will not work if  $\hat{\alpha}_{OLS}$  is also asymptotically efficient, which may happen if the third-moment restriction is redundant and the error is conditionally homoskedastic so that the optimal instrument reduces to the one implied by OLS. Also, the test may be inconsistent (i.e., asymptotically have power less than 1) if  $\hat{\alpha}_{OIV}$  happens to be consistent under conditional non-symmetry too.

Part 2. Under the assumption that  $e|x \sim \mathcal{N}(0, \sigma^2)$ , irrespective of whether  $\sigma^2$  is known or not, the QML estimator  $\hat{\alpha}_{QML}$  coincides with the OLS estimator and thus has the same asymptotic distribution

$$\sqrt{n}(\hat{\alpha}_{QML} - \alpha) \xrightarrow{d} \mathcal{N} \left( 0, \frac{\mathbb{E} [x^2 (y - \alpha x)^2]}{(\mathbb{E} [x^2])^2} \right).$$

## 15.3 Optimal instrumentation of consumption function

As the error has a martingale difference structure, the optimal instrument is (up to a factor of proportionality)

$$\zeta_t \propto \frac{1}{\mathbb{E} [e_t^2 | y_{t-1}, c_{t-1}, \dots]} \left( \begin{array}{c} 1 \\ \mathbb{E} [y_t^\gamma | y_{t-1}, c_{t-1}, \dots] \\ \mathbb{E} [y_t^\gamma \ln y_t | y_{t-1}, c_{t-1}, \dots] \end{array} \right).$$

One could first estimate  $\alpha$ ,  $\delta$  and  $\gamma$  using the instrument vector (for example)  $(1, y_{t-1}, y_{t-2})$ . Using these, one could get feasible versions of  $e_t^2$  and  $y_t^\gamma$ . Then one could do one's best to parametrically fit feasible  $e_t^2$ ,  $y_t^\gamma$  and  $y_t^\gamma \ln y_t$  to recent variables in  $\{y_{t-1}, c_{t-1}, y_{t-2}, c_{t-2}, \dots\}$ , employing autoregressive structures (e.g., fancy GARCH for  $e_t^2$ ). The fitted values then could be used to form feasible  $\zeta_t$  to be eventually used to IV-estimate the parameters. Naturally, because auxiliary parameterizations are used, such instrument will be only "nearly optimal". In particular, the asymptotic variance has to be computed in a "sandwich" form.

## 15.4 Optimal instrument in AR-ARCH model

Let us for convenience view a typical element of  $\mathcal{Z}_t$  as  $\sum_{i=1}^{\infty} \omega_i \varepsilon_{t-i}$ , and let the optimal instrument be  $\zeta_t = \sum_{i=1}^{\infty} a_i \varepsilon_{t-i}$ . The optimality condition is

$$\mathbb{E}[v_t x_{t-1}] = \mathbb{E}[v_t \zeta_t \varepsilon_t^2] \quad \text{for all } v_t \in \mathcal{Z}_t.$$

Since it should hold for any  $v_t \in \mathcal{Z}_t$ , let us make it hold for  $v_t = \varepsilon_{t-j}$ ,  $j = 1, 2, \dots$ . Then we get a system of equations of the type

$$\mathbb{E}[\varepsilon_{t-j} x_{t-1}] = \mathbb{E}\left[\varepsilon_{t-j} \left(\sum_{i=1}^{\infty} a_i \varepsilon_{t-i}\right) \varepsilon_t^2\right].$$

The left-hand side is just  $\rho^{j-1}$  because  $x_{t-1} = \sum_{i=1}^{\infty} \rho^{i-1} \varepsilon_{t-i}$  and because  $\mathbb{E}[\varepsilon_t^2] = 1$ . In the right-hand side, all terms are zeros due to conditional symmetry of  $\varepsilon_t$ , except  $a_j \mathbb{E}\left[\varepsilon_{t-j}^2 \varepsilon_t^2\right]$ . Therefore,

$$a_j = \frac{\rho^{j-1}}{1 + \alpha^j (\kappa - 1)},$$

where  $\kappa = \mathbb{E}[\varepsilon_t^4]$ . This follows from the *ARCH*(1) structure:

$$\mathbb{E}[\varepsilon_{t-j}^2 \varepsilon_t^2] = \mathbb{E}[\varepsilon_{t-j}^2 \mathbb{E}[\varepsilon_t^2 | I_{t-1}]] = \mathbb{E}[\varepsilon_{t-j}^2 ((1 - \alpha) + \alpha \varepsilon_{t-1}^2)] = (1 - \alpha) + \alpha \mathbb{E}[\varepsilon_{t-j+1}^2 \varepsilon_t^2],$$

so that we can recursively obtain

$$\mathbb{E}[\varepsilon_{t-j}^2 \varepsilon_t^2] = 1 - \alpha^j + \alpha^j \kappa.$$

Thus the optimal instrument is

$$\begin{aligned} \zeta_t &= \sum_{i=1}^{\infty} \frac{\rho^{i-1}}{1 + \alpha^i (\kappa - 1)} \varepsilon_{t-i} = \\ &= \frac{x_{t-1}}{1 + \alpha(\kappa - 1)} + (\kappa - 1)(1 - \alpha) \sum_{i=2}^{\infty} \frac{(\alpha\rho)^{i-1}}{(1 + \alpha^i (\kappa - 1))(1 + \alpha^{i-1} (\kappa - 1))} x_{t-i}. \end{aligned}$$

To construct a feasible estimator, set  $\hat{\rho}$  to be the OLS estimator of  $\rho$ ,  $\hat{\alpha}$  to be the OLS estimator of  $\alpha$  in the model  $\hat{\varepsilon}_t^2 - 1 = \alpha(\hat{\varepsilon}_{t-1}^2 - 1) + v_t$ , and compute  $\hat{\kappa} = T^{-1} \sum_{t=2}^T \hat{\varepsilon}_t^4$ .

The optimal instrument based on  $\mathbb{E}[\varepsilon_t | I_{t-1}] = 0$  uses a large set of allowable instruments, relative to which our  $\mathcal{Z}_t$  is extremely thin. Therefore, we can expect big losses in efficiency in comparison with what we could get. In fact, calculations for empirically relevant sets of parameter values reveal that this intuition is correct. Weighting by the skedastic function is much more powerful than trying to capture heteroskedasticity by using an infinite history of the basic instrument in a linear fashion.

## 15.5 Optimal instrument in AR with nonlinear error

1. We look for an optimal combination of  $\varepsilon_{t-1}, \varepsilon_{t-2}, \dots$ , say

$$\zeta_t = \sum_{i=1}^{\infty} \phi_i^* \varepsilon_{t-i},$$

when the class of allowable instruments is

$$\mathcal{Z}_t = \left\{ \zeta_t = \sum_{i=1}^{\infty} \phi_i \varepsilon_{t-i} \right\}.$$

The optimality condition for each  $\varepsilon_{t-r}$ ,  $r \geq 1$ , is

$$\forall r \geq 1 \quad \mathbb{E} [\varepsilon_{t-r} x_{t-1}] = \mathbb{E} [\varepsilon_{t-r} \zeta_t \varepsilon_t^2],$$

or

$$\forall r \geq 1 \quad \rho^{r-1} = \mathbb{E} \left[ \varepsilon_{t-r} \left( \sum_{i=1}^{\infty} \phi_i^* \varepsilon_{t-i} \right) \varepsilon_t^2 \right].$$

When  $r > 1$ , this implies  $\rho^{r-1} = \phi_r^*$ , while when  $r = 1$ , we have  $1 = \phi_1^* \mathbb{E} [\eta_{t-1}^2 \eta_{t-2}^2 \eta_t^2 \eta_{t-1}^2] = \phi_1^* \mathbb{E} [\eta^4]$ . The optimal instrument then is

$$\begin{aligned} \zeta_t &= \mathbb{E} [\eta^4]^{-1} \varepsilon_{t-1} + \sum_{i=2}^{\infty} \rho^{r-1} \varepsilon_{t-i} = \left( \mathbb{E} [\eta^4]^{-1} - 1 \right) \varepsilon_{t-1} + \sum_{i=1}^{\infty} \rho^{r-1} \varepsilon_{t-i} \\ &= \left( \mathbb{E} [\eta^4]^{-1} - 1 \right) (x_{t-1} - \rho x_{t-2}) + x_{t-1} = \mathbb{E} [\eta^4]^{-1} x_{t-1} + \left( \mathbb{E} [\eta^4]^{-1} - 1 \right) \rho x_{t-2}, \end{aligned}$$

and employs only two lags of  $x_t$ . To construct a feasible version, one needs to first consistently estimate  $\rho$  (say, by OLS), and  $\mathbb{E} [\eta^4]$  (say, by a sample analog to  $\mathbb{E} [\varepsilon_{t-1}^2 \varepsilon_t^2]$ ).

2. The optimal instrument based on this conditional moment restriction is Chamberlain's one:

$$\zeta_t \propto \frac{x_{t-1}}{\mathbb{E} [\varepsilon_t^2 | I_{t-1}]},$$

where

$$\begin{aligned} \mathbb{E} [\varepsilon_t^2 | I_{t-1}] &= \mathbb{E} [\eta_t^2 \eta_{t-1}^2 | x_{t-1}, x_{t-2}, \dots] = \mathbb{E} [\mathbb{E} [\eta_t^2 | \eta_{t-1}, x_{t-1}, x_{t-2}, \dots] \eta_{t-1}^2 | x_{t-1}, x_{t-2}, \dots] \\ &= \mathbb{E} [\eta_{t-1}^2 | x_{t-1}, x_{t-2}, \dots] = \mathbb{E} \left[ \frac{\varepsilon_{t-1}^2}{\eta_{t-2}^2} | x_{t-1}, x_{t-2}, \dots \right] = \dots = \frac{\varepsilon_{t-1}^2 \varepsilon_{t-3}^2 \varepsilon_{t-5}^2 \dots}{\varepsilon_{t-2}^2 \varepsilon_{t-4}^2 \varepsilon_{t-6}^2 \dots}. \end{aligned}$$

This, apparently, requires knowledge of all past errors, and not only of their finite number. Hence, it is problematic to construct a feasible version from a finite sample.

## 15.6 Optimal IV estimation of a constant

From the DGP it follows that the moment function is conditionally (on  $y_{t-p-1}, y_{t-p-2}, \dots$ ) homoskedastic. Therefore, the optimal instrument is Hansen's (1985)

$$\Theta(L)\zeta_t = \mathbb{E} [\Theta(L^{-1})^{-1} 1 | y_{t-p-1}, y_{t-p-2}, \dots],$$

or

$$\Theta(L)\zeta_t = \Theta(1)^{-1}.$$

This is a deterministic recursion. Since the instrument we are looking for should be stationary,  $\zeta_t$  has to be a constant. Since the value of the constant does not matter, the optimal instrument may be taken as unity.

## 15.7 Negative binomial distribution and PML

Yes, it does. The density belongs to the linear exponential class, with  $C(m) = \log m - \log(a + m)$ . The form of the skedastic function is immaterial for the consistency property to hold.

## 15.8 Nesting and PML

Indeed, in the example,  $h(z, \mu, \varsigma)$  reduces to exponential distribution when  $\varsigma = 1$ , and the exponential pseudodensity consistently estimates  $\theta_0$ . When  $\varsigma$  is estimated, the vector of parameters is  $(\mu, \varsigma)'$ , and the pseudoscore is

$$\begin{aligned} \frac{\partial \log h(z, \mu, \varsigma)}{\partial (\mu, \varsigma)'} &= \frac{\partial (\log \varsigma + \varsigma \log \Gamma(1 + \varsigma^{-1}) - \varsigma \log \mu + (\varsigma - 1) \log z - (\Gamma(1 + \varsigma^{-1}) z/\mu)^\varsigma)}{\partial (\mu, \varsigma)'} \\ &= \begin{pmatrix} ((\Gamma(1 + \varsigma^{-1}) z/\mu)^\varsigma - 1) \varsigma/\mu \\ \dots \end{pmatrix}. \end{aligned}$$

Observe that the pseudoscore for  $\mu$  has zero expectation if and only if  $\mathbb{E}[(\Gamma(1 + \varsigma^{-1}) z/\mu)^\varsigma] = 1$ . This obviously holds when  $\varsigma = 1$ , but may not hold for other  $\varsigma$ . For instance, when  $\varsigma = 2$ , we know that  $\mathbb{E}[z^2] = \mu^2 \Gamma(2.5) / \Gamma(1.5)^2 = 1.5\mu^2 / \Gamma(1.5)$ , which contradicts the zero expected pseudoscore. The pseudotrue value  $\varsigma_*$  of  $\varsigma$  can be obtained by solving the system of zero expected pseudoscore, but it is very unlikely that it equals 1. The result can be explained by the presence of an extraneous parameter whose pseudotrue value has nothing to do with the problem and whose estimation adversely affects estimation of the quantity of interest.

## 15.9 Misspecification in variance

The log pseudodensity on which the proposed PML1 estimator relies has the form

$$\log(y, m) = -\log \sqrt{2\pi} - \frac{1}{2} \log m^2 + \frac{(y - m)^2}{2m^2},$$

which does *not* belong to the linear exponential family of densities (the term  $y^2/2m^2$  does not fit). Therefore, the PML1 estimator is not consistent except by improbable chance.

The inconsistency can be shown directly. Consider a special case of no regressors and estimation of mean:

$$y \sim \mathcal{N}(\theta_0, \sigma^2),$$

while the pseudodensity is

$$y \sim \mathcal{N}(\theta, \theta^2).$$

Then the pseudotrue value of  $\theta$  is

$$\theta_* = \arg \max_{\theta} \mathbb{E} \left[ -\frac{1}{2} \log \theta^2 - \frac{(y - \theta)^2}{2\theta^2} \right] = \arg \max_{\theta} \left\{ -\frac{1}{2} \log \theta^2 - \frac{\sigma^2 + (\theta_0 - \theta)^2}{2\theta^2} \right\}.$$

It is easy to see by differentiating that  $\theta_*$  is not  $\theta_0$  until by chance  $\theta_0^2 = \sigma^2$ .

## 15.10 Modified Poisson regression and PML estimators

Part 1. The mean regression function is  $\mathbb{E}[y|x] = \mathbb{E}[\mathbb{E}[y|x, \varepsilon]|x] = \mathbb{E}[\exp(x'\beta + \varepsilon)|x] = \exp(x'\beta)$ . The skedastic function is  $\mathbb{V}[y|x] = \mathbb{E}[(y - \mathbb{E}[y|x])^2|x] = \mathbb{E}[y^2|x] - \mathbb{E}[y|x]^2$ . Because

$$\begin{aligned}\mathbb{E}[y^2|x] &= \mathbb{E}[\mathbb{E}[y^2|x, \varepsilon]|x] = \mathbb{E}[\exp(2x'\beta + 2\varepsilon) + \exp(x'\beta + \varepsilon)|x] \\ &= \exp(2x'\beta)\mathbb{E}[(\exp \varepsilon)^2|x] + \exp(x'\beta) = (\sigma^2 + 1)\exp(2x'\beta) + \exp(x'\beta),\end{aligned}$$

we have  $\mathbb{V}[y|x] = \sigma^2 \exp(2x'\beta) + \exp(x'\beta)$ .

Part 2. Use the formula for asymptotic variance of NLLS estimator  $V_{NLLS} = Q_{gg}^{-1}Q_{gge^2}Q_{gg}^{-1}$ , where  $Q_{gg} = \mathbb{E}[\partial g(x, \beta)/\partial \beta \cdot \partial g(x, \beta)/\partial \beta']$  and  $Q_{gge^2} = \mathbb{E}[\partial g(x, \beta)/\partial \beta \cdot \partial g(x, \beta)/\partial \beta' (y - g(x, \beta))^2]$ . In our problem  $g(x, \beta) = \exp(x'\beta)$  and  $Q_{gg} = \mathbb{E}[xx' \exp(2x'\beta)]$ ,

$$\begin{aligned}Q_{gge^2} &= \mathbb{E}[xx' \exp(2x'\beta)(y - \exp(x'\beta))^2] = \mathbb{E}[xx' \exp(2x'\beta)\mathbb{V}[y|x]] \\ &= \mathbb{E}[xx' \exp(2x'\beta)(\sigma^2 \exp(2x'\beta) + \exp(x'\beta))] = \mathbb{E}[xx' \exp(3x'\beta)] + \sigma^2 \mathbb{E}[xx' \exp(4x'\beta)].\end{aligned}$$

To find the expectations we use the formula  $\mathbb{E}[xx' \exp(nx'\beta)] = \exp(\frac{n^2}{2}\beta'\beta)(I + n^2\beta\beta')$ . Now, we have  $Q_{gg} = \exp(2\beta'\beta)(I + 4\beta\beta')$  and  $Q_{gge^2} = \exp(\frac{9}{2}\beta'\beta)(I + 9\beta\beta') + \sigma^2 \exp(8\beta'\beta)(I + 16\beta\beta')$ . Finally,

$$V_{NLLS} = (I + 4\beta\beta')^{-1} \left( \exp(\frac{1}{2}\beta'\beta)(I + 9\beta\beta') + \sigma^2 \exp(4\beta'\beta)(I + 16\beta\beta') \right) (I + 4\beta\beta')^{-1}.$$

The formula for asymptotic variance of WNLLS estimator is  $V_{WNLLS} = Q_{gg/\sigma^2}^{-1}$ , where  $Q_{gg/\sigma^2} = \mathbb{E}[\mathbb{V}[y|x]^{-1} \partial g(x, \beta)/\partial \beta \cdot \partial g(x, \beta)/\partial \beta']$ . In this problem

$$Q_{gg/\sigma^2} = \mathbb{E}[xx' \exp(2x'\beta)(\sigma^2 \exp(2x'\beta) + \exp(x'\beta))^{-1}],$$

which can be rearranged as

$$V_{WNLLS} = \sigma^2 \left( I - \mathbb{E} \left[ \frac{xx'}{1 + \sigma^2 \exp(x'\beta)} \right] \right)^{-1}.$$

Part 3. We use the formula for asymptotic variance of PML estimator:  $V_{PML} = \mathcal{J}^{-1} \mathcal{I} \mathcal{J}^{-1}$ , where

$$\begin{aligned}\mathcal{J} &= \mathbb{E} \left[ \frac{\partial C}{\partial m} \Big|_{m(x, \beta_0)} \frac{\partial m(x, \beta_0)}{\partial \beta} \frac{\partial m(x, \beta_0)}{\partial \beta'} \right], \\ \mathcal{I} &= \mathbb{E} \left[ \left( \frac{\partial C}{\partial m} \Big|_{m(x, \beta_0)} \right)^2 \sigma^2(x, \beta_0) \frac{\partial m(x, \beta_0)}{\partial \beta} \frac{\partial m(x, \beta_0)}{\partial \beta'} \right].\end{aligned}$$

In this problem  $m(x, \beta) = \exp(x'\beta)$  and  $\sigma^2(x, \beta) = \sigma^2 \exp(2x'\beta) + \exp(x'\beta)$ .

- (a) For the normal distribution  $C(m) = m$ , therefore  $\partial C/\partial m = 1$  and  $V_{NPML} = V_{NLLS}$ .
- (b) For the Poisson distribution  $C(m) = \log m$ , therefore  $\partial C/\partial m = 1/m$ ,

$$\begin{aligned}\mathcal{J} &= \mathbb{E}[\exp(-x'\beta)xx' \exp(2x'\beta)] = \exp(\frac{1}{2}\beta'\beta)(I + \beta\beta'), \\ \mathcal{I} &= \mathbb{E}[\exp(-2x'\beta)(\sigma^2 \exp(2x'\beta) + \exp(x'\beta))xx' \exp(2x'\beta)] \\ &= \exp(\frac{1}{2}\beta'\beta)(I + \beta\beta') + \sigma^2 \exp(2\beta'\beta)(I + 4\beta\beta').\end{aligned}$$

Finally,

$$V_{PPML} = (I + \beta\beta')^{-1} \left( \exp(-\frac{1}{2}\beta'\beta)(I + \beta\beta') + \sigma^2 \exp(\beta'\beta)(I + 4\beta\beta') \right) (I + \beta\beta')^{-1}.$$

(c) For the Gamma distribution  $C(m) = -\alpha/m$ , therefore  $\partial C/\partial m = \alpha/m^2$ ,

$$\begin{aligned} \mathcal{J} &= \mathbb{E}[\alpha \exp(-2x'\beta)xx' \exp(2x'\beta)] = \alpha I, \\ \mathcal{I} &= \alpha^2 \mathbb{E}[\exp(-4x'\beta)(\sigma^2 \exp(2x'\beta) + \exp(x'\beta))xx' \exp(2x'\beta)] \\ &= \alpha^2 \sigma^2 I + \alpha^2 \exp(\frac{1}{2}\beta'\beta)(I + \beta\beta'). \end{aligned}$$

Finally,

$$V_{GPML} = \sigma^2 I + \exp(\frac{1}{2}\beta'\beta)(I + \beta\beta').$$

Part 4. We have the following variances:

$$\begin{aligned} V_{NLLS} &= (I + 4\beta\beta')^{-1} \left( \exp(\frac{1}{2}\beta'\beta)(I + 9\beta\beta') + \sigma^2 \exp(4\beta'\beta)(I + 16\beta\beta') \right) (I + 4\beta\beta')^{-1}, \\ V_{WNLLS} &= \sigma^2 \left( I - \mathbb{E} \left[ \frac{xx'}{1 + \sigma^2 \exp(x'\beta)} \right] \right)^{-1}, \\ V_{NPML} &= V_{NLLS}, \\ V_{PPML} &= (I + \beta\beta')^{-1} \left( \exp(-\frac{1}{2}\beta'\beta)(I + \beta\beta') + \sigma^2 \exp(\beta'\beta)(I + 4\beta\beta') \right) (I + \beta\beta')^{-1}, \\ V_{GPML} &= \sigma^2 I + \exp(\frac{1}{2}\beta'\beta)(I + \beta\beta'). \end{aligned}$$

From the theory we know that  $V_{WNLLS} \leq V_{NLLS}$ . Next, we know that in the class of PML estimators the efficiency bound is achieved when  $\partial C/\partial m|_{m(x,\beta_0)}$  is proportional to  $\sigma^{-2}(x, \beta_0)$ , then the bound is

$$\mathbb{E} \left[ \frac{\partial m(x, \beta)}{\partial \beta} \frac{\partial m(x, \beta)}{\partial \beta'} \frac{1}{\mathbb{V}[y|x]} \right]$$

which is equal to  $V_{WNLLS}$  in our case. So, we have  $V_{WNLLS} \leq V_{PPML}$  and  $V_{WNLLS} \leq V_{GPML}$ . The comparison of other variances is not straightforward. Consider the one-dimensional case. Then we have

$$\begin{aligned} V_{NLLS} &= \frac{e^{\beta^2/2}(1 + 9\beta^2) + \sigma^2 e^{4\beta^2}(1 + 16\beta^2)}{(1 + 4\beta^2)^2}, \\ V_{WNLLS} &= \sigma^2 \left( 1 - \mathbb{E} \left[ \frac{x^2}{1 + \sigma^2 \exp(x\beta)} \right] \right)^{-1}, \\ V_{NPML} &= V_{NLLS}, \\ V_{PPML} &= \frac{e^{\beta^2/2}(1 + \beta^2) + \sigma^2 e^{\beta^2}(1 + 4\beta^2)}{(1 + \beta^2)^2}, \\ V_{GPML} &= \sigma^2 + e^{\beta^2/2}(1 + \beta^2). \end{aligned}$$

We can calculate these (except  $V_{WNLLS}$ ) for various parameter sets. For example, for  $\sigma^2 = 0.01$  and  $\beta^2 = 0.4$   $V_{NLLS} < V_{PPML} < V_{GPML}$ , for  $\sigma^2 = 0.01$  and  $\beta^2 = 0.1$   $V_{PPML} < V_{NLLS} < V_{GPML}$ , for  $\sigma^2 = 1$  and  $\beta^2 = 0.4$   $V_{GPML} < V_{PPML} < V_{NLLS}$ , for  $\sigma^2 = 0.5$  and  $\beta^2 = 0.4$   $V_{PPML} < V_{GPML} < V_{NLLS}$ . However, it appears impossible to make  $V_{NLLS} < V_{GPML} < V_{PPML}$  or  $V_{GPML} < V_{NLLS} < V_{PPML}$ .

## 15.11 Optimal instrument and regression on constant

Part 1. We have the following moment function:  $m(x, y, \theta) = (y - \alpha, (y - \alpha)^2 - \sigma^2 x^2)'$  with  $\theta = (\alpha, \sigma^2)'$ . The optimal unconditional moment restriction is  $\mathbb{E}[A^*(x)m(x, y, \theta)] = 0$ , where  $A^*(x) = D'(x)\Omega(x)^{-1}$ ,  $D(x) = \mathbb{E}[\partial m(x, y, \theta)/\partial \theta' | x]$ ,  $\Omega(x) = \mathbb{E}[m(x, y, \theta)m(x, y, \theta)' | x]$ .

(a) For the first moment restriction  $m_1(x, y, \theta) = y - \alpha$  we have  $D(x) = -1$  and  $\Omega(x) = \mathbb{E}[(y - \alpha)^2 | x] = \sigma^2 x^2$ , therefore the optimal moment restriction is

$$\mathbb{E}\left[\frac{y - \alpha}{x^2}\right] = 0.$$

(b) For the moment function  $m(x, y, \theta)$  we have

$$D(x) = \begin{pmatrix} -1 & 0 \\ 0 & -x^2 \end{pmatrix}, \quad \Omega(x) = \begin{pmatrix} \sigma^2 x^2 & 0 \\ 0 & \mu_4(x) - x^4 \sigma^4 \end{pmatrix},$$

where  $\mu_4(x) = \mathbb{E}[(y - \alpha)^4 | x]$ . The optimal weighting matrix is

$$A^*(x) = \begin{pmatrix} \frac{1}{\sigma^2 x^2} & 0 \\ 0 & \frac{x^2}{\mu_4(x) - x^4 \sigma^4} \end{pmatrix}.$$

The optimal moment restriction is

$$\mathbb{E}\left[\begin{pmatrix} \frac{y - \alpha}{x^2} \\ \frac{(y - \alpha)^2 - \sigma^2 x^2}{\mu_4(x) - x^4 \sigma^4} x^2 \end{pmatrix}\right] = 0.$$

Part 2. (a) The GMM estimator is the solution of

$$\frac{1}{n} \sum_i \frac{y_i - \hat{\alpha}}{x_i^2} = 0 \quad \Rightarrow \quad \hat{\alpha} = \frac{\sum_i y_i}{\sum_i \frac{1}{x_i^2}}.$$

The estimator for  $\sigma^2$  can be drawn from the sample analog of the condition  $\mathbb{E}[(y - \alpha)^2] = \sigma^2 \mathbb{E}[x^2]$ :

$$\hat{\sigma}^2 = \frac{\sum_i (y_i - \hat{\alpha})^2}{\sum_i x_i^2}.$$

(b) The GMM estimator is the solution of

$$\sum_i \begin{pmatrix} \frac{y_i - \hat{\alpha}}{x_i^2} \\ \frac{(y_i - \hat{\alpha})^2 - \hat{\sigma}^2 x_i^2}{\hat{\mu}_4(x_i) - x_i^4 \hat{\sigma}^4} x_i^2 \end{pmatrix} = 0.$$

We have the same estimator for  $\alpha$ :

$$\hat{\alpha} = \frac{\sum_i y_i}{\sum_i \frac{1}{x_i^2}},$$

$\hat{\sigma}^2$  is the solution of

$$\sum_i \frac{(y_i - \hat{\alpha})^2 - \hat{\sigma}^2 x_i^2}{\hat{\mu}_4(x_i) - x_i^4 \hat{\sigma}^4} x_i^2 = 0,$$

where  $\hat{\mu}_4(x)$  is non-parametric estimator for  $\mu_4(x) = \mathbb{E}[(y - \alpha)^4 | x]$ , for example, a nearest neighbor or a series estimator.

Part 3. The general formula for the variance of the optimal estimator is

$$V = (\mathbb{E}[D'(x)\Omega(x)^{-1}D(x)])^{-1}.$$

(a)  $V_{\hat{\alpha}} = \sigma^2 (\mathbb{E}[x^{-2}])^{-1}$ . Use standard asymptotic techniques to find

$$V_{\hat{\sigma}^2} = \frac{\mathbb{E}[(y - \alpha)^4]}{(\mathbb{E}[x^2])^2} - \sigma^4.$$

(b)

$$V_{(\hat{\alpha}, \hat{\sigma}^2)} = \left( \mathbb{E} \left[ \begin{pmatrix} \frac{1}{\sigma^2 x^2} & 0 \\ 0 & \frac{x^4}{\mu_4(x) - x^4 \sigma^4} \end{pmatrix} \right] \right)^{-1} = \begin{pmatrix} \sigma^2 (\mathbb{E}[x^{-2}])^{-1} & 0 \\ 0 & \left( \mathbb{E} \left[ \frac{x^4}{\mu_4(x) - x^4 \sigma^4} \right] \right)^{-1} \end{pmatrix}.$$

When we use the optimal instrument, our estimator is more efficient, therefore  $V_{\hat{\sigma}^2} > V_{\hat{\sigma}^2}$ .

Estimators of asymptotic variance can be found through sample analogs:

$$\hat{V}_{\hat{\alpha}} = \hat{\sigma}^2 \left( \frac{1}{n} \sum_i \frac{1}{x_i^2} \right)^{-1}, \quad V_{\hat{\sigma}^2} = n \frac{\sum_i (y_i - \hat{\alpha})^4}{(\sum_i x_i^2)^2} - \hat{\sigma}^4, \quad V_{\hat{\sigma}^2} = n \left( \sum_i \frac{x_i^4}{\hat{\mu}_4(x_i) - x_i^4 \hat{\sigma}^4} \right)^{-1}.$$

Part 4. The normal distribution PML2 estimator is the solution of the following problem:

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\sigma}^2 \end{pmatrix}_{PML2} = \arg \max_{\alpha, \sigma^2} \left\{ \text{const} - \frac{n}{2} \log \sigma^2 - \frac{1}{\sigma^2} \sum_i \frac{(y_i - \alpha)^2}{2x_i^2} \right\}.$$

Solving gives

$$\hat{\alpha}_{PML2} = \hat{\alpha} = \frac{\sum_i y_i}{\sum_i \frac{1}{x_i^2}}, \quad \hat{\sigma}_{PML2}^2 = \frac{1}{n} \sum_i \frac{(y_i - \hat{\alpha})^2}{x_i^2}$$

Since we have the same estimator for  $\alpha$ , we have the same variance  $V_{\hat{\alpha}} = \sigma^2 (\mathbb{E}[x_i^{-2}])^{-1}$ . It can be shown that

$$V_{\hat{\sigma}^2} = \mathbb{E} \left[ \frac{\mu_4(x)}{x^4} \right] - \sigma^4.$$



# 16. EMPIRICAL LIKELIHOOD

## 16.1 Common mean

1. We have the following moment function:  $m(x, y, \theta) = (x - \theta, y - \theta)'$ . The EL estimator is the solution of the following optimization problem.

$$\sum_i \log p_i \rightarrow \max_{p_i, \theta}$$

subject to

$$\sum_i p_i m(x_i, y_i, \theta) = 0, \quad \sum_i p_i = 1.$$

Let  $\lambda$  be a Lagrange multiplier for the restriction  $\sum_i p_i m(x_i, y_i, \theta) = 0$ , then the solution of the problem satisfies

$$\begin{aligned} p_i &= \frac{1}{n} \frac{1}{1 + \lambda' m(x_i, y_i, \theta)}, \\ 0 &= \frac{1}{n} \sum_i \frac{1}{1 + \lambda' m(x_i, y_i, \theta)} m(x_i, y_i, \theta), \\ 0 &= \frac{1}{n} \sum_i \frac{1}{1 + \lambda' m(x_i, y_i, \theta)} \left( \frac{\partial m(x_i, y_i, \theta)}{\partial \theta'} \right)' \lambda. \end{aligned}$$

In our case,  $\lambda = (\lambda_1, \lambda_2)$ , and the system is

$$\begin{aligned} p_i &= \frac{1}{1 + \lambda_1(x_i - \theta) + \lambda_2(y_i - \theta)}, \\ 0 &= \frac{1}{n} \sum_i \frac{1}{1 + \lambda_1(x_i - \theta) + \lambda_2(y_i - \theta)} \begin{pmatrix} x_i - \theta \\ y_i - \theta \end{pmatrix}, \\ 0 &= \frac{1}{n} \sum_i \frac{-\lambda_1 - \lambda_2}{1 + \lambda_1(x_i - \theta) + \lambda_2(y_i - \theta)}. \end{aligned}$$

The asymptotic distribution of the estimators is

$$\sqrt{n}(\hat{\theta}_{EL} - \theta) \xrightarrow{d} N(0, V), \quad \sqrt{n} \begin{pmatrix} \lambda_1 \\ \lambda_2 \end{pmatrix} \xrightarrow{d} \mathcal{N}(0, U),$$

where  $V = (Q'_{\partial m} Q_{mm}^{-1} Q_{\partial m})^{-1}$ ,  $U = Q_{mm}^{-1} - Q_{mm}^{-1} Q_{\partial m} V Q'_{\partial m} Q_{mm}^{-1}$ . In our case  $Q_{\partial m} = \begin{pmatrix} -1 \\ -1 \end{pmatrix}$

and  $Q_{mm} = \begin{pmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{pmatrix}$ , therefore

$$V = \frac{\sigma_x^2 \sigma_y^2 - \sigma_{xy}^2}{\sigma_y^2 + \sigma_x^2 - 2\sigma_{xy}}, \quad U = \frac{1}{\sigma_y^2 + \sigma_x^2 - 2\sigma_{xy}} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}.$$

Estimators for  $V$  and  $U$  based on consistent estimators for  $\sigma_x^2$ ,  $\sigma_y^2$  and  $\sigma_{xy}$  can be constructed from sample moments.

2. The last equation of the system gives  $\lambda_1 = -\lambda_2 = \lambda$ , so we have

$$p_i = \frac{1}{1 + \lambda(x_i - y_i)}, \quad 0 = \frac{1}{n} \sum_i \frac{1}{1 + \lambda(x_i - y_i)} \begin{pmatrix} x_i - \theta \\ y_i - \theta \end{pmatrix}.$$

The EL estimator is

$$\hat{\theta}_{EL} = \sum_i \frac{x_i}{1 + \lambda_{EL}(x_i - y_i)} = \sum_i \frac{y_i}{1 + \lambda_{EL}(x_i - y_i)},$$

where  $\lambda_{EL}$  is the solution of

$$\sum_i \frac{x_i - y_i}{1 + \lambda(x_i - y_i)} = 0.$$

Consider the linearized EL estimator. Linearization with respect to  $\lambda$  around 0 gives

$$p_i = 1 - \lambda(x_i - y_i), \quad 0 = \frac{1}{n} \sum_i (1 - \lambda(x_i - y_i)) \begin{pmatrix} x_i - \theta \\ y_i - \theta \end{pmatrix},$$

and helps to find an approximate but explicit solution

$$\tilde{\lambda}_{AEL} = \frac{\sum_i (x_i - y_i)}{\sum_i (x_i - y_i)^2}, \quad \tilde{\theta}_{AEL} = \frac{\sum_i (1 - \lambda(x_i - y_i))x_i}{\sum_i (1 - \lambda(x_i - y_i))} = \frac{\sum_i (1 - \lambda(x_i - y_i))y_i}{\sum_i (1 - \lambda(x_i - y_i))}.$$

Observe that  $\tilde{\lambda}_{AEL}$  is a normalized distance between the sample means of  $x$ 's and  $y$ 's,  $\tilde{\theta}_{AEL}$  is a weighted sample mean. The weights are such that the weighted mean of  $x$ 's equals the weighted mean of  $y$ 's. So, the moment restriction is satisfied in the sample. Moreover, the weight of observation  $i$  depends on the distance between  $x_i$  and  $y_i$  and on how the signs of  $x_i - y_i$  and  $\bar{x} - \bar{y}$  relate to each other. If they have the same sign, then such observation says against the hypothesis that the means are equal, thus the weight corresponding to this observation is relatively small. If they have opposite signs, such observation supports the hypothesis that means are equal, thus the weight corresponding to this observation is relatively large.

3. The technique is the same as in the EL problem. The Lagrangian is

$$L = - \sum_i p_i \log p_i + \mu \left( \sum_i p_i - 1 \right) + \lambda' \sum_i p_i m(x_i, y_i, \theta).$$

The first-order conditions are

$$-\frac{1}{n}(\log p_i + 1) + \mu + \lambda' m(x_i, y_i, \theta) = 0, \quad \lambda' \sum_i p_i \frac{\partial m(x_i, y_i, \theta)}{\partial \theta'} = 0.$$

The first equation together with the condition  $\sum_i p_i = 1$  gives

$$p_i = \frac{e^{\lambda' m(x_i, y_i, \theta)}}{\sum_i e^{\lambda' m(x_i, y_i, \theta)}}.$$

Also, we have

$$0 = \sum_i p_i m(x_i, y_i, \theta), \quad 0 = \sum_i p_i \left( \frac{\partial m(x_i, y_i, \theta)}{\partial \theta'} \right)' \lambda.$$

The system for  $\theta$  and  $\lambda$  that gives the ET estimator is

$$0 = \sum_i e^{\lambda' m(x_i, y_i, \theta)} m(x_i, y_i, \theta), \quad 0 = \sum_i e^{\lambda' m(x_i, y_i, \theta)} \left( \frac{\partial m(x_i, y_i, \theta)}{\partial \theta'} \right)' \lambda.$$

In our simple case, this system is

$$0 = \sum_i e^{\lambda_1(x_i - \theta) + \lambda_2(y_i - \theta)} \begin{pmatrix} x_i - \theta \\ y_i - \theta \end{pmatrix}, \quad 0 = \sum_i e^{\lambda_1(x_i - \theta) + \lambda_2(y_i - \theta)} (\lambda_1 + \lambda_2).$$

Here we have  $\lambda_1 = -\lambda_2 = \lambda$  again. The ET estimator is

$$\hat{\theta}_{ET} = \frac{\sum_i x_i e^{\lambda(x_i - y_i)}}{\sum_i e^{\lambda(x_i - y_i)}} = \frac{\sum_i y_i e^{\lambda(x_i - y_i)}}{\sum_i e^{\lambda(x_i - y_i)}},$$

where  $\lambda$  is the solution of

$$\sum_i (x_i - y_i) e^{\lambda(x_i - y_i)} = 0.$$

Note, that linearization of this system gives the same result as in EL case.

Since ET estimators are asymptotically equivalent to EL estimators (the proof of this fact is trivial: the first-order Taylor expansion of the ET system gives the same result as that of the EL system), there is no need to calculate the asymptotic variances, they are the same as in part 1.

## 16.2 Kullback–Leibler Information Criterion

1. Minimization of

$$\mathcal{KLIC}(e : \pi) = \mathbb{E}_e \left[ \log \frac{e}{\pi} \right] = \sum_i \frac{1}{n} \log \frac{1}{n\pi_i}$$

is equivalent to maximization of  $\sum_i \log \pi_i$  which gives the EL estimator.

2. Minimization of

$$\mathcal{KLIC}(\pi : e) = \mathbb{E}_\pi \left[ \log \frac{\pi}{e} \right] = \sum_i \pi_i \log \frac{\pi_i}{1/n}$$

gives the ET estimator.

3. The knowledge of probabilities  $p_i$  gives the following modification of EL problem:

$$\sum_i p_i \log \frac{p_i}{\pi_i} \rightarrow \min_{\pi_i, \theta} \quad \text{s.t.} \quad \sum \pi_i = 1, \quad \sum \pi_i m(z_i, \theta) = 0.$$

The solution of this problem satisfies the following system:

$$\begin{aligned} \pi_i &= \frac{p_i}{1 + \lambda' m(z_i, \theta)}, \\ 0 &= \sum_i \frac{p_i}{1 + \lambda' m(z_i, \theta)} m(z_i, \theta), \\ 0 &= \sum_i \frac{p_i}{1 + \lambda' m(z_i, \theta)} \left( \frac{\partial m(z_i, \theta)}{\partial \theta'} \right)' \lambda. \end{aligned}$$

The knowledge of probabilities  $p_i$  gives the following modification of ET problem

$$\sum_i \pi_i \log \frac{\pi_i}{p_i} \rightarrow \min_{\pi_i, \theta} \quad \text{s.t.} \quad \sum \pi_i = 1, \quad \sum \pi_i m(z_i, \theta) = 0.$$

The solution of this problem satisfies the following system

$$\begin{aligned}\pi_i &= \frac{p_i e^{\lambda' m(z_i, \theta)}}{\sum_j p_j e^{\lambda' m(z_j, \theta)}}, \\ 0 &= \sum_i p_i e^{\lambda' m(z_i, \theta)} m(z_i, \theta), \\ 0 &= \sum_i p_i e^{\lambda' m(z_i, \theta)} \left( \frac{\partial m(z_i, \theta)}{\partial \theta'} \right)' \lambda.\end{aligned}$$

4. The problem

$$\mathcal{KLIC}(e : f) = \mathbb{E}_e \left[ \log \frac{e}{f} \right] = \sum_i \frac{1}{n} \log \frac{1/n}{f(z_i, \theta)} \rightarrow \min_{\theta}$$

is equivalent to

$$\sum_i \log f(z_i, \theta) \rightarrow \max_{\theta},$$

which gives the Maximum Likelihood estimator.

## 16.3 Empirical likelihood as IV estimation

The efficient GMM estimator is

$$\beta_{GMM} = (\mathcal{Z}'_{GMM} \mathcal{X})^{-1} \mathcal{Z}'_{GMM} \mathcal{Y},$$

where  $\mathcal{Z}_{GMM}$  contains implied GMM instruments

$$\mathcal{Z}_{GMM} = \mathcal{X}' \mathcal{Z} \left( \mathcal{Z}' \hat{\Omega} \mathcal{Z} \right)^{-1} \mathcal{Z}',$$

and  $\hat{\Omega}$  contains squared residuals on the main diagonal.

The FOC to the EL problem are

$$\begin{aligned}0 &= \frac{1}{n} \sum_i \frac{z_i (y_i - x'_i \beta_{EL})}{1 + \lambda'_{EL} z_i (y_i - x'_i \beta_{EL})} \equiv \sum_i \pi_i z_i (y_i - x'_i \beta_{EL}), \\ 0 &= \frac{1}{n} \sum_i \frac{1}{1 + \lambda'_{EL} z_i (y_i - x'_i \beta_{EL})} (-z_i x'_i)' \lambda_{EL} \equiv - \sum_i \pi_i x_i z'_i \lambda_{EL}.\end{aligned}$$

From the first equation after premultiplication by  $\sum_i \pi_i x_i z'_i \left( \mathcal{Z}' \hat{\Omega} \mathcal{Z} \right)^{-1}$  it follows that

$$\beta_{EL} = (\mathcal{Z}'_{EL} \mathcal{X})^{-1} \mathcal{Z}'_{EL} \mathcal{Y},$$

where  $\mathcal{Z}_{EL}$  contains nonfeasible (because they depend on yet unknown parameters  $\beta_{EL}$  and  $\lambda_{EL}$ ) EL instruments

$$\mathcal{Z}_{EL} = \mathcal{X}' \hat{\Pi} \mathcal{Z} \left( \mathcal{Z}' \hat{\Omega} \mathcal{Z} \right)^{-1} \mathcal{Z}' \hat{\Pi},$$

where  $\hat{\Pi} \equiv \text{diag}(\pi_1, \dots, \pi_n)$ .

If we compare the expressions for  $\beta_{GMM}$  and  $\beta_{EL}$ , we see that in the construction of  $\beta_{EL}$  some expectations are estimated using EL probability weights rather than the empirical distribution. Using probability weights yields more efficient estimates, hence  $\beta_{EL}$  is expected to exhibit better finite sample properties than  $\beta_{GMM}$ .

# 17. ADVANCED ASYMPTOTIC THEORY

## 17.1 Maximum likelihood and asymptotic bias

(a) The ML estimator is  $\hat{\lambda} = \bar{y}_T^{-1}$  for which the second order expansion is

$$\begin{aligned}\hat{\lambda} &= \frac{1}{\mathbb{E}[y] \left(1 + (\mathbb{E}[y])^{-1} \frac{1}{T} \sum_{t=1}^T (y_t - \mathbb{E}[y])\right)} \\ &= \lambda \left(1 - \lambda \frac{1}{\sqrt{T}} \frac{1}{\sqrt{T}} \sum_{t=1}^T (y_t - \lambda^{-1}) + \lambda^2 \frac{1}{T} \left(\frac{1}{\sqrt{T}} \sum_{t=1}^T (y_t - \lambda^{-1})\right)^2 + o_p\left(\frac{1}{T}\right)\right).\end{aligned}$$

Therefore, the second order bias of  $\hat{\lambda}$  is

$$\mathbb{B}_2(\hat{\lambda}) = \frac{1}{T} \lambda^3 \mathbb{V}[y] = \frac{1}{T} \lambda.$$

(b) We can use the general formula for the second order bias of extremum estimators. For the ML,  $\Psi(y, \lambda) = \log f(y, \lambda) = \log \lambda - \lambda y$ , so

$$\mathbb{E} \left[ \left( \frac{\partial \Psi}{\partial \lambda} \right)^2 \right] = \lambda^{-2}, \quad \mathbb{E} \left[ \frac{\partial^2 \Psi}{\partial \lambda^2} \right] = -\lambda^{-2}, \quad \mathbb{E} \left[ \frac{\partial^3 \Psi}{\partial \lambda^3} \right] = 2\lambda^{-3}, \quad \mathbb{E} \left[ \frac{\partial^2 \Psi}{\partial \lambda^2} \frac{\partial \Psi}{\partial \lambda} \right] = 0,$$

so the second order bias of  $\hat{\lambda}$  is

$$\mathbb{B}_2(\hat{\lambda}) = \frac{1}{T} \left( (\lambda^{-2})^{-2} \cdot 0 + \frac{1}{2} (\lambda^{-2})^{-3} \cdot 2\lambda^{-3} \cdot \lambda^{-2} \right) = \frac{1}{T} \lambda.$$

The bias corrected ML estimator of  $\lambda$  is

$$\hat{\lambda}^* = \hat{\lambda} - \frac{1}{T} \hat{\lambda} = \frac{T-1}{T} \frac{1}{\bar{y}_T}.$$

## 17.2 Empirical likelihood and asymptotic bias

Solving the standard EL problem, we end up with the system in a most convenient form

$$\begin{aligned}\hat{\theta}_{EL} &= \frac{1}{n} \sum_{i=1}^n \frac{x_i}{1 + \hat{\lambda}(x_i - y_i)}, \\ 0 &= \frac{1}{n} \sum_{i=1}^n \frac{x_i - y_i}{1 + \hat{\lambda}(x_i - y_i)},\end{aligned}$$

where  $\hat{\lambda}$  is one of original Lagrange multiplies (the other equals  $-\hat{\lambda}$ ). From the first equation, the second order expansion for  $\hat{\theta}_{EL}$  is

$$\begin{aligned}\hat{\theta}_{EL} &= \frac{1}{n} \sum_{i=1}^n x_i (1 - \hat{\lambda}(x_i - y_i) + \hat{\lambda}^2(x_i - y_i)^2) + o_p\left(\frac{1}{n}\right) \\ &= \theta + \frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}} \sum_{i=1}^n (x_i - \theta) - \frac{1}{n} \sqrt{n} \hat{\lambda} \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i (x_i - y_i) \\ &\quad + \frac{1}{n} (\sqrt{n} \hat{\lambda})^2 \mathbb{E}[x(x-y)^2] + o_p\left(\frac{1}{n}\right).\end{aligned}$$

We need a first order expansion for  $\hat{\lambda}$ , which from the second equation is

$$\sqrt{n} \hat{\lambda} = \frac{1}{\sqrt{n}} \frac{\sum_{i=1}^n (x_i - y_i)}{n^{-1} \sum_{i=1}^n (x_i - y_i)^2} + o_p(1) = \frac{1}{\mathbb{E}[(x-y)^2]} \frac{1}{\sqrt{n}} \sum_{i=1}^n (x_i - y_i) + o_p(1).$$

Then, continuing,

$$\begin{aligned}\hat{\theta}_{EL} &= \theta + \frac{1}{\sqrt{n}} \frac{1}{\sqrt{n}} \sum_{i=1}^n (x_i - \theta) - \frac{1}{n} \left( \frac{1}{\mathbb{E}[(x-y)^2]} \frac{1}{\sqrt{n}} \sum_{i=1}^n (x_i - y_i) \right) \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i (x_i - y_i) \\ &\quad + \frac{1}{n} \left( \frac{1}{\mathbb{E}[(x-y)^2]} \frac{1}{\sqrt{n}} \sum_{i=1}^n (x_i - y_i) \right)^2 \mathbb{E}[x(x-y)^2] + o_p\left(\frac{1}{n}\right).\end{aligned}$$

The second order bias of  $\hat{\theta}_{EL}$  then is

$$\begin{aligned}\mathbb{B}_2(\hat{\theta}_{EL}) &= \frac{1}{n} \mathbb{E} \left[ -\frac{1}{\mathbb{E}[(x-y)^2]} \frac{1}{\sqrt{n}} \sum_{i=1}^n (x_i - y_i) \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i (x_i - y_i) \right. \\ &\quad \left. + \left( \frac{1}{\mathbb{E}[(x-y)^2]} \frac{1}{\sqrt{n}} \sum_{i=1}^n (x_i - y_i) \right)^2 \mathbb{E}[x(x-y)^2] \right] \\ &= \frac{1}{n} \left( -\frac{\mathbb{E}[x(x-y)^2]}{\mathbb{E}[(x-y)^2]} + \frac{\mathbb{E}[(x-y)^2] \mathbb{E}[x(x-y)^2]}{(\mathbb{E}[(x-y)^2])^2} \right) \\ &= 0.\end{aligned}$$

## 17.3 Asymptotically irrelevant instruments

1. The formula for the 2SLS estimator is

$$\hat{\beta}_{2SLS} = \beta + \frac{\sum_i x_i z_i' (\sum_i z_i z_i')^{-1} \sum_i z_i e_i}{\sum_i x_i z_i' (\sum_i z_i z_i')^{-1} \sum_i x_i z_i}.$$

According to the LLN,  $n^{-1} \sum_i z_i z_i' \xrightarrow{p} Q_{zz} = \mathbb{E}[zz']$ . According to the CLT,

$$\frac{1}{\sqrt{n}} \sum_i \begin{pmatrix} z_i x_i \\ z_i e_i \end{pmatrix} \xrightarrow{p} \begin{pmatrix} \xi \\ \zeta \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_x^2 & \rho \sigma \sigma_x \\ \rho \sigma \sigma_x & \sigma^2 \end{pmatrix} \otimes Q_{zz} \right),$$

where  $\sigma_x^2 = \mathbb{E}[x^2]$  and  $\sigma^2 = \mathbb{E}[e^2]$  (it is additionally assumed that there is convergence in probability). Hence, we have

$$\hat{\beta}_{2SLS} = \beta + \frac{(n^{-1/2} \sum_i x_i z_i') (n^{-1} \sum_i z_i z_i')^{-1} (n^{-1/2} \sum_i z_i e_i)}{(n^{-1/2} \sum_i x_i z_i') (n^{-1} \sum_i z_i z_i')^{-1} (n^{-1/2} \sum_i x_i z_i)} \xrightarrow{p} \beta + \frac{\xi' Q_{zz}^{-1} \zeta}{\xi' Q_{zz}^{-1} \xi}.$$

2. Under weak instruments,

$$\hat{\beta}_{2SLS} \xrightarrow{p} \beta + \frac{(Q_{zz}c + \psi_{zv})' Q_{zz}^{-1} \psi_{zu}}{(Q_{zz}c + \psi_{zv})' Q_{zz}^{-1} (Q_{zz}c + \psi_{zv})},$$

where  $c$  is a constant in the weak instrument assumption. If  $c = 0$ , this formula coincides with the previous one, with  $\psi_{zv} = \xi$  and  $\psi_{zu} = \zeta$ .

3. The expected value of the probability limit of the 2SLS estimator is

$$\begin{aligned} \mathbb{E} \left[ p \lim \hat{\beta}_{2SLS} \right] &= \beta + \mathbb{E} \left[ \frac{\xi' Q_{zz}^{-1} \zeta}{\xi' Q_{zz}^{-1} \xi} \right] = \beta + \rho \frac{\sigma}{\sigma_x} \mathbb{E} \left[ \frac{\xi' Q_{zz}^{-1} \xi}{\xi' Q_{zz}^{-1} \xi} \right] \\ &= \beta + \frac{\rho \sigma \sigma_x}{\sigma_x^2} = p \lim \hat{\beta}_{OLS}, \end{aligned}$$

where we use joint normality to deduce that  $\mathbb{E}[\zeta|\xi] = (\rho\sigma/\sigma_x)\xi$ .

## 17.4 Weakly endogenous regressors

The OLS estimator satisfies

$$\hat{\beta} - \beta = (n^{-1} \mathcal{X}' \mathcal{X})^{-1} n^{-1} \mathcal{X}' \mathcal{E} = \frac{c}{\sqrt{n}} + (n^{-1} \mathcal{X}' \mathcal{X})^{-1} n^{-1} \mathcal{X}' \mathcal{U} \xrightarrow{p} 0,$$

and

$$\sqrt{n} (\hat{\beta} - \beta) = c + (n^{-1} \mathcal{X}' \mathcal{X})^{-1} n^{-1/2} \mathcal{X}' \mathcal{U} \xrightarrow{p} c + Q^{-1} \xi \sim N(c, \sigma_u^2 Q^{-1}).$$

The Wald test statistic satisfies

$$\begin{aligned} \mathcal{W} &= \frac{n (R\hat{\beta} - r)' (R(\mathcal{X}' \mathcal{X})^{-1} R')^{-1} (R\hat{\beta} - r)}{(y - \mathcal{X}\hat{\beta})' (y - \mathcal{X}\hat{\beta})} \\ &\xrightarrow{p} \frac{(c + Q^{-1}\xi)' R' (RQ^{-1}R')^{-1} R(c + Q^{-1}\xi)}{\sigma_u^2} \sim \chi_q^2(\delta), \end{aligned}$$

where

$$\delta = \frac{c' R' (RQ^{-1}R')^{-1} R c}{\sigma_u^2}$$

is the noncentrality parameter.

## 17.5 Weakly invalid instruments

1. Consider the projection of  $e$  on  $z$ :

$$e = z'\omega + v,$$

where  $v$  is orthogonal to  $z$ . Thus

$$n^{-1/2}Z'\mathcal{E} = n^{-1/2}Z'(Z\omega + \mathcal{V}) = (n^{-1}Z'Z)c_\omega + n^{-1/2}Z'\mathcal{V}.$$

Let us assume that

$$(n^{-1}Z'Z, n^{-1}Z'\mathcal{X}, n^{-1/2}Z'\mathcal{V}) \xrightarrow{p} (Q_{zz}, Q_{zx}, \xi),$$

where  $Q_{zz} \equiv \mathbb{E}[zz']$ ,  $Q_{zx} \equiv \mathbb{E}[zx']$  and  $\xi \sim \mathcal{N}(0, \sigma_v^2 Q_{zz})$ . The 2SLS estimator satisfies

$$\begin{aligned} \sqrt{n}(\hat{\beta} - \beta) &= \frac{n^{-1}\mathcal{X}'Z(n^{-1}Z'Z)^{-1}n^{-1/2}Z'\mathcal{E}}{n^{-1}\mathcal{X}'Z(n^{-1}Z'Z)^{-1}n^{-1}Z'\mathcal{X}} \\ &\xrightarrow{p} \frac{Q'_{zx}(c_\omega + Q_{zz}^{-1}\xi)}{Q'_{zx}Q_{zz}^{-1}Q_{zx}} \sim \mathcal{N}\left(\frac{Q'_{zx}c_\omega}{Q'_{zx}Q_{zz}^{-1}Q_{zx}}, \frac{\sigma_v^2}{Q'_{zx}Q_{zz}^{-1}Q_{zx}}\right). \end{aligned}$$

Since  $\hat{\beta}$  is consistent,

$$\begin{aligned} \hat{\sigma}_e^2 &\equiv n^{-1}\hat{\mathcal{U}}'\hat{\mathcal{U}} = n^{-1}(\mathcal{Y} - \mathcal{X}\hat{\beta})'(\mathcal{Y} - \mathcal{X}\hat{\beta}) - 2(\hat{\beta} - \beta)n^{-1}\mathcal{X}'(\mathcal{Y} - \mathcal{X}\hat{\beta}) + (\hat{\beta} - \beta)^2 n^{-1}\mathcal{X}'\mathcal{X} \\ &= n^{-1}(Z\omega + \mathcal{V})'(Z\omega + \mathcal{V}) - 2(\hat{\beta} - \beta)(n^{-1}\mathcal{X}'Z\omega + n^{-1}\mathcal{X}'\mathcal{V}) + (\hat{\beta} - \beta)^2 n^{-1}\mathcal{X}'\mathcal{X} \\ &\xrightarrow{p} \sigma_v^2. \end{aligned}$$

Thus the  $t$  ratio satisfies

$$t_\beta \equiv \frac{\hat{\beta} - \beta}{\sqrt{\hat{\sigma}_e^2 (\mathcal{X}'Z(Z'Z)^{-1}Z'\mathcal{X})^{-1}}} \xrightarrow{p} \frac{\frac{Q'_{zx}(c_\omega + Q_{zz}^{-1}\xi)}{Q'_{zx}Q_{zz}^{-1}Q_{zx}}}{\sqrt{\sigma_v^2 (Q'_{zx}Q_{zz}^{-1}Q_{zx})^{-1}}} \sim \mathcal{N}(\delta_t, 1),$$

where

$$\delta_t = \frac{Q'_{zx}c_\omega}{\sigma_v \sqrt{Q'_{zx}Q_{zz}^{-1}Q_{zx}}}$$

is the noncentrality parameter. Finally note that

$$\begin{aligned} n^{-1/2}Z'\hat{\mathcal{U}} &= n^{-1/2}Z'\mathcal{E} - \sqrt{n}(\hat{\beta} - \beta)n^{-1}Z'\mathcal{X} \\ &\xrightarrow{p} Q_{zz}(c_\omega + Q_{zz}^{-1}\xi) - \frac{Q'_{zx}(c_\omega + Q_{zz}^{-1}\xi)}{Q'_{zx}Q_{zz}^{-1}Q_{zx}}Q_{zx} \\ &\sim \mathcal{N}\left(Q_{zz}c_\omega - \frac{Q'_{zx}c_\omega Q_{zx}}{Q'_{zx}Q_{zz}^{-1}Q_{zx}}, \sigma_v^2 \left(Q_{zz} - \frac{Q_{zx}Q'_{zx}}{Q'_{zx}Q_{zz}^{-1}Q_{zx}}\right)\right), \end{aligned}$$

so,

$$\mathcal{J} = \frac{n^{-1/2}\hat{\mathcal{U}}'Z(n^{-1}Z'Z)^{-1}n^{-1/2}Z'\hat{\mathcal{U}}}{\hat{\sigma}_e^2} \xrightarrow{d} \chi_{\ell-1}^2(\delta_J),$$



where

$$\begin{aligned}\delta_J &= \left( Q_{zz}c_\omega - \frac{Q'_{zx}c_\omega Q_{zx}}{Q'_{zx}Q_{zz}^{-1}Q_{zx}} \right)' \frac{Q_{zz}^{-1}}{\sigma_v^2} \left( Q_{zz}c_\omega - \frac{Q'_{zx}c_\omega Q_{zx}}{Q'_{zx}Q_{zz}^{-1}Q_{zx}} \right) \\ &= \frac{1}{\sigma_v^2} c'_\omega \left( Q_{zz} - \frac{Q_{zx}Q'_{zx}}{Q'_{zx}Q_{zz}^{-1}Q_{zx}} \right) c_\omega\end{aligned}$$

is the noncentrality parameter. In particular, when  $\ell = 1$ ,

$$\begin{aligned}\sqrt{n}(\hat{\beta} - \beta) &\xrightarrow{p} \frac{c_\omega + Q_{zz}^{-1}\xi}{Q_{zz}^{-1}Q_{zx}} \sim \mathcal{N}\left(\frac{Q_{zz}}{Q_{zx}}c_\omega, \frac{Q_{zz}}{Q_{zx}^2}\sigma_v^2\right), \\ t_\beta^2 &\xrightarrow{p} \frac{(c_\omega + Q_{zz}^{-1}\xi)^2}{\sigma_u^2} Q_{zz} \sim \chi_1^2\left(Q_{zz} \frac{c_\omega^2}{\sigma_v^2}\right), \\ n^{-1/2}\mathcal{Z}'\hat{u} &\xrightarrow{p} 0 \quad \Rightarrow \quad \mathcal{J} \xrightarrow{p} 0.\end{aligned}$$

2. Consider also the projection of  $x$  on  $z$ :

$$x = z'\pi + w,$$

where  $w$  is orthogonal to  $z$ . Thus

$$n^{-1/2}\mathcal{Z}'\mathcal{X} = n^{-1/2}\mathcal{Z}'(\mathcal{Z}\pi + \mathcal{W}) = (n^{-1}\mathcal{Z}'\mathcal{Z})c_\pi + n^{-1/2}\mathcal{Z}'\mathcal{W}.$$

We have

$$\left( n^{-1}\mathcal{Z}'\mathcal{Z}, n^{-1/2}\mathcal{Z}'\mathcal{W}, n^{-1/2}\mathcal{Z}'\mathcal{V} \right) \xrightarrow{p} (Q_{zz}, \zeta, \xi),$$

where  $\begin{pmatrix} \zeta \\ \xi \end{pmatrix} \sim \mathcal{N}\left(0, \begin{pmatrix} \sigma_w^2 & \sigma_{wv} \\ \sigma_{wv} & \sigma_v^2 \end{pmatrix} \otimes Q_{zz}\right)$ . The 2SLS estimator satisfies

$$\begin{aligned}\hat{\beta} - \beta &= \frac{n^{-1/2}\mathcal{X}'\mathcal{Z}(n^{-1}\mathcal{Z}'\mathcal{Z})^{-1}n^{-1/2}\mathcal{Z}'\mathcal{E}}{n^{-1/2}\mathcal{X}'\mathcal{Z}(n^{-1}\mathcal{Z}'\mathcal{Z})^{-1}n^{-1/2}\mathcal{Z}'\mathcal{X}} \\ &\xrightarrow{p} \frac{(Q_{zz}c_\pi + \zeta)' Q_{zz}^{-1}(Q_{zz}c_\omega + \xi)}{(Q_{zz}c_\pi + \zeta)' Q_{zz}^{-1}(Q_{zz}c_\pi + \zeta)} \equiv \frac{\sigma_v(\lambda_\pi + z_w)'(\lambda_\omega + z_v)}{\sigma_w(\lambda_\pi + z_w)'(\lambda_\pi + z_w)} \equiv \frac{\sigma_v \nu_2}{\sigma_w \nu_1},\end{aligned}$$

where  $\begin{pmatrix} z_w \\ z_v \end{pmatrix} \sim \mathcal{N}\left(0, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \otimes I_\ell\right)$  and  $\rho \equiv \frac{\sigma_{wv}}{\sigma_w\sigma_v}$ . Note that  $\hat{\beta}$  is inconsistent. Thus,

$$\begin{aligned}\hat{\sigma}_e^2 &\equiv n^{-1}\hat{u}'\hat{u} = n^{-1}(\mathcal{Y} - \mathcal{X}\hat{\beta})'(\mathcal{Y} - \mathcal{X}\hat{\beta}) - 2(\hat{\beta} - \beta)n^{-1}\mathcal{X}'(\mathcal{Y} - \mathcal{X}\hat{\beta}) + (\hat{\beta} - \beta)^2 n^{-1}\mathcal{X}'\mathcal{X} \\ &= n^{-1}(\mathcal{Z}\omega + \mathcal{V})'(\mathcal{Z}\omega + \mathcal{V}) - 2(\hat{\beta} - \beta)n^{-1}(\mathcal{Z}\pi + \mathcal{W})'(\mathcal{Z}\omega + \mathcal{V}) + (\hat{\beta} - \beta)^2 n^{-1}\mathcal{X}'\mathcal{X} \\ &\xrightarrow{p} \sigma_v^2 - 2\left(\frac{\sigma_v \nu_2}{\sigma_w \nu_1}\right)\sigma_{wv} + \left(\frac{\sigma_v \nu_2}{\sigma_w \nu_1}\right)^2 \sigma_w^2 = \sigma_v^2 \left(1 - 2\rho\frac{\nu_2}{\nu_1} + \left(\frac{\nu_2}{\nu_1}\right)^2\right).\end{aligned}$$

Thus the  $t$  ratio satisfies

$$t_\beta \equiv \frac{\hat{\beta} - \beta}{\sqrt{\hat{\sigma}_e^2 (\mathcal{X}'\mathcal{Z}(\mathcal{Z}'\mathcal{Z})^{-1}\mathcal{Z}'\mathcal{X})^{-1}}} \xrightarrow{p} \frac{\nu_2/\sqrt{\nu_1}}{\sqrt{1 - 2\rho\nu_2/\nu_1 + (\nu_2/\nu_1)^2}}.$$

Finally note that

$$\begin{aligned}
n^{-1/2} \mathbf{Z}' \widehat{\mathbf{U}} &= n^{-1/2} \mathbf{Z}' \boldsymbol{\varepsilon} - (\widehat{\beta} - \beta) n^{-1/2} \mathbf{Z}' \mathcal{X} \\
&\xrightarrow{p} Q_{zz}^{1/2} \sigma_v (\lambda_\omega + z_v) - \frac{\sigma_v \nu_2}{\sigma_w \nu_1} Q_{zz}^{1/2} \sigma_w (\lambda_\pi + z_w) \\
&= \sigma_v Q_{zz}^{1/2} \left( (\lambda_\omega + z_v) - \frac{\nu_2}{\nu_1} (\lambda_\pi + z_w) \right) \equiv \psi,
\end{aligned}$$

so,

$$\begin{aligned}
\mathcal{J} &= \frac{n^{-1/2} \widehat{\mathbf{U}}' \mathbf{Z} (n^{-1} \mathbf{Z}' \mathbf{Z})^{-1} n^{-1/2} \mathbf{Z}' \widehat{\mathbf{U}}}{\widehat{\sigma}_e^2} \\
&\xrightarrow{p} \frac{\left( (\lambda_\omega + z_v) - \frac{\nu_2}{\nu_1} (\lambda_\pi + z_w) \right)' \left( (\lambda_\omega + z_v) - \frac{\nu_2}{\nu_1} (\lambda_\pi + z_w) \right)}{1 - 2\rho \frac{\nu_2}{\nu_1} + \left( \frac{\nu_2}{\nu_1} \right)^2} \\
&= \frac{\nu_3 - \frac{\nu_2^2}{\nu_1}}{1 - 2\rho \frac{\nu_2}{\nu_1} + \left( \frac{\nu_2}{\nu_1} \right)^2},
\end{aligned}$$

where  $\nu_3 \equiv (\lambda_\omega + z_v)' (\lambda_\omega + z_v)$ . In particular, when  $\ell = 1$ ,

$$\begin{aligned}
\widehat{\beta} - \beta &\xrightarrow{p} \frac{\sigma_v \lambda_\omega + z_v}{\sigma_w \lambda_\pi + z_w}, \\
\widehat{\sigma}_e^2 &\xrightarrow{p} \sigma_v^2 \left( 1 - 2\rho \frac{\lambda_\omega + z_v}{\lambda_\pi + z_w} + \left( \frac{\lambda_\omega + z_v}{\lambda_\pi + z_w} \right)^2 \right), \\
t_\beta &\xrightarrow{p} \frac{\frac{\lambda_\omega + z_v}{\lambda_\pi + z_w}}{\sqrt{1 - 2\rho \frac{\lambda_\omega + z_v}{\lambda_\pi + z_w} + \left( \frac{\lambda_\omega + z_v}{\lambda_\pi + z_w} \right)^2}}, \\
\psi = 0 &\Rightarrow \mathcal{J} \xrightarrow{p} 0.
\end{aligned}$$