

# Bavarian Graduate Program in Economics

Prof. Dr. Martin Kukuk

Universität Würzburg

## Frontiers in Econometrics

Das vorliegende Manuskript ist für Studierende, die den Graduiertenkurs *Frontiers in Econometrics* besuchen, als Begleitmaterial gedacht, um den Mitschrieb zu erleichtern. Ich kann nicht für die Vollständigkeit und die Fehlerfreiheit des Manuskripts garantieren. Allein die Ausführungen in den Vorlesungen sind relevant.

Juli 2007

# Contents

<b>1</b>	<b>The Geometry of Least Squares</b>	<b>3</b>
1.1	Vector spaces . . . . .	3
1.2	Geometry of OLS estimation . . . . .	5
1.3	Frisch-Waugh-Lovell Theorem . . . . .	8
1.4	Goodness of Fit . . . . .	10
1.5	Influential Points and Leverage . . . . .	10
<b>2</b>	<b>Statistical Properties of OLS</b>	<b>12</b>
2.1	Variance-Covariance matrix of $\hat{\beta}$ . . . . .	13
2.2	Gauss-Markov-Theorem . . . . .	15
2.3	Some statistical results . . . . .	16
2.4	Estimation of $\sigma^2$ . . . . .	18
2.5	Hypothesis testing . . . . .	18
<b>3</b>	<b>Asymptotic Theory</b>	<b>20</b>
3.1	Probability Limit . . . . .	20
3.2	Consistency of OLS estimator . . . . .	21
3.3	Asymptotic normality . . . . .	23
<b>4</b>	<b>Generalized Method of Moments</b>	<b>26</b>
4.1	Method-of-Moments estimation . . . . .	26
4.2	Errors-in-Variables . . . . .	26
4.3	Generalized Method of Moments . . . . .	29

<b>5</b>	<b>Maximum Likelihood</b>	<b>31</b>
5.1	Hypothesis tests . . . . .	35
<b>6</b>	<b>Generalized Least Squares</b>	<b>37</b>
6.1	GLS estimator . . . . .	37
6.2	Feasible GLS . . . . .	39
6.3	Heteroscedasticity . . . . .	40
6.3.1	Testing for Heteroscedasticity . . . . .	41
6.3.2	Heteroscedasticity-Consistent Covariance Matrix (HCCM) . . . . .	42
6.4	Autoregressive Processes . . . . .	43
6.4.1	AR(1) Process . . . . .	43
6.4.2	Higher-Order Autoregressive Processes . . . . .	44
6.4.3	Testing for Serial Correlation . . . . .	45
6.4.4	Estimation . . . . .	45
6.5	Panel Data . . . . .	47
6.5.1	Instrumental Variables Estimator . . . . .	50
6.5.2	Dynamic Linear Models . . . . .	50
6.6	Systems of Regression Equations . . . . .	52
6.6.1	Seemingly Unrelated Regressions . . . . .	52
6.6.2	Linear Simultaneous Equation Models . . . . .	55
<b>7</b>	<b>Discrete and Limited Dependent Variables</b>	<b>56</b>
7.1	Binary Response Models . . . . .	56
7.2	Ordered Discrete Response Models . . . . .	58

	4
7.3	Unordered Discrete Response Models . . . . . 59
7.4	Count Data Models . . . . . 61
7.5	Truncated and Censored Models . . . . . 62
7.5.1	Truncated Regression Models . . . . . 62
7.5.2	Censored Regression Models . . . . . 63
7.6	Sample Selectivity Models . . . . . 64
7.7	Duration Models . . . . . 65
<b>8</b>	<b>Time Series Analysis . . . . . 69</b>
8.1	Stationary Univariate Processes . . . . . 69
8.2	Stationary Multivariate Processes, VAR . . . . . 71
8.3	Nonstationary Time Series . . . . . 72
8.3.1	Deterministic and Stochastic Trends . . . . . 72
8.3.2	Estimation of AR(1) Models . . . . . 74
8.3.3	Unit Roots Tests . . . . . 75
8.4	Cointegration . . . . . 76
8.4.1	Error Correction model . . . . . 76
8.4.2	Spurious Regression . . . . . 77
8.4.3	Cointegration of two Variables . . . . . 78
8.4.4	Cointegration in a VAR . . . . . 79

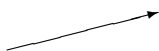
# 1 The Geometry of Least Squares

## 1.1 Vector spaces

Representation of a vector

$$\mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix} \in E^n$$

as



### Operations

- Addition  $a + b = c$
- Scalar product  $\alpha \cdot \mathbf{a}$  with  $\alpha \in \mathbb{R}$
- Inner product  $\mathbf{a}'\mathbf{b} = \mathbf{b}'\mathbf{a} = \sum_{i=1}^n a_i \cdot b_i$
- Length or norm

$$\|\mathbf{x}\| = (\mathbf{x}'\mathbf{x})^{1/2} = \left(\sum x_i^2\right)^{1/2}$$

This is inspired by Pythagoras' theorem:  $a^2 + b^2 = c^2$



Vectors in  $E^2$  using cartesian coordinates:



Geometric interpretation of innerproduct



General case:



$$\left( \frac{1}{\|\mathbf{x}\|} \mathbf{x} \right)' \left( \frac{1}{\|\mathbf{y}\|} \mathbf{y} \right) = \cos \theta \quad \Leftrightarrow \quad \mathbf{x}' \mathbf{y} = \|\mathbf{x}\| \cdot \|\mathbf{y}\| \cdot \cos \theta$$

Orthogonality:  $\mathbf{x}' \mathbf{y} = 0 = \cos(\pi/2)$

## Subspaces

Subspace  $S(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k)$  is a set of all **linear combinations**:

$$S(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k) = \left\{ \mathbf{z} \in E^n \mid \mathbf{z} = \sum_{i=1}^k b_i \cdot \mathbf{x}_i, b_i \in \mathbb{R}, \mathbf{x}_i \in E^n \right\}$$

Denoting  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k]$  and  $\mathbf{b} = (b_1, \dots, b_k)'$ , then

$$S(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k) = S(\mathbf{X}) = \left\{ \mathbf{z} \in E^n \mid \mathbf{z} = \mathbf{X} \mathbf{b}, \mathbf{b} \in E^k, \mathbf{x}_i \in E^n \right\}$$

## Linear Independence

Columns of  $\mathbf{X}$  are linear independent if  $\mathbf{X} \mathbf{b} = \mathbf{0}$  is only possible with  $\mathbf{b} = \mathbf{0}$ .

Linear dependent columns of  $\mathbf{X}$ :

$$\begin{aligned} \mathbf{X} \mathbf{b} &= \mathbf{0} \quad \text{with } \mathbf{b} \neq \mathbf{0} \\ \Leftrightarrow \mathbf{X}' \mathbf{X} \mathbf{b} &= \mathbf{X}' \mathbf{0} = \mathbf{0} \\ \Leftrightarrow (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{X} \mathbf{b} &= \mathbf{b} = \mathbf{0} \end{aligned}$$

This is a contradiction. Thus if columns of  $\mathbf{X}$  are linear dependent then  $\mathbf{X}'\mathbf{X}$  is not invertible.

## 1.2 Geometry of OLS estimation

Considering the model:

$$y_t = \underset{(1 \times k)}{\mathbf{X}_t} \cdot \underset{(k \times 1)}{\boldsymbol{\beta}} + u_t \quad ; \quad t = 1, \dots, n$$

In matrix notation

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$$

$\boldsymbol{\beta}$  is unknown,  $\mathbf{y}$  and  $\mathbf{X}$  are observed,  $\mathbf{u} = \mathbf{u}(\boldsymbol{\beta})$  unobservable.

The estimation principle of Least Squares:

$$\min_{\boldsymbol{\beta}} \sum_{t=1}^n u_t^2 = \min_{\boldsymbol{\beta}} \mathbf{u}(\boldsymbol{\beta})' \mathbf{u}(\boldsymbol{\beta}) = \min_{\boldsymbol{\beta}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

In words: Find  $\hat{\mathbf{u}} = \mathbf{u}(\hat{\boldsymbol{\beta}})$  that has least squared length and therefore least length  $\|\mathbf{u}(\hat{\boldsymbol{\beta}})\|$ .

$\mathbf{u}(\hat{\boldsymbol{\beta}})$  with least length is found if we drop a perpendicular on  $S(\mathbf{X})$ .

We see that  $\hat{\mathbf{u}}$  is orthogonal on  $\mathbf{X}\hat{\boldsymbol{\beta}}$  or more general orthogonal on every vector in  $S(\mathbf{X})$ .

Geometrically, least squares estimation can be decomposed in two steps:

- Find point in  $S(\mathbf{X})$  closest to  $\mathbf{y}$  :  $\mathbf{X}\hat{\boldsymbol{\beta}}$
- Knowing  $\mathbf{X}\hat{\boldsymbol{\beta}}$ , determine  $\hat{\boldsymbol{\beta}}$

Applying Pythagoras' theorem we obtain

$$\begin{aligned}\|\mathbf{y}\| &= \|\mathbf{X}\hat{\boldsymbol{\beta}}\| + \|\hat{\mathbf{u}}\| \\ \Leftrightarrow \mathbf{y}'\mathbf{y} &= \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\mathbf{u}}'\hat{\mathbf{u}}\end{aligned}$$

The total sum of squares is equal to the explained sum of squares plus the sum of squared residuals.

### OLS estimator

Vector  $\hat{\mathbf{u}}$  is orthogonal to all vectors in  $S(\mathbf{X})$  and especially on all the elements of  $\mathbf{X}$ :

$$\begin{aligned}\mathbf{X}'\hat{\mathbf{u}} &= \underset{(k \times 1)}{\mathbf{0}} \quad \Leftrightarrow \quad \mathbf{X}'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{0} \\ \Leftrightarrow \mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} &= \mathbf{0} \quad \Leftrightarrow \quad \mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y}\end{aligned}$$



If  $\mathbf{X}$  has linear independent columns then  $\mathbf{X}'\mathbf{X}$  is invertible and we can premultiply the last equation with  $(\mathbf{X}'\mathbf{X})^{-1}$ :

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

If  $\mathbf{X}$  has linear dependent columns we still can project on to  $S(\mathbf{X})$  but solving for  $\hat{\boldsymbol{\beta}}$  is not unique.

### Orthogonal projection

First step of OLS is a mapping of  $\mathbf{y}$  on to  $S(\mathbf{X})$ . We call it a projection of  $\mathbf{y}$  into a point of subspace  $S(\mathbf{X})$ . This mapping is  $\mathbf{y} \longrightarrow \mathbf{X}\hat{\boldsymbol{\beta}}$  or

$$\underset{(n \times n)}{\mathbf{P}_X} \cdot \underset{(n \times 1)}{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$$

Inserting the OLS estimator  $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ :

$$\mathbf{P}_X \cdot \mathbf{y} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

Thus we obtain the orthogonal projector:

$$\mathbf{P}_X = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

Since every point (or vector) in subspace  $S(\mathbf{X})$  is projected on to itself, we have

$$\mathbf{P}_X \cdot \mathbf{X}\boldsymbol{\beta} = \mathbf{X}\boldsymbol{\beta} \quad \text{or specifically } \mathbf{P}_X \cdot \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}\hat{\boldsymbol{\beta}}$$

Thus we get

$$\mathbf{P}_X \mathbf{P}_X \mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{P}_X \mathbf{y}$$

or  $\mathbf{P}_X \cdot \mathbf{P}_X = \mathbf{P}_X$ . Analytically:

$$\mathbf{P}_X \cdot \mathbf{P}_X = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \cdot \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' = \mathbf{P}_X$$

A matrix having this property is called **idempotent**.

Projecting  $\mathbf{y}$  on to  $S(\mathbf{X})$  yields the fitted values  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ . Analogously, projecting on to

$$S^\perp(\mathbf{X}) = \left\{ \mathbf{w} \in E^n \mid \mathbf{w}'\mathbf{z} = 0 \quad \forall \mathbf{z} \in S(\mathbf{X}) \right\}$$

yields the residuals  $\hat{\mathbf{u}}$  :

$$\begin{aligned} \mathbf{M}_X &= \mathbf{I} - \mathbf{P}_X = \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \\ \Rightarrow \mathbf{M}_X \cdot \mathbf{y} &= (\mathbf{I} - \mathbf{P}_X) \cdot \mathbf{y} = \mathbf{y} - \mathbf{P}_X \cdot \mathbf{y} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \hat{\mathbf{u}} \end{aligned}$$

$\mathbf{M}_X$  annihilates each vector of  $S(\mathbf{X})$ :

$$\mathbf{M}_X \mathbf{X}\boldsymbol{\beta} = (\mathbf{I} - \mathbf{P}_X)\mathbf{X}\boldsymbol{\beta} = \mathbf{X}\boldsymbol{\beta} - \mathbf{X}\boldsymbol{\beta} = \mathbf{0}$$

$\mathbf{P}_X$  annihilates each vector of  $S^\perp(\mathbf{X})$ :

$$\mathbf{P}_X \cdot \mathbf{w} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \underbrace{\mathbf{X}'\mathbf{w}}_{=0}$$

$\mathbf{P}_X$  and  $\mathbf{M}_X$  are called orthogonal matrices:

$$\mathbf{M}_X \cdot \mathbf{P}_X = (\mathbf{I} - \mathbf{P}_X)\mathbf{P}_X = \mathbf{P}_X - \mathbf{P}_X\mathbf{P}_X = \mathbf{P}_X - \mathbf{P}_X = \mathbf{0}$$

### 1.3 Frisch-Waugh-Lovell Theorem

We have seen that  $\mathbf{P}_X + \mathbf{M}_X = \mathbf{I}$  and thus

$$\mathbf{y} = \mathbf{P}_X \mathbf{y} + \mathbf{M}_X \mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\mathbf{u}}$$

Partitioning regressor matrix  $\mathbf{X}$  and also vector  $\boldsymbol{\beta}$

$$\mathbf{X}_{(n \times k)} = \begin{bmatrix} \mathbf{X}_1 & \vdots & \mathbf{X}_2 \\ (n \times k_1) & & (n \times k_2) \end{bmatrix} \quad \boldsymbol{\beta} = \begin{pmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{pmatrix}$$

we can write our model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$  as

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u} \quad (1)$$

We now define  $\mathbf{P}_{\mathbf{X}_1} = \mathbf{X}_1(\mathbf{X}'_1\mathbf{X}_1)^{-1}\mathbf{X}'_1$  and  $\mathbf{M}_{\mathbf{X}_1} = \mathbf{I} - \mathbf{P}_{\mathbf{X}_1}$  and analogously  $\mathbf{P}_{\mathbf{X}_2}, \mathbf{M}_{\mathbf{X}_2}$ .

### FWL Theorem

1. OLS estimates for  $\boldsymbol{\beta}_2$  from regression (1) and from  $\mathbf{M}_{\mathbf{X}_1}\mathbf{y} = \mathbf{M}_{\mathbf{X}_1}\mathbf{X}_2\boldsymbol{\beta}_2 + \text{residual}$  are numerically identical.
2. The residuals from from regression (1) and from  $\mathbf{M}_{\mathbf{X}_1}\mathbf{y} = \mathbf{M}_{\mathbf{X}_1}\mathbf{X}_2\boldsymbol{\beta}_2 + \text{residual}$  are numerically identical.

### Proof

1.  $\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\beta}}_2$  fulfill  $\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{M}_{\mathbf{X}}\mathbf{y} = \mathbf{X}_1\hat{\boldsymbol{\beta}}_1 + \mathbf{X}_2\hat{\boldsymbol{\beta}}_2 + \mathbf{M}_{\mathbf{X}}\mathbf{y}$ . Premultiplying with  $\mathbf{X}'_2\mathbf{M}_{\mathbf{X}_1}$  yields:

$$\begin{aligned} \mathbf{X}'_2\mathbf{M}_{\mathbf{X}_1}\mathbf{y} &= \mathbf{X}'_2\underbrace{\mathbf{M}_{\mathbf{X}_1}\mathbf{X}_1}_{\mathbf{O}}\hat{\boldsymbol{\beta}}_1 + \mathbf{X}'_2\mathbf{M}_{\mathbf{X}_1}\mathbf{X}_2\hat{\boldsymbol{\beta}}_2 + \underbrace{\mathbf{X}'_2\mathbf{M}_{\mathbf{X}_1}\mathbf{M}_{\mathbf{X}}}_{\mathbf{O}}\mathbf{y} \\ \Leftrightarrow \mathbf{X}'_2\mathbf{M}_{\mathbf{X}_1}\mathbf{y} &= \mathbf{X}'_2\mathbf{M}_{\mathbf{X}_1}\mathbf{X}_2\hat{\boldsymbol{\beta}}_2 \\ \Leftrightarrow \hat{\boldsymbol{\beta}}_2 &= (\mathbf{X}'_2\mathbf{M}_{\mathbf{X}_1}\mathbf{X}_2)^{-1}\mathbf{X}'_2\mathbf{M}_{\mathbf{X}_1}\mathbf{y} \end{aligned}$$

OLS estimator for regression  $\mathbf{M}_{\mathbf{X}_1}\mathbf{y} = \mathbf{M}_{\mathbf{X}_1}\mathbf{X}_2\boldsymbol{\beta}_2 + \text{residual}$  :

$$\hat{\boldsymbol{\beta}}_2 = (\mathbf{X}'_2\mathbf{M}'_{\mathbf{X}_1}\mathbf{M}_{\mathbf{X}_1}\mathbf{X}_2)^{-1}\mathbf{X}'_2\mathbf{M}'_{\mathbf{X}_1}\mathbf{M}_{\mathbf{X}_1}\mathbf{y}$$

We used the result

$$\mathbf{M}_{\mathbf{X}} \cdot \mathbf{M}_{\mathbf{X}_1} = \mathbf{M}_{\mathbf{X}}(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1}) = \mathbf{M}_{\mathbf{X}} - \mathbf{M}_{\mathbf{X}}\mathbf{X}_1(\mathbf{X}'_1\mathbf{X}_1)^{-1}\mathbf{X}'_1 = \mathbf{M}_{\mathbf{X}}$$

and therefore  $\mathbf{M}_{\mathbf{X}} \cdot \mathbf{M}_{\mathbf{X}_1}\mathbf{X}_2 = \mathbf{M}_{\mathbf{X}}\mathbf{X}_2 = \mathbf{O}$ .

2. Premultiplying  $\mathbf{y} = \mathbf{X}_1\hat{\boldsymbol{\beta}}_1 + \mathbf{X}_2\hat{\boldsymbol{\beta}}_2 + \mathbf{M}_{\mathbf{X}}\mathbf{y}$  with  $\mathbf{M}_{\mathbf{X}_1}$  yields:

$$\mathbf{M}_{\mathbf{X}_1}\mathbf{y} = \mathbf{M}_{\mathbf{X}_1}\mathbf{X}_1\hat{\boldsymbol{\beta}}_1 + \mathbf{M}_{\mathbf{X}_1}\mathbf{X}_2\hat{\boldsymbol{\beta}}_2 + \underbrace{\mathbf{M}_{\mathbf{X}_1}\mathbf{M}_{\mathbf{X}}}_{\mathbf{O}}\mathbf{y}$$

Thus the regression of  $\mathbf{M}_{\mathbf{X}_1}\mathbf{y}$  on  $\mathbf{M}_{\mathbf{X}_1}\mathbf{X}_2$  obtains residuals  $\mathbf{M}_{\mathbf{X}}\mathbf{y}$  just as the regression of  $\mathbf{y}$  on  $\mathbf{X}$ .

## 1.4 Goodness of Fit

Due to  $\mathbf{y} = \mathbf{P}_X \mathbf{y} + \mathbf{M}_X \mathbf{y}$  we could construct a goodness of fit measure

$$\frac{\|\mathbf{P}_X \mathbf{y}\|^2}{\|\mathbf{y}\|^2}$$

Suppose  $\mathbf{X}$  includes a constant, i.e.  $\boldsymbol{\iota}$  is a column of  $\mathbf{X}$ , we can increase

$$\frac{\|\mathbf{P}_X \mathbf{y} + \alpha \boldsymbol{\iota}\|^2}{\|\mathbf{y} + \alpha \boldsymbol{\iota}\|^2}$$

close to 1 by increasing  $\alpha$ . If we apply FWL and center both sides  $\mathbf{M}_\iota \mathbf{y} = \mathbf{P}_X \mathbf{M}_\iota \mathbf{y} + \mathbf{M}_X \mathbf{M}_\iota \mathbf{y}$  we get

$$R^2 = \frac{\|\mathbf{P}_X \mathbf{M}_\iota \mathbf{y}\|^2}{\|\mathbf{M}_\iota \mathbf{y}\|^2}$$

It measures the additional fit of a regression model over a benchmark model with just a constant.

## 1.5 Influential Points and Leverage

In simple regression we can assess high leverage points by inspecting the scatter plot of  $(x_t, y_t)$  observations.

To evaluate highly influential points in a multiple regression model we can run a regression in order to see whether observation  $t$  is influential

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \alpha \mathbf{e}_t + \mathbf{u}$$

with  $\mathbf{e}_t$  having 1 as the  $t^{\text{th}}$  components and zeros elsewhere. Applying FWL yields

$$\mathbf{M}_t \mathbf{y} = \mathbf{M}_t \mathbf{X} \boldsymbol{\beta} + \mathbf{M}_t \mathbf{u}$$

with

$$\begin{aligned} \mathbf{M}_t &= \mathbf{I} - \mathbf{e}_t (\mathbf{e}_t' \mathbf{e}_t)^{-1} \mathbf{e}_t' = \mathbf{I} - \mathbf{e}_t \mathbf{e}_t' \\ \Rightarrow \mathbf{M}_t \mathbf{y} &= \mathbf{y} - \mathbf{e}_t \mathbf{e}_t' \mathbf{y} = \mathbf{y} - y_t \mathbf{e}_t \end{aligned}$$

Furthermore,

$$\hat{\alpha} = \frac{\mathbf{e}_t' \mathbf{M}_t \mathbf{X} \mathbf{y}}{\mathbf{e}_t' \mathbf{M}_t \mathbf{X} \mathbf{e}_t} = \frac{\hat{u}_t}{1 - \mathbf{e}_t' \mathbf{P}_X \mathbf{e}_t} = \frac{\hat{u}_t}{1 - h_t}$$

Thus  $h_t$  is the  $t^{\text{th}}$  diagonal element of  $\mathbf{P}_X$ . It holds that  $1/n \leq h_t \leq 1$  and  $\bar{h}_t = n^{-1} \sum h_t = k/n$ .

The effects of the  $t^{\text{th}}$  observation on the estimates of  $\hat{\boldsymbol{\beta}}$  are analyzed by running a regression without  $t^{\text{th}}$  observation: (Define  $\mathbf{Z} = [\mathbf{X} : \mathbf{e}_t]$ )

$$\begin{aligned} \mathbf{y} &= \mathbf{X} \hat{\boldsymbol{\beta}}^{(t)} + \hat{\alpha} \mathbf{e}_t + \mathbf{M}_Z \mathbf{u} \\ \Leftrightarrow \mathbf{P}_X \mathbf{y} &= \mathbf{X} \hat{\boldsymbol{\beta}}^{(t)} + \hat{\alpha} \mathbf{P}_X \mathbf{e}_t + \mathbf{P}_X \mathbf{M}_Z \mathbf{u} \\ \Leftrightarrow \mathbf{X} (\hat{\boldsymbol{\beta}}^{(t)} - \hat{\boldsymbol{\beta}}) &= -\hat{\alpha} \mathbf{P}_X \mathbf{e}_t \\ \Leftrightarrow \mathbf{e}_t' \mathbf{X} (\hat{\boldsymbol{\beta}}^{(t)} - \hat{\boldsymbol{\beta}}) &= -\hat{\alpha} \mathbf{e}_t' \mathbf{P}_X \mathbf{e}_t \\ \Leftrightarrow \mathbf{X}_t \hat{\boldsymbol{\beta}}^{(t)} - \mathbf{X}_t \hat{\boldsymbol{\beta}} &= -\hat{\alpha} h_t = -\frac{\hat{u}_t}{1 - h_t} h_t = -\frac{h_t}{1 - h_t} \hat{u}_t \end{aligned}$$

## 2 Statistical Properties of OLS

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u} \quad \mathbf{u} \sim IID(\mathbf{0}, \sigma^2 \cdot \mathbf{I}) \quad (2)$$

Using the LS estimation principle:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}$$

Assuming  $\mathbf{X}$  as **non-stochastic** or **fix**:

$$E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\mathbf{u}) = \boldsymbol{\beta}$$

- This assumption is obviously appropriate in experimental situations.
- Social sciences seldomly have experimental designs.
- Assumption is mostly used for analytical simplicity.

It can be replaced by the assumption that  $\mathbf{X}$  is exogenous.

- Randomness of  $\mathbf{X}$  is *outside* our regression model (2).
- Error term is therefore independent of  $\mathbf{X}$ :

$$\mathbf{0} = E(\mathbf{u} \cdot \mathbf{h}(\mathbf{X})') = E_{\mathbf{X}}[E(\mathbf{u} \cdot \mathbf{h}(\mathbf{X})' | \mathbf{X})] = E_{\mathbf{X}} \left[ \underbrace{E(\mathbf{u} | \mathbf{X})}_{\mathbf{0}} \cdot \mathbf{h}(\mathbf{X})' \right]$$

- Thus assumption  $E(\mathbf{u} | \mathbf{X}) = \mathbf{0}$  suffices to ensure tractability:

$$\begin{aligned} E(\hat{\boldsymbol{\beta}}) &= \boldsymbol{\beta} + E((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}) \\ &= \boldsymbol{\beta} + E_{\mathbf{X}}[E((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u} | \mathbf{X})] \\ &= \boldsymbol{\beta} + E_{\mathbf{X}} \left[ (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \underbrace{E(\mathbf{u} | \mathbf{X})}_{\mathbf{0}} \right] = \boldsymbol{\beta} \end{aligned}$$

- $E(\mathbf{u} | \mathbf{X}) = \mathbf{0}$  is especially in time series regression contexts quite restrictive.

For instance, if

$$y_t = \mu + \rho \cdot y_{t-1} + u_t \quad t = 1, \dots, T$$

In matrix notation

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$$

with

$$\mathbf{X} = \begin{pmatrix} 1 & y_0 \\ \vdots & \vdots \\ 1 & y_{T-1} \end{pmatrix}; \quad \boldsymbol{\beta} = \begin{pmatrix} \mu \\ \rho \end{pmatrix}$$

$$E(\mathbf{u}|\mathbf{X}) = \mathbf{0} \text{ implies } E(u_t | \underbrace{y_0, \dots, y_t, y_{t+1}, \dots, y_{T-1}}_{\text{past and future values of } \mathbf{X}_t}) = 0.$$

It seems more appropriate in that context to assume  $E(u_t|\mathbf{X}_t) = 0$  with  $\mathbf{X}_t$  possibly containing components of  $\mathbf{X}_{t-1}, \mathbf{X}_{t-2}, \dots$  (lagged values, past history).

But this does not guarantee  $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$ .

- For cross-section analysis

$$E(u_i | \mathbf{X}_1, \dots, \mathbf{X}_i, \dots, \mathbf{X}_n) = 0$$

is not too restrictive.

## 2.1 Variance-Covariance matrix of $\hat{\boldsymbol{\beta}}$

$$\begin{aligned} V(\mathbf{u}) &= E((\mathbf{u} - E(\mathbf{u})) \cdot (\mathbf{u} - E(\mathbf{u}))') \\ &= E(\mathbf{u} \cdot \mathbf{u}') - E(\mathbf{u}) \cdot E(\mathbf{u})' \end{aligned}$$

In our context assuming  $E(\mathbf{u}) = \mathbf{0}$  and  $V(\mathbf{u}) = \sigma^2 \mathbf{I}$  we have  $E(\mathbf{u} \cdot \mathbf{u}') = \sigma^2 \mathbf{I}$ .

We utilize the very important property for linear transformations:

$$\begin{aligned} E(\mathbf{Y}) &= \boldsymbol{\mu}; \quad V(\mathbf{Y}) = E((\mathbf{Y} - E(\mathbf{Y}))(\mathbf{Y} - E(\mathbf{Y}))') = \boldsymbol{\Sigma} \\ \mathbf{Z} &= \mathbf{a} + \mathbf{B}\mathbf{Y} \\ \Rightarrow E(\mathbf{Z}) &= \mathbf{a} + \mathbf{B}\boldsymbol{\mu} \\ V(\mathbf{Z}) &= V(\mathbf{a} + \mathbf{B}\mathbf{Y}) = \mathbf{B} \cdot V(\mathbf{Y}) \cdot \mathbf{B}' = \mathbf{B}\boldsymbol{\Sigma}\mathbf{B}' \end{aligned}$$

- $\mathbf{X}$  is fix

$$\begin{aligned} V(\hat{\boldsymbol{\beta}}) &= V(\boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'V(\mathbf{u})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\sigma^2\mathbf{I}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\ &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1} \end{aligned}$$

- $\mathbf{X}$  is exogenous

$$\begin{aligned} V(\hat{\boldsymbol{\beta}}) &= E((\hat{\boldsymbol{\beta}} - E(\hat{\boldsymbol{\beta}})) \cdot (\hat{\boldsymbol{\beta}} - E(\hat{\boldsymbol{\beta}}))' | \mathbf{X}) \\ &= E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u} \cdot \mathbf{u}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} | \mathbf{X}] = \sigma^2(\mathbf{X}'\mathbf{X})^{-1} \end{aligned}$$

We can rewrite this covariance matrix as:

$$V(\hat{\boldsymbol{\beta}}) = \frac{\sigma^2}{n} \cdot \left( \frac{1}{n} \mathbf{X}'\mathbf{X} \right)^{-1}$$

Thus the precision of estimator  $\hat{\boldsymbol{\beta}}$  in terms of variation around  $\boldsymbol{\beta}$  depends on:

- $\sigma^2 = \text{Var}(u_i)$
- $n$  (sample size) (Just like the variance of  $\bar{y}$ )
- regressor matrix  $\mathbf{X}$ :

We are, without loss of generality, interested in the variance of  $\underbrace{\beta_1}_{1 \times 1}$ :

$$\mathbf{y} = \underbrace{\mathbf{X}_1}_{n \times 1} \beta_1 + \underbrace{\mathbf{X}_2}_{n \times k-1} \boldsymbol{\beta}_2 + \mathbf{u}$$

Applying FWL-Theorem: (Defining  $\mathbf{M}_2 = \mathbf{I} - \mathbf{X}_2(\mathbf{X}_2'\mathbf{X}_2)^{-1}\mathbf{X}_2'$ )

$$\begin{aligned} \mathbf{M}_2\mathbf{y} &= \mathbf{M}_2\mathbf{X}_1\beta_1 + \underbrace{\mathbf{M}_2\mathbf{X}_2}_{=0}\boldsymbol{\beta}_2 + res. \\ \Rightarrow \hat{\beta}_1 &= (\mathbf{X}_1'\mathbf{M}_2\mathbf{X}_1)^{-1}\mathbf{X}_1'\mathbf{M}_2\mathbf{y} = \frac{\mathbf{X}_1'\mathbf{M}_2\mathbf{y}}{\mathbf{X}_1'\mathbf{M}_2\mathbf{X}_1} \\ \Rightarrow \text{Var}(\hat{\beta}_1) &= \sigma^2 \cdot \frac{1}{\mathbf{X}_1'\mathbf{M}_2\mathbf{X}_1} \end{aligned}$$

Since  $\mathbf{M}_2\mathbf{X}_1$  are the regression residuals of the regression of  $\mathbf{X}_1$  on  $\mathbf{X}_2$ :

$$\begin{aligned} \mathbf{X}_1 &= \mathbf{X}_2 \cdot \boldsymbol{\theta} + res. \\ \mathbf{X}_1'\mathbf{M}_2\mathbf{X}_1 &= \mathbf{X}_1'\mathbf{M}_2'\mathbf{M}_2\mathbf{X}_1 = (\mathbf{M}_2\mathbf{X}_1)'(\mathbf{M}_2\mathbf{X}_1) \end{aligned}$$



**Multicollinearity:**  $\mathbf{X}_1$  is well explained by  $\mathbf{X}_2$

→ SSR small

→  $\|\mathbf{M}_2\mathbf{X}_1\|$  small

→  $\text{Var}(\hat{\beta}_1)$  large

## 2.2 Gauss-Markov-Theorem

$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$  is a linear function of  $\mathbf{y}$  and unbiased.

A class of linear function of  $\mathbf{y}$  estimators  $\tilde{\beta}$  being unbiased look like:

$$\tilde{\beta} = \mathbf{a} + \mathbf{A}\mathbf{y}$$

Since  $E(\tilde{\beta}) = \mathbf{a} + \mathbf{A} \cdot E(\mathbf{y}|\mathbf{X}) = \mathbf{a} + \mathbf{A}\mathbf{X}\beta$ , unbiasedness requires:

- $\mathbf{a} = \mathbf{0}$
- $\mathbf{A}\mathbf{X} = \mathbf{I}$

Gauss-Markov-Theorem states that within that class of linear unbiased estimators the OLS estimator is the most efficient.

### Proof

$$V(\tilde{\beta}) = V(\mathbf{A}\mathbf{y}) = V(\mathbf{A}\mathbf{X}\beta + \mathbf{A}\mathbf{u}) = \mathbf{A}\sigma^2\mathbf{I}\mathbf{A}' = \sigma^2\mathbf{A}\mathbf{A}'$$

Thus

$$\begin{aligned} V(\tilde{\beta}) - V(\hat{\beta}) &= \sigma^2(\mathbf{A}\mathbf{A}' - (\mathbf{X}'\mathbf{X})^{-1}) = \sigma^2(\mathbf{A}\mathbf{A}' - \underbrace{\mathbf{A}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{A}'}_{=\mathbf{I}}) \\ &= \sigma^2\mathbf{A} \cdot (\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\mathbf{A}' = \sigma^2\mathbf{A}\mathbf{M}\mathbf{A}' \\ &= \sigma^2\mathbf{A}\mathbf{M}\mathbf{M}'\mathbf{A}' = \sigma^2(\mathbf{A}\mathbf{M})(\mathbf{A}\mathbf{M})' \end{aligned}$$

This matrix is positive semidefinit, since

$$\begin{aligned} \mathbf{z}'(\mathbf{A}\mathbf{M})(\mathbf{A}\mathbf{M})'\mathbf{z} &= (\mathbf{z}'\mathbf{A}\mathbf{M}) \cdot (\mathbf{M}'\mathbf{A}'\mathbf{z}) = (\mathbf{M}'\mathbf{A}'\mathbf{z})' \underbrace{(\mathbf{M}'\mathbf{A}'\mathbf{z})}_{\mathbf{w}} \\ &= \mathbf{w}'\mathbf{w} = \sum w_i^2 \geq 0 \quad \forall \mathbf{z} \neq \mathbf{0} \end{aligned}$$

Taking for instance the  $i$ -th unit vector  $\mathbf{e}'_i = (0, \dots, 0, 1, 0, \dots, 0)' = \mathbf{z}'$ :

$$\mathbf{e}'_i(\mathbf{A}\mathbf{M})(\mathbf{A}\mathbf{M})'\mathbf{e}_i \geq 0 \Rightarrow \text{Var}(\tilde{\beta}_i) - \text{Var}(\hat{\beta}_i) \geq 0$$

The variances of all components of  $\tilde{\beta}$  are not larger than corresponding components in alternative  $\hat{\beta}$ .

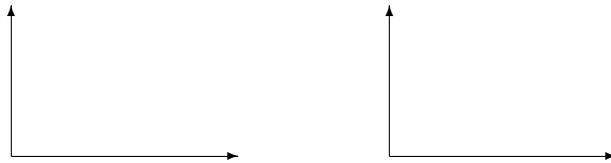
### 2.3 Some statistical results

- Transformation Theorem

$y$  has cdf  $F(y)$  and pdf  $f(y) = F'(y)$

Transformation  $z = h(y)$  with inverse  $y = h^{-1}(z)$

There must be a one-to-one relation  $F(y) = G(z)$



$$\text{Density } g(z) = G'(z) = \frac{dF(h^{-1}(z))}{dz} = f(h^{-1}(z)) \cdot \left| \frac{dh^{-1}(z)}{dz} \right|$$

- Standard normal distribution:  $y \sim N(0; 1)$

$$\phi(y) = \frac{1}{\sqrt{2\pi}} \cdot \exp\left(-\frac{1}{2}y^2\right)$$

- Linear transformation  $z = a + by$

$$\begin{aligned} \Rightarrow y &= \frac{z - a}{b} \quad \Rightarrow \quad \frac{dy}{dz} = \frac{1}{b} \\ f(z) &= \frac{1}{\sqrt{2\pi}} \frac{1}{b} \exp\left(-\frac{1}{2} \left(\frac{z - a}{b}\right)^2\right) = \frac{1}{b} \phi\left(\frac{z - a}{b}\right) \end{aligned}$$

- The sum of two independent  $N(0; 1)$  random variables is again normally distributed.

- Multivariate normal distribution

– independent standard normal components:  $\mathbf{y} \sim N(0; \mathbf{I})$

– linear transformation  $\mathbf{z} = \boldsymbol{\mu} + \mathbf{B}\mathbf{y} \stackrel{B \text{ full rank}}{\Leftrightarrow} \mathbf{y} = \mathbf{B}^{-1}(\mathbf{z} - \boldsymbol{\mu})$

Transformation theorem (multivariate version) tells us that  $\mathbf{z}$  is multivariate normal with  $E(\mathbf{z}) = \boldsymbol{\mu}$  and  $V(\mathbf{z}) = \mathbf{B}\mathbf{I}\mathbf{B}' = \mathbf{B}\mathbf{B}' = \boldsymbol{\Sigma}$ .

We revert that argument:

$V(\mathbf{z}) = \boldsymbol{\Sigma}$  for which we find a decomposition  $\boldsymbol{\Sigma} = \mathbf{B}\mathbf{B}'$

$\Rightarrow \mathbf{B}^{-1}(\mathbf{z} - \mathbf{y})$  is then  $N(\mathbf{0}; \mathbf{I})$

- Chi-Square ( $\chi^2$ ) distribution

$y_1, \dots, y_n$  are independent standard normally distributed, i.e.  $\mathbf{y} \sim N(\mathbf{0}; \mathbf{I})$

$$z = \sum^n y_i^2 = \|\mathbf{y}\|^2 = \mathbf{y}'\mathbf{y} \sim \chi_n^2$$

It follows that  $E(z) = n$  and  $\text{Var}(z) = 2 \cdot n$

- If  $\mathbf{y} \sim N(\mathbf{0}; \boldsymbol{\Sigma})$ , then

$$\mathbf{y}'\boldsymbol{\Sigma}^{-1}\mathbf{y} \sim \chi_n^2$$

### Proof

Assume  $\boldsymbol{\Sigma} = \mathbf{B}\mathbf{B}'$  as above. Then  $\boldsymbol{\Sigma}^{-1} = (\mathbf{B}')^{-1}\mathbf{B}^{-1}$

$$\Rightarrow \mathbf{y}'(\mathbf{B}')^{-1}\mathbf{B}^{-1}\mathbf{y} = \mathbf{w}'\mathbf{w}$$

Since  $\mathbf{B}^{-1}\mathbf{y}$  is  $N(\cdot)$  with  $V(\mathbf{B}^{-1}\mathbf{y}) = \mathbf{B}^{-1}\boldsymbol{\Sigma}(\mathbf{B}^{-1})' = \mathbf{B}^{-1}\mathbf{B}\mathbf{B}'(\mathbf{B}')^{-1} = \mathbf{I}$

$$\mathbf{w}'\mathbf{w} \sim \chi_n^2$$

- The orthogonal projection matrix  $\mathbf{P}_X = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$  projects orthogonally onto  $S(\mathbf{X})$  with  $\mathbf{X}$  a  $k$ -dimensional linear subspace

$$\mathbf{y} \sim N(\mathbf{0}; \mathbf{I}) \Rightarrow \mathbf{y}'\mathbf{P}_X\mathbf{y} \sim \chi_k^2$$

### Proof

$$\underbrace{\mathbf{y}'\mathbf{X}}_{\mathbf{z}'}(\mathbf{X}'\mathbf{X})^{-1}\underbrace{\mathbf{X}'\mathbf{y}}_z = \mathbf{z}'(\mathbf{X}'\mathbf{X})^{-1}z$$

Since  $\mathbf{y} \sim N(\mathbf{0}; \mathbf{I})$  it follows that  $\mathbf{z} = \mathbf{X}'\mathbf{y} \sim N(\mathbf{0}; \mathbf{X}'\mathbf{I}\mathbf{X})$  and therefore

$$\mathbf{z}'(\mathbf{X}'\mathbf{X})^{-1}z \sim \chi_k^2$$

- F-distribution

If  $y_1 \sim \chi_{m_1}^2$  and  $y_2 \sim \chi_{m_2}^2$  and  $y_1, y_2$  are independent, then

$$\frac{y_1/m_1}{y_2/m_2} \sim F(m_1, m_2)$$

- $\mathbf{A}$  and  $\mathbf{B}$  are each symmetric and idempotent,  $\mathbf{z} \sim N(\mathbf{0}; \mathbf{I})$ . Then  $\mathbf{z}'\mathbf{A}\mathbf{z}$  and  $\mathbf{z}'\mathbf{B}\mathbf{z}$  are independent if and only if  $\mathbf{A} \cdot \mathbf{B} = \mathbf{O}$

**Proof**

$$\begin{aligned} \mathbf{z}'\mathbf{A}\mathbf{z} &= \mathbf{z}'\mathbf{A}'\mathbf{A}\mathbf{z} = \mathbf{z}'_1\mathbf{z}_1 \quad \text{with } \mathbf{z}_1 \sim N(\mathbf{0}; \underbrace{\mathbf{A}\mathbf{I}\mathbf{A}'}_{\mathbf{A}}) \text{ and} \\ \mathbf{z}'\mathbf{B}\mathbf{z} &= \mathbf{z}'\mathbf{B}'\mathbf{B}\mathbf{z} = \mathbf{z}'_2\mathbf{z}_2 \quad \text{with } \mathbf{z}_2 \sim N(\mathbf{0}; \underbrace{\mathbf{B}\mathbf{I}\mathbf{B}'}_{\mathbf{B}}) \\ E(\mathbf{z}_1 \cdot \mathbf{z}'_2) &= E(\mathbf{A}\mathbf{z}(\mathbf{B}\mathbf{z})') = E(\mathbf{A}\mathbf{z}\mathbf{z}'\mathbf{B}') = \mathbf{A}\mathbf{I}\mathbf{B}' \end{aligned}$$

## 2.4 Estimation of $\sigma^2$

Since  $\hat{\mathbf{u}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$  is always orthogonal on  $S(\mathbf{X})$  it is shorter than  $\mathbf{u}$ :

$$\begin{aligned} \|\hat{\mathbf{u}}\| &\leq \|\mathbf{u}\| \Leftrightarrow \|\hat{\mathbf{u}}\|^2 \leq \|\mathbf{u}\|^2 \\ \Rightarrow E(\|\hat{\mathbf{u}}\|^2) &= E(\hat{\mathbf{u}}'\hat{\mathbf{u}}) \leq E(\mathbf{u}'\mathbf{u}) = n \cdot \sigma^2 \end{aligned}$$

Since  $\hat{\mathbf{u}} = \mathbf{M}_\mathbf{X}\mathbf{y} = \mathbf{M}_\mathbf{X}(\mathbf{X}\boldsymbol{\beta} + \mathbf{u}) = \mathbf{M}_\mathbf{X}\mathbf{X}\boldsymbol{\beta} + \mathbf{M}_\mathbf{X}\mathbf{u}$  we get

$$V(\hat{\mathbf{u}}) = \mathbf{M}_\mathbf{X}\sigma^2\mathbf{I}\mathbf{M}_\mathbf{X} = \sigma^2\mathbf{M}_\mathbf{X}$$

$$\begin{aligned} E(\hat{\mathbf{u}}'\hat{\mathbf{u}}) &= E(\mathbf{u}'\mathbf{M}_\mathbf{X}\mathbf{u}) = \text{tr}(E(\mathbf{u}'\mathbf{M}_\mathbf{X}\mathbf{u})) \\ &= E(\text{tr}(\mathbf{u}'\mathbf{M}_\mathbf{X}\mathbf{u})) = E(\text{tr}(\mathbf{M}_\mathbf{X}\mathbf{u}\mathbf{u}')) = \text{tr}(E(\mathbf{M}_\mathbf{X}\mathbf{u}\mathbf{u}')) \\ &= \text{tr}(\mathbf{M}_\mathbf{X} \cdot \sigma^2\mathbf{I}) = \sigma^2(\text{tr}(\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')) \\ &= \sigma^2(\text{tr}(\mathbf{I}) - \text{tr}((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X})) = \sigma^2(n - k) \end{aligned}$$

We therefore define  $s^2 = \frac{1}{n-k} \sum_{t=1}^n \hat{u}_t^2 = \frac{\hat{\mathbf{u}}'\hat{\mathbf{u}}}{n-k}$ .

## 2.5 Hypothesis testing

Our model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$ ;  $\mathbf{u} \sim IID(\mathbf{0}, \sigma^2\mathbf{I})$  yields  $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}$  with  $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$  and  $V(\hat{\boldsymbol{\beta}}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$ .

If we assume  $\mathbf{u} \sim N(\mathbf{0}; \sigma^2\mathbf{I})$ , then

$$\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}; \sigma^2(\mathbf{X}'\mathbf{X})^{-1})$$

Thus

$$\frac{(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})'(\mathbf{X}'\mathbf{X})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})}{\sigma^2} \sim \chi_k^2$$

However, it depends on  $\sigma^2$  which is unknown. Rewriting this yields

$$\begin{aligned} \frac{(\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{X}\boldsymbol{\beta})}{\sigma^2} &= \frac{(\mathbf{P}\mathbf{y} - \mathbf{P}(E(\mathbf{y}|\mathbf{X})))'(\mathbf{P}\mathbf{y} - \mathbf{P}(E(\mathbf{y}|\mathbf{X})))}{\sigma^2} \\ &= \frac{(\mathbf{P}\mathbf{u})'\mathbf{P}\mathbf{u}}{\sigma^2} = \frac{\mathbf{u}'\mathbf{P}\mathbf{u}}{\sigma^2} \sim \chi_k^2 \end{aligned}$$

We just saw in the previous section, that

$$s^2 = \frac{\hat{\mathbf{u}}'\hat{\mathbf{u}}}{n-k} = \frac{\mathbf{u}'\mathbf{M}\mathbf{u}}{n-k}$$

Note that  $\mathbf{u} \sim N(\mathbf{0}; \sigma^2\mathbf{I}) \Rightarrow \mathbf{u}/\sigma \sim N(\mathbf{0}; \mathbf{I})$  and

$$\frac{\hat{\mathbf{u}}'\hat{\mathbf{u}}}{\sigma^2} = \frac{\mathbf{u}'\mathbf{M}\mathbf{u}}{\sigma^2} \sim \chi_{n-k}^2$$

Since  $\mathbf{P}_X \cdot \mathbf{M}_X = \mathbf{P}_X(\mathbf{I} - \mathbf{P}_X) = \mathbf{O}$ ,  $\mathbf{u}'\mathbf{M}\mathbf{u}$  is independent of  $\mathbf{u}'\mathbf{P}\mathbf{u}$  and therefore

$$\frac{\frac{\mathbf{u}'\mathbf{P}\mathbf{u}}{k}}{\frac{\mathbf{u}'\mathbf{M}\mathbf{u}}{n-k}} = \frac{(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})'\mathbf{X}'\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})}{s^2} \sim F(k, n-k)$$

$\sigma^2$  cancels out.

### 3 Asymptotic Theory

#### 3.1 Probability Limit

We consider an  $n$ -vector of random components  $\mathbf{y}^{(n)}$  and a vector valued function  $\mathbf{a}(\mathbf{y}^{(n)})$ .

We call the limiting (random) vector  $\mathbf{a}_0$  if  $\mathbf{a}(\mathbf{y}^{(n)})$  tends in probability to that  $\mathbf{a}_0$ :

$$\lim_{n \rightarrow \infty} P(\|\mathbf{a}(\mathbf{y}^{(n)}) - \mathbf{a}_0\| < \varepsilon) = 1$$

Shortly we write  $\text{plim}(\mathbf{a}(\mathbf{y}^{(n)})) = \mathbf{a}_0$  and say  $\mathbf{a}(\mathbf{y}^{(n)})$  is consistent estimator of  $\mathbf{a}_0$ .  $\mathbf{a}_0$  might be stochastic or nonstochastic.

**Example:**

$$y_i \sim IID(\mu, \sigma^2) \quad i = 1, \dots, n$$

We consider the arithmetic mean  $a(\mathbf{y}^{(n)}) = \frac{1}{n} \mathbf{y}^{(n)'} \cdot \mathbf{1} = \frac{1}{n} \sum_{i=1}^n y_i$ .

Thus  $a(\mathbf{y}^{(n)})$  is a scalar function and the Euclidian norm  $\|\cdot\|$  reduces to  $|\cdot|$ .

The first two moments of  $a(\mathbf{y}^{(n)}) = \bar{y}$ :

$$\begin{aligned} E(\bar{y}^{(n)}) &= E\left(\frac{1}{n} \sum y_i\right) = \frac{1}{n} \sum E(y_i) = \mu \\ \text{Var}(\bar{y}^{(n)}) &= \text{Var}\left(\frac{1}{n} \sum y_i\right) \stackrel{\text{independent}}{=} \left(\frac{1}{n}\right)^2 n \text{Var}(y_i) = \frac{1}{n^2} n \sigma^2 = \frac{\sigma^2}{n} \end{aligned}$$

In order to prove consistency of  $\bar{y}$  for  $\mu$ , we need Chebychev's inequality:

$$P(|Z - E(Z)| > \varepsilon) \leq \frac{\text{Var}(Z)}{\varepsilon^2} \quad \forall \varepsilon > 0$$

Applying it to  $\bar{y}$ :

$$P(|\bar{y}^{(n)} - \mu| > \varepsilon) \leq \frac{\sigma^2}{n \cdot \varepsilon^2}$$

And taking limitis:

$$\lim_{n \rightarrow \infty} P(|\bar{y}^{(n)} - \mu| > \varepsilon) \leq 0 \Leftrightarrow \lim_{n \rightarrow \infty} P(|\bar{y} - \mu| < \varepsilon) = 1$$

or

$$\text{plim}(\bar{y}^{(n)}) = \mu$$

The example above is called the **Law of Large Numbers (LLN)**. More generally

$$\text{plim} \left( \frac{1}{n} \sum h(y_i) \right) = E(h(y_i))$$

Notes:

- We assumed *IID* which can be replaced by independence plus the same first moments and bounded variance
- Independence can be relaxed to weak dependence plus weak form of heteroskedasticity:

$$\text{plim} \left( \frac{1}{n} \sum h(y_i) \right) = \lim \left( \frac{1}{n} \sum E(h(y_i)) \right)$$

**Property**

$$\text{plim} \mathbf{a}(\mathbf{y}^{(n)}) = \mathbf{a}_0$$

$\mathbf{a}_0$  is nonstochastic and  $g(\cdot)$  is a smooth function. Then

$$\text{plim} g(\mathbf{a}(\mathbf{y}^{(n)})) = g(\mathbf{a}_0)$$

**Example**

We know that the sequence of random variables  $\bar{y}^{(1)}, \bar{y}^{(2)}, \dots, \bar{y}^{(n)}$  has  $\text{plim} \bar{y}^{(n)} = \mu$ .

Thus

$$\text{plim} g(\bar{y}^{(n)}) = g(\mu)$$

### 3.2 Consistency of OLS estimator

In order to focus on  $\hat{\beta} = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}$  we have to see that

$$\mathbf{X}'\mathbf{X} = \sum \mathbf{X}'_t \cdot \mathbf{X}_t \quad \mathbf{X}_t \text{ is } t^{\text{th}} \text{ row of } \mathbf{X}$$

Thus  $ij$  element of that matrix

$$(\mathbf{X}'\mathbf{X})_{ij} = \sum^n \mathbf{X}_{ti} \cdot \mathbf{X}_{tj}$$

which is probably growing with increasing  $n$ . Therefore, we consider

$$\left(\frac{1}{n}\mathbf{X}'\mathbf{X}\right)_{ij} = \frac{1}{n} \sum_{t=1}^n \mathbf{X}_{ti} \cdot \mathbf{X}_{tj}$$

for which it is not unreasonable to assume that

$$\text{plim} \left(\frac{1}{n}\mathbf{X}'\mathbf{X}\right)_{ij} = E(\mathbf{x}_i \cdot \mathbf{x}_j)$$

as a form of LLN. Specifically, we assume

$$\text{plim} \left(\frac{1}{n}\mathbf{X}'\mathbf{X}\right) = \mathbf{S}_{\mathbf{X}'\mathbf{X}}$$

This excludes regressors like a time trend, since e.g.

$$\frac{1}{T} \sum_{t=1}^T t = \frac{1}{T} \cdot \frac{T(T+1)}{2}$$

or

$$\frac{1}{T} \sum_{t=1}^T t^2$$

are not converging.

Thus we rewrite

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \boldsymbol{\beta} + \left(\frac{1}{n}\mathbf{X}'\mathbf{X}\right)^{-1} \cdot \frac{1}{n}\mathbf{X}'\mathbf{u} \\ \Rightarrow \text{plim}(\hat{\boldsymbol{\beta}}) &= \boldsymbol{\beta} + (\mathbf{S}_{\mathbf{X}'\mathbf{X}})^{-1} \cdot \text{plim} \frac{1}{n}\mathbf{X}'\mathbf{u} \end{aligned}$$

If we assume that  $E(u_t|\mathbf{X}_t) = 0$ , which is much weaker than  $E(\mathbf{u}|\mathbf{X}) = \mathbf{0}$  implying  $E(u_t|\mathbf{X}) = 0$ , we obtain

$$\begin{aligned} \mathbf{X}'_t E(u_t|\mathbf{X}_t) &= \mathbf{X}'_t \cdot 0 \Leftrightarrow E(\mathbf{X}'_t u_t|\mathbf{X}_t) = \mathbf{0} \\ &\Leftrightarrow E(\mathbf{X}'_t \cdot u_t) = \mathbf{0} \end{aligned}$$

As a consequence,

$$\text{plim} \left(\frac{1}{n} \sum_{t=1}^n \mathbf{X}'_t u_t\right) = \text{plim} \left(\frac{1}{n}\mathbf{X}'\mathbf{u}\right) \stackrel{\text{LLN}}{=} E(\mathbf{X}'_t u_t) = \mathbf{0}$$

As a result,

$$\text{plim}(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$$

Note that taking expectations and plim are two different concepts.

### Examples



- Consider

$$y_t = \beta_1 + \beta_2 \cdot \frac{1}{t} + u_t \quad u_t \sim IID(0; \sigma^2)$$

OLS is unbiased. But  $\hat{\beta}_2$  is not consistent, since as  $n \rightarrow \infty \lim \frac{1}{t} = 0$ . Thus more observations provide less and less information about  $\beta_2$ .

- Estimating population parameter  $\mu = E(y_t)$

$$\begin{aligned} \hat{\mu}_1 &= \frac{1}{n-1} \sum_{t=1}^n y_t & \text{plim } \hat{\mu}_1 &= \mu \\ \hat{\mu}_2 &= 0,01y_1 + \frac{0,99}{n-1} \sum_{t=2}^n y_t \\ E(\hat{\mu}_2) &= 0,01\mu + 0,99 \cdot \frac{1}{n-1} (n-1) \cdot \mu = \mu \\ \text{plim } \hat{\mu}_2 &= 0,99\mu + 0,01y_1 & \text{which is random} \end{aligned}$$

### 3.3 Asymptotic normality

The second fundamental result is called **central limit theorem (CLT)**. Again, there are various types of CLT's.

Considering  $\bar{y}^{(n)}$ , for which we know that  $\text{Var}(\bar{y}^{(n)}) = \frac{\sigma^2}{n}$ . In the limit, this variance is 0. Thus  $\bar{y}^{(n)}$  has a degenerate pdf. But note that

$$\sqrt{n} \bar{y}^{(n)} \quad \text{has} \quad \text{Var}(\sqrt{n} \bar{y}^{(n)}) = n \cdot \text{Var}(\bar{y}^{(n)}) = \sigma^2$$

(Variance stabilizing transformation)

Thus  $\sqrt{n}$  is called the rate of convergence

#### Central Limit Theorem

$$\begin{aligned} & y_t \sim IID(\mu; \sigma^2) \\ \lim_{n \rightarrow \infty} P \left( \frac{\bar{y}^{(n)} - \mu}{\sigma} \sqrt{n} < z \right) &= \Phi(z) \end{aligned}$$

We can also write

$$\sqrt{n} \left( \frac{\bar{y}^{(n)} - \mu}{\sigma} \right) \xrightarrow{d} N(0; 1)$$

or we write

$$\bar{y}^{(n)} \stackrel{a}{\sim} N\left(\mu; \frac{\sigma^2}{n}\right)$$

Reformulating  $\bar{y}^{(n)}$  in the above version of CLT gives the CLT which is most common in econometrics:

### CLT

Let  $y_t$  be a sequence of random variables,  $t = 1, \dots, \infty$  with  $E(y_t) = 0$ , then

$$\text{plim} \frac{1}{\sqrt{n}} \sum y_t = y_0 \sim N\left(0; \lim_{n \rightarrow \infty} \frac{1}{n} \sum \text{Var}(y_t)\right)$$

In a multivariate version:

Suppose we have a sequence of random  $m$ -vectors ( $m$  fixed)  $\mathbf{y}_t$  with  $E(\mathbf{y}_t) = \mathbf{0}$ , then

$$\text{plim} \frac{1}{\sqrt{n}} \sum \mathbf{y}_t = \mathbf{y}_0 \sim N(\mathbf{0}; \lim_{n \rightarrow \infty} \frac{1}{n} \sum V(\mathbf{y}_t))$$

Now we turn to our regression model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u} \quad \mathbf{u} \sim IID(0; \sigma^2 \mathbf{I})$$

with  $E(u_t | \mathbf{X}_t) = 0$  and  $E(u_t^2 | \mathbf{X}_t) = \sigma^2$ ,  $\text{plim} \left(\frac{1}{n} \mathbf{X}' \mathbf{X}\right) = \mathbf{S}_{\mathbf{X}' \mathbf{X}}$ .

Considering

$$\mathbf{v} = \frac{1}{\sqrt{n}} \mathbf{X}' \mathbf{u} = \frac{1}{\sqrt{n}} \sum_{t=1}^n u_t \cdot \mathbf{X}'_t$$

By assumption

$$\begin{aligned} E(u_t | \mathbf{X}_t) = 0 &\Rightarrow E_X[\mathbf{X}'_t E(u_t | \mathbf{X}_t)] = E_X E(u_t \mathbf{X}'_t | \mathbf{X}_t) = \mathbf{0} \\ &\Leftrightarrow E(u_t \mathbf{X}'_t) = \mathbf{0} \end{aligned}$$

so that we can apply the multivariate CLT:

$$\mathbf{v} \stackrel{a}{\sim} N\left(\mathbf{0}; \lim_{n \rightarrow \infty} \frac{1}{n} \sum \text{Var}(u_t \mathbf{X}'_t)\right) = N\left(\mathbf{0}; \lim_{n \rightarrow \infty} \frac{1}{n} \sum E(u_t^2 \mathbf{X}'_t \mathbf{X}_t)\right)$$

The covariance matrix is

(recall that  $E(u_t^2 | \mathbf{X}_t) = \sigma^2 \Leftrightarrow E_X[\mathbf{X}'_t \mathbf{X}_t E(u_t^2 | \mathbf{X}_t)] = \sigma^2 E(\mathbf{X}'_t \mathbf{X}_t)$ )

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \sum E(u_t^2 \mathbf{X}'_t \mathbf{X}_t) &= \lim_{n \rightarrow \infty} \sigma^2 \underbrace{\frac{1}{n} \sum E(\mathbf{X}'_t \mathbf{X}_t)}_{\text{LLN}} \\ &= \text{plim} \sigma^2 \frac{1}{n} \sum \mathbf{X}'_t \mathbf{X}_t \\ &= \text{plim} \sigma^2 \frac{1}{n} \mathbf{X}' \mathbf{X} = \sigma^2 \cdot \mathbf{S}_{\mathbf{X}' \mathbf{X}} \end{aligned}$$

Thus

$$\mathbf{v} \stackrel{a}{\sim} N(\mathbf{0}; \sigma^2 \cdot \mathbf{S}_{\mathbf{X}'\mathbf{X}})$$

and therefore

$$\begin{aligned} \text{plim} \left( \frac{1}{n} \mathbf{X}'\mathbf{X} \right)^{-1} \cdot \text{plim} \mathbf{v} &= \mathbf{S}_{\mathbf{X}'\mathbf{X}}^{-1} \cdot \text{plim} \mathbf{v} \\ &= \mathbf{S}_{\mathbf{X}'\mathbf{X}}^{-1} \mathbf{v}_0 \sim N(\mathbf{0}; \sigma^2 \mathbf{S}_{\mathbf{X}'\mathbf{X}}^{-1} \mathbf{S}_{\mathbf{X}'\mathbf{X}} \mathbf{S}_{\mathbf{X}'\mathbf{X}}^{-1}) = N(\mathbf{0}; \sigma^2 \mathbf{S}_{\mathbf{X}'\mathbf{X}}^{-1}) \end{aligned}$$

Thus we obtain for  $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta} + \left( \frac{1}{n} \mathbf{X}'\mathbf{X} \right)^{-1} \cdot \frac{1}{n} \mathbf{X}'\mathbf{u}$  and

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = \left( \frac{1}{n} \mathbf{X}'\mathbf{X} \right)^{-1} \underbrace{\frac{1}{\sqrt{n}} \mathbf{X}'\mathbf{u}}_v \sim N(\mathbf{0}; \sigma^2 \mathbf{S}_{\mathbf{X}'\mathbf{X}}^{-1}) \text{ or } \hat{\boldsymbol{\beta}} \stackrel{a}{\sim} N \left( \boldsymbol{\beta}, \frac{\sigma^2}{n} \mathbf{S}_{\mathbf{X}'\mathbf{X}}^{-1} \right)$$

## 4 Generalized Method of Moments

### 4.1 Method-of-Moments estimation

The basic idea of the method of moments is to replace population means by sample means.

For our model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u} \quad \mathbf{u} \sim IID$$

we assume  $E(u_t|\mathbf{X}_t) = 0$  which implies

$$\begin{aligned} \mathbf{X}'_t \cdot E(u_t|\mathbf{X}_t) &= \mathbf{X}'_t \cdot 0 = \mathbf{0} \Leftrightarrow E(\mathbf{X}'_t u_t|\mathbf{X}_t) = \mathbf{0} \\ \Leftrightarrow E_{X_t} E(\mathbf{X}'_t u_t|\mathbf{X}_t) &= E(\mathbf{X}'_t u_t) = E_{X_t}(\mathbf{0}) = \mathbf{0} \end{aligned}$$

Thus replacing  $E(\mathbf{X}'_t u_t)$  by its sample counterpart

$$\frac{1}{n} \sum_{t=1}^n \mathbf{X}'_t u_t = \frac{1}{n} \mathbf{X}' \mathbf{u} = \frac{1}{n} \mathbf{X}' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \stackrel{!}{=} \mathbf{0}$$

Multiplication with  $n$  yields the LS normal equations and therefore we obtain

$$\hat{\boldsymbol{\beta}}_{MM} = \hat{\boldsymbol{\beta}}_{LS} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$$

### 4.2 Errors-in-Variables

Suppose the true model is

$$\mathbf{y} = \mathbf{X}^* \boldsymbol{\beta} + \mathbf{u}^* \quad \mathbf{u}^* \sim IID(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (3)$$

A problem arises whenever  $\mathbf{X}^*$  is not observed directly. Instead we observe  $\mathbf{X}^*$  with some measurement error:

$$\mathbf{X} = \mathbf{X}^* + \mathbf{V}$$

Inserting this into model (3) yields

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} - \mathbf{V}\boldsymbol{\beta} + \mathbf{u}^*$$

Thus the new error term is  $\mathbf{u} = \mathbf{u}^* - \mathbf{V}\boldsymbol{\beta}$  and it is quite obvious that  $u_t = u_t^* - \mathbf{V}_t\boldsymbol{\beta}$  is correlated with  $\mathbf{X}_t = \mathbf{X}_t^* + \mathbf{V}_t$ . Thus  $E(\mathbf{X}_t' \mathbf{u}_t) \neq \mathbf{0}$ .

However, if we find instruments  $\mathbf{W}_t$  that are uncorrelated with  $u_t$  implying that they are uncorrelated with  $\mathbf{V}_t$  but highly correlated with  $\mathbf{X}_t$  we obtain the MM estimator:

$$\begin{aligned} \frac{1}{n} \sum \mathbf{W}_t' \cdot u_t &= \frac{1}{n} \mathbf{W}' \mathbf{u} = \frac{1}{n} \mathbf{W}' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \stackrel{!}{=} 0 \\ \Leftrightarrow \mathbf{W}' \mathbf{y} &= \mathbf{W}' \mathbf{X}\boldsymbol{\beta} \end{aligned}$$

Having exactly as many instruments as regressors and  $\mathbf{W}' \mathbf{X}$  having full rank, we get

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{IV} &= (\mathbf{W}' \mathbf{X})^{-1} \mathbf{W}' \mathbf{y} \\ &= (\mathbf{W}' \mathbf{X})^{-1} \mathbf{W}' \mathbf{X}\boldsymbol{\beta} + (\mathbf{W}' \mathbf{X})^{-1} \mathbf{W}' \mathbf{u} \\ &= \boldsymbol{\beta} + (\mathbf{W}' \mathbf{X})^{-1} \mathbf{W}' \mathbf{u} \end{aligned}$$

This is a linear transformation in  $\mathbf{y}$  and  $\mathbf{y}$  is projected onto  $\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{W}' \mathbf{X})^{-1} \mathbf{W}' \mathbf{y}$  so that the projector now is

$$\mathbf{P} = \mathbf{X}(\mathbf{W}' \mathbf{X})^{-1} \mathbf{W}'$$

For consistency we write  $\hat{\boldsymbol{\beta}}_{IV} = \boldsymbol{\beta} + (\mathbf{W}' \mathbf{X})^{-1} (\mathbf{W}' \mathbf{u})$  and see that if

$$\text{plim} \frac{1}{n} \mathbf{W}' \mathbf{X} = \mathbf{S}_{\mathbf{W}' \mathbf{X}}$$

we obtain consistency of  $\hat{\boldsymbol{\beta}}_{IV}$  since

$$\text{plim} \frac{1}{n} \mathbf{W}' \mathbf{u} = \text{plim} \frac{1}{n} \sum \mathbf{W}_t' u_t \stackrel{\text{LLN}}{=} \lim \frac{1}{n} \sum E(\mathbf{W}_t' u_t) \stackrel{\text{assumption } E(u_t | \mathbf{W}_t) = 0}{=} \mathbf{0}$$

Examining  $\sqrt{n}(\hat{\boldsymbol{\beta}}_{IV} - \boldsymbol{\beta}) = n^{-1}(\mathbf{W}' \mathbf{X})^{-1} n^{-\frac{1}{2}} \mathbf{W}' \mathbf{u}$  we can apply a multivariate CLT as above and obtain

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_{IV} - \boldsymbol{\beta}) \stackrel{a}{\sim} N(\mathbf{0}; \sigma^2 \mathbf{S}_{\mathbf{W}' \mathbf{X}}^{-1} \mathbf{S}_{\mathbf{W}' \mathbf{W}} (\mathbf{S}_{\mathbf{W}' \mathbf{X}}^{-1})')$$

where  $\text{plim} \frac{1}{n} \mathbf{W}' \mathbf{W} = \mathbf{S}_{\mathbf{W}' \mathbf{W}}$ . Thus the asymptotic variance is

$$\begin{aligned} \sigma^2 \mathbf{S}_{\mathbf{W}' \mathbf{X}}^{-1} \mathbf{S}_{\mathbf{W}' \mathbf{W}} (\mathbf{S}_{\mathbf{W}' \mathbf{X}}^{-1})' &= \sigma^2 \text{plim} \left( \frac{1}{n} \mathbf{W}' \mathbf{X} \right)^{-1} \frac{1}{n} \mathbf{W}' \mathbf{W} \left( \frac{1}{n} \mathbf{X}' \mathbf{W} \right)^{-1} \\ &= \sigma^2 \text{plim} \left( \frac{1}{n} \mathbf{X}' \mathbf{W} \left( \frac{1}{n} \mathbf{W}' \mathbf{W} \right)^{-1} \frac{1}{n} \mathbf{W}' \mathbf{X} \right)^{-1} \end{aligned}$$

$$= \sigma^2 \text{plim} \left( \frac{1}{n} \mathbf{X}' \mathbf{P}_W \mathbf{X} \right)^{-1} \quad (4)$$

If we could choose instruments  $\mathbf{W}$  we could use those ones implying the least  $V(\sqrt{n}(\hat{\boldsymbol{\beta}}_{IV} - \boldsymbol{\beta}))$ .

If we write for the explanatory variables

$$\mathbf{X} = \bar{\mathbf{X}} + \mathbf{z} \quad E(z_t | \Omega_t) = 0$$

where  $\Omega_t$  is the information set in  $t$  and  $\bar{\mathbf{X}}_t$  denoting the  $t$ -th row of  $\bar{\mathbf{X}}$ .

These conditional expectations of  $\mathbf{X}_t$  are the optimal instruments. To see that we first note that

$$\begin{aligned} \text{plim} \frac{1}{n} \mathbf{X}' \mathbf{W} &= \lim \frac{1}{n} E(\mathbf{X}' \mathbf{W}) = \lim \frac{1}{n} E((\bar{\mathbf{X}} + \mathbf{z})' \mathbf{W}) \\ &= \lim \frac{1}{n} E(\bar{\mathbf{X}}' \mathbf{W}) = \text{plim} \frac{1}{n} \bar{\mathbf{X}}' \mathbf{W} \end{aligned}$$

Thus the covariance matrix becomes to  $\sigma^2 \text{plim} \left( \frac{1}{n} \bar{\mathbf{X}}' \mathbf{P}_W \bar{\mathbf{X}} \right)^{-1}$ .

Now setting  $\mathbf{W} = \bar{\mathbf{X}}$  we get

$$\sigma^2 \text{plim} \left( \frac{1}{n} \bar{\mathbf{X}}' \mathbf{W} (\mathbf{W}' \mathbf{W})^{-1} \mathbf{W}' \bar{\mathbf{X}} \right)^{-1} = \sigma^2 \text{plim} \left( \frac{1}{n} \bar{\mathbf{X}}' \bar{\mathbf{X}} \right)^{-1} \quad (5)$$

Considering the differences between (4) and (5) or rather their inverses (to show pos. semi-definiteness of  $\mathbf{A} - \mathbf{B}$  it is equivalent to prove pos. semi-definiteness of  $\mathbf{B}^{-1} - \mathbf{A}^{-1}$ )

$$\sigma^2 \text{plim} \left[ \frac{1}{n} (\bar{\mathbf{X}}' \bar{\mathbf{X}} - \bar{\mathbf{X}}' \mathbf{P}_W \bar{\mathbf{X}}) \right] = \sigma^2 \text{plim} \frac{1}{n} \bar{\mathbf{X}}' \mathbf{M}_W \bar{\mathbf{X}}$$

which is a positiv semi-definit matrix showing the optimality of using  $\bar{\mathbf{X}}$  as instruments.

So far we assumed that  $\mathbf{W}$  has the same number of columns as  $\mathbf{X}$ . What if we have more instruments than regressors?

The basic idea is to use  $k$  linear combinations of  $\mathbf{W}$  (e.g.  $\mathbf{W} \cdot \mathbf{J}$ ) to obtain a just identified problem:

$$(\mathbf{W} \mathbf{J})' (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}) \stackrel{!}{=} \mathbf{0} \Leftrightarrow \mathbf{J}' \mathbf{W}' \mathbf{y} = \mathbf{J}' \mathbf{W}' \mathbf{X} \hat{\boldsymbol{\beta}} \Leftrightarrow \hat{\boldsymbol{\beta}}_{GIVE} = (\mathbf{J}' \mathbf{W}' \mathbf{X})^{-1} \mathbf{J}' \mathbf{W}' \mathbf{y}$$

This generalized estimator  $\sqrt{n}(\hat{\boldsymbol{\beta}}_{GIVE} - \boldsymbol{\beta})$  has asymptotic covariance matrix

$$\sigma^2 \text{plim } n(\mathbf{J}'\mathbf{W}'\mathbf{X})^{-1}\mathbf{J}'\mathbf{W}'\mathbf{W}\mathbf{J}(\mathbf{X}'\mathbf{W}\mathbf{J})^{-1}$$

The sandwich form tells us that it is not the most efficient one. If we choose  $\mathbf{J} = (\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}'\mathbf{X}$  the asymptotic covariance becomes:

$$\begin{aligned} \sigma^2 \text{plim} & \left( \frac{1}{n}(\mathbf{X}'\mathbf{W}(\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}(\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}' \right. \\ & \cdot \left. \mathbf{W}(\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}'\mathbf{X}(\mathbf{X}'\mathbf{W}(\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}'\mathbf{X})^{-1} \right) \\ & = \sigma^2 \text{plim} \left( \frac{1}{n}\mathbf{X}'\mathbf{W}(\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}'\mathbf{X} \right)^{-1} = \sigma^2 \text{plim} \left( \frac{1}{n}\mathbf{X}'\mathbf{P}_W\mathbf{X} \right)^{-1} \end{aligned}$$

Thus our generalized IV estimator is

$$\hat{\boldsymbol{\beta}}_{GIVE} = (\mathbf{X}'\mathbf{P}_W\mathbf{X})^{-1}\mathbf{X}'\mathbf{P}_W\mathbf{y} \quad (6)$$

The form of the optimal  $\mathbf{J}$  is not surprising.  $\mathbf{W}\mathbf{J}$  are the orthogonal projections of  $\mathbf{X}$  onto  $S(\mathbf{W})$ . Thus each column of  $\mathbf{W}\mathbf{J}$  has shortest distance to the corresponding column of  $\mathbf{X}$ .

From equation (6) it is quite obvious that  $\hat{\boldsymbol{\beta}}_{GIVE}$  is also called the 2SLS estimator. It is important to note that  $\hat{\boldsymbol{\beta}}_{GIVE}$  is also the solution to  $\min Q(\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'\mathbf{P}_W(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$  with respect to  $\boldsymbol{\beta}$ .

### 4.3 Generalized Method of Moments

In the linear regression models we have the unconditional expectation

$$E(u_t) = 0$$

which we replace by the sample counterpart

$$\frac{1}{n} \sum_{t=1}^n u_t = \frac{1}{n} \sum_{t=1}^n (y_t - \mathbf{X}_t\boldsymbol{\beta}) = 0$$

We have one equation with  $k$  unknown parameters.

Introducing  $k$  instruments fulfilling the predetermined condition we get

$$\frac{1}{n} \sum \mathbf{W}'_t (y_t - \mathbf{X}_t\boldsymbol{\beta}) = \frac{1}{n} \mathbf{W}'\mathbf{y} - \mathbf{W}'\mathbf{X}\boldsymbol{\beta} \stackrel{!}{=} 0$$

In this case, the moment condition is linear in  $\beta$ . The GMM allows  $u_t = u_t(\beta)$  to be nonlinear.

By the same token, we need at least  $k$  instruments in  $\mathbf{W}_t$  to solve

$$\frac{1}{n} \sum \mathbf{W}_t \cdot \mathbf{u}_t(\beta) \stackrel{!}{=} \mathbf{0}$$

Or more generally we allow for the just identified case

$$\mathbf{g}_n(\beta) = \frac{1}{n} \sum \mathbf{f}(y_t, \mathbf{X}_t, \mathbf{W}_t, \beta) \stackrel{!}{=} \mathbf{0}$$

If we have more instruments than parameters the GMM estimator minimizes the function

$$Q(\beta) = \mathbf{g}_n(\beta)' \Psi \mathbf{g}_n(\beta) \tag{7}$$

with weighting matrix  $\Psi$ .

The first order conditions of minimizing (7) turn out to be

$$\mathbf{G}' \Psi \mathbf{g} \stackrel{!}{=} \mathbf{0} \quad \text{with } \mathbf{G} = \frac{\partial \mathbf{g}}{\partial \beta'}$$

Since  $E(\mathbf{f}) = \mathbf{0}$  the covariance matrix  $V(\mathbf{f}) = E(\mathbf{f} \mathbf{f}')$ . Thus the GMM estimator has asymptotic covariance matrix

$$\mathbf{G}' \Psi V(\mathbf{f}) \Psi' \mathbf{G}$$

Again, the sandwich form disappears if we choose  $\Psi = \Psi'$  optimally as  $\Psi_{opt} = V(\mathbf{f})^{-1}$ .



## 5 Maximum Likelihood

We consider now fully specified parametric models, i.e. we know the joint pdf

$$f(\mathbf{y}, \boldsymbol{\theta})$$

For given  $\boldsymbol{\theta}$  we obtain the pdf of vector  $\mathbf{y}$ . However, if we interpret  $f(\mathbf{y}, \boldsymbol{\theta})$  for a given vector  $\mathbf{y}$  it is now called a likelihood function for a given data set.

Stochastically independent  $y_t$  thus yield:

$$f(\mathbf{y}, \boldsymbol{\theta}) = \prod_{t=1}^n f(y_t, \boldsymbol{\theta})$$

The parameter vector  $\hat{\boldsymbol{\theta}}$  that maximizes  $f(\mathbf{y}, \boldsymbol{\theta})$  is called MLE.

For computational and analytical ease a strictly monotonic transformation of  $f(\mathbf{y}, \boldsymbol{\theta})$  is considered:

$$\ell(\mathbf{y}, \boldsymbol{\theta}) \equiv \log f(\mathbf{y}, \boldsymbol{\theta}) = \sum_{t=1}^n \ell_t(y_t, \boldsymbol{\theta})$$

$\ell_t(y_t, \boldsymbol{\theta})$  is called the  $t^{\text{th}}$  contribution to loglikelihood function.

### Example 1 Exponential Distribution

$$f(y_t, \theta) = \theta \exp(-\theta y_t) \quad y_t > 0; \theta > 0$$

$$\ell(\mathbf{y}, \theta) = \sum_{t=1}^n (\log \theta - \theta y_t) = n \log \theta - \theta \sum y_t$$

FOC:

$$\frac{d\ell}{d\theta} = \frac{n}{\theta} - \sum y_t = 0$$

$$\Leftrightarrow \hat{\theta} = \frac{n}{\sum y_t} = \frac{1}{\bar{y}}$$

### Example 2 Regression with normal errors

$$\mathbf{u} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

The elements  $u_t$  are independently distributed as  $N(0; \sigma^2)$ . Thus  $y_t$  conditional on  $\mathbf{X}_t$  has the same distribution law:

$$f_t(y_t, \boldsymbol{\beta}, \sigma | \mathbf{X}_t) = \frac{1}{\sigma} \phi \left( \frac{y_t - \mathbf{X}_t \boldsymbol{\beta}}{\sigma} \right) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left( -\frac{1}{2} \left( \frac{y_t - \mathbf{X}_t \boldsymbol{\beta}}{\sigma} \right)^2 \right)$$

$$\ell_t(y_t, \boldsymbol{\beta}, \sigma | \mathbf{X}_t) = -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (y_t - \mathbf{X}_t \boldsymbol{\beta})^2$$

$$\ell(\mathbf{y}, \boldsymbol{\beta}, \sigma | \mathbf{X}) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log \sigma^2 - \frac{n}{2\sigma^2} (\mathbf{y} - \mathbf{X} \boldsymbol{\beta})' (\mathbf{y} - \mathbf{X} \boldsymbol{\beta})$$

FOC:

$$\begin{aligned} \frac{\partial \ell}{\partial \boldsymbol{\beta}} &= -\frac{1}{2\sigma^2} (-2\mathbf{X}' \mathbf{y} + 2\mathbf{X}' \mathbf{X} \boldsymbol{\beta}) \stackrel{!}{=} \mathbf{0} \\ \frac{\partial \ell}{\partial \sigma^2} &= -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2(\sigma^2)^2} (\mathbf{y} - \mathbf{X} \boldsymbol{\beta})' (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}) \stackrel{!}{=} 0 \end{aligned}$$

The first set of equations obviously yields  $\hat{\boldsymbol{\beta}}_{ML} = \hat{\boldsymbol{\beta}}_{OLS}$ . Inserting this into the second equation obtains

$$\hat{\sigma}_{ML}^2 = \frac{\hat{\mathbf{u}}' \hat{\mathbf{u}}}{n}$$

As the examples showed most MLE can be characterized by

$$\left. \frac{\partial \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right|_{\hat{\boldsymbol{\theta}}} = \mathbf{0}$$

Since  $\ell(\boldsymbol{\theta}) = \sum \ell_t(\boldsymbol{\theta})$ , it follows that we can write the FOC:

$$\mathbf{g}(\boldsymbol{\theta}) = \frac{\partial \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum \frac{\partial \ell_t(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum \mathbf{g}_t(\boldsymbol{\theta}) \stackrel{!}{=} \mathbf{0}$$

$\mathbf{g}(\boldsymbol{\theta})$  is called gradient vector or score vector. The true parameter vector is denoted as  $\boldsymbol{\theta}_0$  which generated the data set  $\mathbf{y}$ . Thus

$$\begin{aligned} E(\mathbf{g}(\boldsymbol{\theta}_0)) &= E \left( \frac{\partial \ell(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \right) = E \left( \frac{\partial \log f(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \right) = E \left( \frac{1}{f(\boldsymbol{\theta}_0)} \cdot \frac{\partial f(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \right) \\ &= \int \frac{1}{f(\boldsymbol{\theta}_0)} \cdot \frac{\partial f(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \cdot f(\boldsymbol{\theta}_0) d\mathbf{y} = \int \frac{\partial f(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} d\mathbf{y} = \mathbf{0} \end{aligned}$$

This last equality follows from the fact that  $f$  is a pdf and thus  $\int f(\mathbf{y}, \boldsymbol{\theta}) d\mathbf{y} = 1$ .

Differentiating both sides yields

$$\frac{\partial}{\partial \boldsymbol{\theta}} \int f(\mathbf{y}, \boldsymbol{\theta}) d\mathbf{y} = \int \frac{\partial f(\mathbf{y}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} d\mathbf{y} = \mathbf{0}$$

Differentiating both sides of  $E(\mathbf{g}(\boldsymbol{\theta}_0)) = \mathbf{0}$  once again and using the fact that

$$\frac{\partial \ell(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} = \frac{\partial \log f(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} = \frac{1}{f(\boldsymbol{\theta}_0)} \cdot \frac{\partial f(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}}$$

we get

$$\begin{aligned} \int \frac{\partial^2 f(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} d\mathbf{y} = \mathbf{0} &\Leftrightarrow \int \frac{\partial^2 \ell(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \cdot f(\boldsymbol{\theta}_0) d\mathbf{y} + \int \frac{\partial \ell}{\partial \boldsymbol{\theta}} \cdot \frac{\partial f}{\partial \boldsymbol{\theta}} d\mathbf{y} = \mathbf{0} \\ &\Leftrightarrow -E \left( \frac{\partial^2 \ell(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \right) = E \left( \frac{\partial \ell(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \cdot \frac{\partial \ell(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}'} \right) \end{aligned}$$

The left hand side of this equation is called information matrix which is usually denoted as  $\mathbf{I}(\boldsymbol{\theta})$  and it is the covariance matrix of the score vector  $\mathbf{g}(\boldsymbol{\theta}_0)$ .

## Consistency

Due to Jensen's inequality we know that

$$E \log \left( \frac{f(\boldsymbol{\theta}^*)}{f(\boldsymbol{\theta}_0)} \right) < \log E \left( \frac{f(\boldsymbol{\theta}^*)}{f(\boldsymbol{\theta}_0)} \right)$$

As we have just seen

$$E \left( \frac{f(\boldsymbol{\theta}^*)}{f(\boldsymbol{\theta}_0)} \right) = \int \frac{f(\boldsymbol{\theta}^*)}{f(\boldsymbol{\theta}_0)} \cdot f(\boldsymbol{\theta}_0) d\mathbf{y} = \int f(\boldsymbol{\theta}^*) d\mathbf{y} = 1$$

Thus we get

$$E(\ell(\boldsymbol{\theta}^*)) - E(\ell(\boldsymbol{\theta}_0)) < 0 \Leftrightarrow E(\ell(\boldsymbol{\theta}^*)) < E(\ell(\boldsymbol{\theta}_0))$$

Applying a LLN to the contributions  $\ell_t$  we obtain  $\text{plim } n^{-1} \ell(\boldsymbol{\theta}) = \lim n^{-1} E(\ell(\boldsymbol{\theta}))$ :

$$\text{plim } \frac{1}{n} \ell(\boldsymbol{\theta}^*) \leq \text{plim } \ell(\boldsymbol{\theta}_0)$$

This holds for all  $\boldsymbol{\theta}^* \neq \boldsymbol{\theta}_0$ . The strict inequality disappears due to taking limits.

Since the MLE  $\hat{\boldsymbol{\theta}}$  maximizes  $\ell(\boldsymbol{\theta})$  it must hold that

$$\text{plim } \frac{1}{n} \ell(\hat{\boldsymbol{\theta}}) \geq \text{plim } \ell(\boldsymbol{\theta}_0)$$

## Asymptotic normality of MLE

Performing a Taylor expansion of the likelihood equation  $\mathbf{g}(\hat{\boldsymbol{\theta}}) = \mathbf{0}$  around  $\boldsymbol{\theta}_0$  yields

$$\mathbf{g}(\hat{\boldsymbol{\theta}}) = \mathbf{g}(\boldsymbol{\theta}_0) + \mathbf{H}(\bar{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) = \mathbf{0} \quad (8)$$

where  $\mathbf{H}(\boldsymbol{\theta}) = \partial^2 \ell / \partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'$  is the Hessian of the loglikelihood function. Since  $\bar{\boldsymbol{\theta}}$  must lie somewhere between  $\hat{\boldsymbol{\theta}}$  and  $\boldsymbol{\theta}_0$  we may write

$$\|\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| \leq \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\|$$

Solving equation (8) and inserting appropriate factors of  $n$  we get:

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) = - \left( \frac{1}{n} \mathbf{H}(\bar{\boldsymbol{\theta}}) \right)^{-1} \frac{1}{\sqrt{n}} \mathbf{g}(\boldsymbol{\theta}_0)$$

Components of the score vector  $\mathbf{g}_i(\boldsymbol{\theta}_0) = \sum_{t=1}^n \frac{\partial \ell_t(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}_i}$  having mean of 0 multiplied with  $n^{-1/2}$  is asymptotically normally distributed according to CLT.

Since  $\hat{\boldsymbol{\theta}}$  is consistent for  $\boldsymbol{\theta}_0$ , vector  $\bar{\boldsymbol{\theta}}$  is also consistent.

Thus  $\text{plim } n^{-1} \mathbf{H}(\bar{\boldsymbol{\theta}}) = \text{plim } n^{-1} \mathbf{H}(\hat{\boldsymbol{\theta}}) = \text{plim } n^{-1} \mathbf{H}(\boldsymbol{\theta}_0) = \mathcal{H}(\boldsymbol{\theta}_0)$ . Therefore, it follows with the help of  $-\mathcal{H}(\boldsymbol{\theta}_0) = \mathcal{I}(\boldsymbol{\theta}_0) = \text{plim } n^{-1} \mathbf{I}(\boldsymbol{\theta}_0)$

$$\frac{1}{\sqrt{n}} \mathbf{g}(\boldsymbol{\theta}_0) \stackrel{a}{\sim} N(\mathbf{0}; \mathcal{I}(\boldsymbol{\theta}_0))$$

Thus

$$\sqrt{n} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \stackrel{a}{\sim} N(\mathbf{0}; \underbrace{\mathcal{H}(\boldsymbol{\theta}_0)^{-1} \mathcal{I}(\boldsymbol{\theta}_0) \mathcal{H}(\boldsymbol{\theta}_0)^{-1}}_{=\mathcal{I}(\boldsymbol{\theta}_0)^{-1}})$$

## Computational Issues

Solving a system of nonlinear equations

$$\mathbf{g}(\hat{\boldsymbol{\theta}}) = \mathbf{0}$$

involves iterative maximization procedure like Newton-Raphson

$$\boldsymbol{\theta}^{(j+1)} = \boldsymbol{\theta}^{(j)} - \mathbf{H}(\boldsymbol{\theta}^{(j)})^{-1} \cdot \mathbf{g}(\boldsymbol{\theta}^{(j)})$$

It involves the Hessian which can be replaced by approximations to it. One approach is to use the information equality and use  $\mathbf{D}$  instead of  $\mathbf{H}$  with

$$\mathbf{D} = \sum \frac{\partial \ell_t}{\partial \boldsymbol{\theta}} \cdot \left( \frac{\partial \ell_t}{\partial \boldsymbol{\theta}} \right)'$$

The inverse of this matrix or of  $-\mathbf{H}(\hat{\boldsymbol{\theta}})$  or of  $\mathbf{I}(\hat{\boldsymbol{\theta}})$  are all common estimates for the asymptotic covariance matrix of  $\hat{\boldsymbol{\theta}}$ .

### Delta Theorem

Suppose a  $k$ -vector  $\boldsymbol{\theta}$  with  $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \stackrel{a}{\sim} N(\mathbf{0}; V(\hat{\boldsymbol{\theta}}))$  and the  $l$ -vector ( $l < k$ )  $\boldsymbol{\gamma} = \mathbf{g}(\boldsymbol{\theta})$  with  $\mathbf{g}(\cdot)$  having  $l$  monotonic and differentiable functions  $g_i$ , then

$$\sqrt{n}(\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}_0) \stackrel{a}{\sim} N(\mathbf{0}; \mathbf{G}_0 V(\hat{\boldsymbol{\theta}}) \mathbf{G}_0')$$

with  $\mathbf{G}_0$  a  $l \times k$  matrix with typical element  $\partial g_i / \partial \theta_j$  evaluated at  $\boldsymbol{\theta}_0$ . In practice we estimate this covariance matrix by  $\mathbf{G}(\hat{\boldsymbol{\theta}}) \widehat{V}(\hat{\boldsymbol{\theta}}) \mathbf{G}(\hat{\boldsymbol{\theta}})'$ .

## 5.1 Hypothesis tests

Testing hypothesis of type  $\mathbf{r}(\boldsymbol{\theta}) = \mathbf{0}$  having  $r$  components we have three asymptotically identical tests:

### Likelihood ratio test (LR)

$$LR = 2(\ell(\hat{\boldsymbol{\theta}}) - \ell(\tilde{\boldsymbol{\theta}}))$$

where  $\tilde{\boldsymbol{\theta}}$  denotes the restricted MLE fulfilling  $\mathbf{r}(\tilde{\boldsymbol{\theta}}) = \mathbf{0}$ . Taylor expansion of  $\ell(\tilde{\boldsymbol{\theta}})$  around  $\hat{\boldsymbol{\theta}}$  obtains

$$\begin{aligned} \ell(\tilde{\boldsymbol{\theta}}) &= \ell(\hat{\boldsymbol{\theta}}) + \mathbf{g}(\hat{\boldsymbol{\theta}})'(\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}) + \frac{1}{2}(\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})' \mathbf{H}(\bar{\boldsymbol{\theta}})(\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}) \\ &\Leftrightarrow 2(\ell(\hat{\boldsymbol{\theta}}) - \ell(\tilde{\boldsymbol{\theta}})) = -(\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})' \mathbf{H}(\bar{\boldsymbol{\theta}})(\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}) \end{aligned}$$

Intuitively, it shows that LR is approximately  $\chi_r^2$  distributed.

### Wald test

$$W = \mathbf{r}(\hat{\boldsymbol{\theta}})'[\mathbf{R}\mathbf{I}(\hat{\boldsymbol{\theta}})^{-1}\mathbf{R}']^{-1}\mathbf{r}(\hat{\boldsymbol{\theta}})$$

with  $r \times k$  matrix  $\mathbf{R}$  having typical element  $\partial \mathbf{r}_i / \partial \theta_j$ . Since  $\hat{\boldsymbol{\theta}}$  is asymptotically normally distributed and under  $H_0 : \mathbf{r}(\boldsymbol{\theta}) = \mathbf{0}$  we can apply the Delta Theorem and get

$$\mathbf{r}(\hat{\boldsymbol{\theta}}) \stackrel{a}{\sim} N(\mathbf{0}; \mathbf{R}\mathbf{I}(\hat{\boldsymbol{\theta}})^{-1}\mathbf{R}')$$

The asymptotic distribution of the quadratic form is thus  $\chi_r^2$ .

### Lagrange multiplier test (LM)

$$LM = \mathbf{g}(\tilde{\boldsymbol{\theta}})' \mathbf{I}(\tilde{\boldsymbol{\theta}})^{-1} \mathbf{g}(\tilde{\boldsymbol{\theta}})$$

The FOC of the restricted optimization problem are

$$\mathbf{g}(\tilde{\boldsymbol{\theta}}) + \mathbf{R}(\tilde{\boldsymbol{\theta}})' \boldsymbol{\lambda} = \mathbf{0}$$

$$\mathbf{r}(\tilde{\boldsymbol{\theta}}) = \mathbf{0}$$

Under the null hypothesis  $\mathbf{r}(\boldsymbol{\theta}_0) = \mathbf{0}$  we would get  $\boldsymbol{\lambda} = \mathbf{0}$  implying  $\mathbf{g}(\tilde{\boldsymbol{\theta}}) = \mathbf{0}$ . Thus the quadratic form  $\mathbf{g}(\tilde{\boldsymbol{\theta}})' \mathbf{I}(\tilde{\boldsymbol{\theta}})^{-1} \mathbf{g}(\tilde{\boldsymbol{\theta}})$  should be asymptotically  $\chi_r^2$  distributed.

## 6 Generalized Least Squares

We now assume a general linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u} \quad ; E(\mathbf{u}) = \mathbf{0} \quad , \quad E(\mathbf{u}\mathbf{u}') = \boldsymbol{\Omega} \quad (9)$$

Clearly, if  $\boldsymbol{\Omega} = \sigma^2\mathbf{I}$  we are back to the classical model.

First we can note that  $\hat{\boldsymbol{\beta}}_{OLS}$  is still unbiased given that either  $\mathbf{X}$  is fix or uncorrelated with  $\mathbf{u}$ :

$$E(\hat{\boldsymbol{\beta}}_{OLS}) = \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\mathbf{u})$$

However, its covariance matrix now is

$$V(\hat{\boldsymbol{\beta}}_{OLS}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\mathbf{u}\mathbf{u}')\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\Omega}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$$

As a consequence, inference based on the covariance matrix  $\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$  is not valid anymore.

### 6.1 GLS estimator

The covariance matrix  $\boldsymbol{\Omega}$  is a positive definite matrix implying that its inverse can be decomposed

$$\boldsymbol{\Omega}^{-1} = \boldsymbol{\Psi}\boldsymbol{\Psi}' \quad \Rightarrow \quad \boldsymbol{\Omega} = (\boldsymbol{\Psi}')^{-1}\boldsymbol{\Psi}^{-1}$$

Premultiplying equation (9) with  $\boldsymbol{\Psi}'$  yields

$$\boldsymbol{\Psi}'\mathbf{y} = \boldsymbol{\Psi}'\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\Psi}'\mathbf{u}$$

Analyzing the transformed error terms we see

$$V(\boldsymbol{\Psi}'\mathbf{u}) = \boldsymbol{\Psi}'V(\mathbf{u})\boldsymbol{\Psi} = \boldsymbol{\Psi}'\boldsymbol{\Omega}\boldsymbol{\Psi} = \boldsymbol{\Psi}'(\boldsymbol{\Psi}')^{-1}\boldsymbol{\Psi}^{-1}\boldsymbol{\Psi} = \mathbf{I}$$

Thus we have error terms fulfilling the classical assumptions. For this model we know due to Gauss-Markov theorem the most efficient estimator

$$\hat{\boldsymbol{\beta}}_{GLS} = (\mathbf{X}'\boldsymbol{\Psi}\boldsymbol{\Psi}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\Psi}\boldsymbol{\Psi}'\mathbf{y} = (\mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{y}$$

The covariance matrix of this estimator follows immediately

$$V(\hat{\beta}_{GLS}) = (\mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{X})^{-1}$$

This estimator can also be obtained if we minimize the criterion function

$$(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'\boldsymbol{\Omega}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

which is the sum of squared residuals from the transformed model. The first order conditions for this optimization problem are

$$\mathbf{X}'\boldsymbol{\Omega}^{-1}(\mathbf{y} - \mathbf{X}\hat{\beta}_{GLS}) = \mathbf{0}$$

If we compare it to the instrumental variables estimator (or more general to GMM) we see that setting  $\mathbf{W} = \boldsymbol{\Omega}^{-1}\mathbf{X}$  the moment conditions

$$\mathbf{W}'(\mathbf{y} - \mathbf{X}\hat{\beta}_{\mathbf{W}}) = \mathbf{0}$$

follow. From that the method of moments yields the estimator

$$\hat{\beta}_{\mathbf{W}} = (\mathbf{W}'\mathbf{X})^{-1}\mathbf{W}'\mathbf{y} = \boldsymbol{\beta} + (\mathbf{W}'\mathbf{X})^{-1}\mathbf{W}'\mathbf{u}$$

Therefore, the covariance matrix of  $\hat{\beta}_{\mathbf{W}}$  is

$$V(\hat{\beta}_{\mathbf{W}}) = (\mathbf{W}'\mathbf{X})^{-1}\mathbf{W}'\boldsymbol{\Omega}\mathbf{W}(\mathbf{X}'\mathbf{W})^{-1}$$

Again, the sandwich form indicates that it is not efficient. Note, that setting  $\mathbf{W} = \mathbf{X}$  we are back to the above analyzed OLS estimator.

To prove efficiency of  $\hat{\beta}_{GLS}$  we first show that  $\mathbf{A} - \mathbf{B}$  is positive definite if and only if  $\mathbf{B}^{-1} - \mathbf{A}^{-1}$  is positive definite. To see that we start with a symmetric positive definite matrix  $\mathbf{A}$  and a full column rank matrix  $\mathbf{C}$ . It follows that  $\mathbf{C}'\mathbf{A}\mathbf{C}$  is positive definite:

$$\mathbf{x}'\mathbf{C}'\mathbf{A}\mathbf{C}\mathbf{x} = (\mathbf{x}'\mathbf{C}')\mathbf{A}(\mathbf{C}\mathbf{x}) = (\mathbf{C}\mathbf{x})'\mathbf{A}(\mathbf{C}\mathbf{x}) = \mathbf{w}'\mathbf{A}\mathbf{w} > 0 \quad \forall \mathbf{x} \neq \mathbf{0}$$

Note that full column rank implies  $\mathbf{0} = \mathbf{w} = \mathbf{C}\mathbf{x} \Rightarrow \mathbf{x} = \mathbf{0}$ . Additionally, we need the fact that if  $\mathbf{A}$  is positive definite then  $\mathbf{A}^{-1}$  is positive definite as well. As before,



we write  $\mathbf{A}^{-1}$  as  $\mathbf{A}^{-1} = \mathbf{\Gamma}\mathbf{\Gamma}'$  implying  $\mathbf{\Gamma}'\mathbf{A}\mathbf{\Gamma} = \mathbf{\Gamma}'(\mathbf{\Gamma}')^{-1}\mathbf{\Gamma}^{-1}\mathbf{\Gamma} = \mathbf{I}$ . Now we consider  $(\mathbf{I} - \mathbf{A})$  which is positive definite if and only if  $(\mathbf{A}^{-1} - \mathbf{I})$  is positive definite. To see that we pre- and postmultiply  $(\mathbf{I} - \mathbf{A})$  with  $\mathbf{\Gamma}'$  and  $\mathbf{\Gamma}$ , respectively:

$$\mathbf{\Gamma}'(\mathbf{I} - \mathbf{A})\mathbf{\Gamma} = \mathbf{\Gamma}'\mathbf{\Gamma} - \mathbf{\Gamma}'\mathbf{A}\mathbf{\Gamma} = \mathbf{\Gamma}'\mathbf{\Gamma} - \mathbf{I}$$

This matrix has the same eigenvalues as  $\mathbf{\Gamma}\mathbf{\Gamma}' - \mathbf{I} = \mathbf{A}^{-1} - \mathbf{I}$ :

$$\begin{aligned} (\mathbf{A}^{-1} - \mathbf{I})\mathbf{x} = \lambda\mathbf{x} \quad & \text{premult. } \mathbf{\Gamma}' \quad \Leftrightarrow \quad (\mathbf{\Gamma}'\mathbf{\Gamma}\mathbf{\Gamma}' - \mathbf{\Gamma}')\mathbf{x} = \lambda(\mathbf{\Gamma}'\mathbf{x}) \\ \Leftrightarrow \quad & (\mathbf{\Gamma}'\mathbf{\Gamma} - \mathbf{I})(\mathbf{\Gamma}'\mathbf{x}) = \lambda(\mathbf{\Gamma}'\mathbf{x}) \end{aligned}$$

Now we turn to  $\mathbf{A}$  and  $\mathbf{B}$  both being symmetric and positive definite. If  $\mathbf{A} - \mathbf{B}$  is positive definite then

$$\mathbf{\Gamma}'(\mathbf{A} - \mathbf{B})\mathbf{\Gamma} = \mathbf{\Gamma}'\mathbf{A}\mathbf{\Gamma} - \mathbf{\Gamma}'\mathbf{B}\mathbf{\Gamma} = \mathbf{I} - \mathbf{\Gamma}'\mathbf{B}\mathbf{\Gamma}$$

is positive definite and also  $\mathbf{\Gamma}^{-1}\mathbf{B}^{-1}(\mathbf{\Gamma}')^{-1} - \mathbf{I}$  which implies positive definiteness of

$$\mathbf{\Gamma}(\mathbf{\Gamma}^{-1}\mathbf{B}^{-1}(\mathbf{\Gamma}')^{-1} - \mathbf{I})\mathbf{\Gamma}' = \mathbf{B}^{-1} - \mathbf{\Gamma}\mathbf{\Gamma}' = \mathbf{B}^{-1} - \mathbf{A}^{-1}$$

We now consider the difference between the covariance matrices of  $\hat{\boldsymbol{\beta}}_{GLS}$  and  $\hat{\boldsymbol{\beta}}_W$ :

$$V(\hat{\boldsymbol{\beta}}_W) - V(\hat{\boldsymbol{\beta}}_{GLS}) = (\mathbf{W}'\mathbf{X})^{-1}\mathbf{W}'\mathbf{\Omega}\mathbf{W}(\mathbf{X}'\mathbf{W})^{-1} - (\mathbf{X}'\mathbf{\Omega}^{-1}\mathbf{X})^{-1}$$

for which we want to show positive semi-definiteness. Thus we consider

$$\begin{aligned} \mathbf{X}'\mathbf{\Omega}^{-1}\mathbf{X} - \mathbf{X}'\mathbf{W}(\mathbf{W}'\mathbf{\Omega}\mathbf{W})^{-1}\mathbf{W}'\mathbf{X} &= \mathbf{X}'(\mathbf{\Psi}\mathbf{\Psi}' - \mathbf{W}(\mathbf{W}'(\mathbf{\Psi}')^{-1}\mathbf{\Psi}^{-1}\mathbf{W})^{-1}\mathbf{W}')\mathbf{X} \\ &= \mathbf{X}'\mathbf{\Psi}(\mathbf{I} - \mathbf{\Psi}^{-1}\mathbf{W}(\mathbf{W}'(\mathbf{\Psi}')^{-1}\mathbf{\Psi}^{-1}\mathbf{W})^{-1}\mathbf{W}'(\mathbf{\Psi}')^{-1})\mathbf{\Psi}'\mathbf{X} \\ &= \mathbf{X}'\mathbf{\Psi}\mathbf{M}_{\mathbf{\Psi}^{-1}\mathbf{W}}\mathbf{\Psi}'\mathbf{X} \end{aligned}$$

Since  $\mathbf{M}$  is idempotent the positive semi-definiteness follows immediately showing the efficiency of  $\hat{\boldsymbol{\beta}}_{GLS}$ .

## 6.2 Feasible GLS

The consistency and normality of  $\hat{\boldsymbol{\beta}}_{GLS}$  follows immediately from the interpretation as an instrumental variables estimator (MM or GMM).

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_{GLS} - \boldsymbol{\beta}) \stackrel{a}{=} (\text{plim } \frac{1}{n}\mathbf{X}'\mathbf{\Omega}^{-1}\mathbf{X})^{-1}(\text{plim } n^{-1/2}\mathbf{X}'\mathbf{\Omega}^{-1}\mathbf{u})$$

If we assume that  $\mathbf{\Omega}$  depends on an  $l$ -dimensional vector  $\boldsymbol{\gamma}$  such that

$$E(u_t^2) = \exp(\mathbf{Z}_t \boldsymbol{\gamma})$$

we obtain  $\mathbf{\Omega}(\boldsymbol{\gamma})$  which might be estimated by inserting an estimate of  $\boldsymbol{\gamma}$ :  $\hat{\mathbf{\Omega}} = \mathbf{\Omega}(\hat{\boldsymbol{\gamma}})$ .

To obtain a consistent estimate of  $\boldsymbol{\gamma}$  we first need a consistent estimate of  $u_t$  which is obviously possible if we initially use OLS estimate  $\hat{\boldsymbol{\beta}}_{OLS}$  and calculate the OLS residuals and then estimate

$$\log(u_t^2) = \mathbf{Z}_t \boldsymbol{\gamma} + \nu_t$$

A typical element of  $\mathbf{\Omega}(\hat{\boldsymbol{\gamma}})$  then is

$$\hat{\omega}_t = (\exp(\mathbf{Z}_t \hat{\boldsymbol{\gamma}}))^{1/2}$$

Thus for the feasible GLS estimator  $\hat{\boldsymbol{\beta}}_{FGLS} = (\mathbf{X}' \hat{\mathbf{\Omega}}^{-1} \mathbf{X})^{-1} \mathbf{X}' \hat{\mathbf{\Omega}}^{-1} \mathbf{y}$  we see that

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_{FGLS} - \boldsymbol{\beta}) \stackrel{a}{=} (\text{plim } \frac{1}{n} \mathbf{X}' \hat{\mathbf{\Omega}}^{-1} \mathbf{X})^{-1} (\text{plim } n^{-1/2} \mathbf{X}' \hat{\mathbf{\Omega}}^{-1} \mathbf{u})$$

If we have

$$\text{plim } \frac{1}{n} \mathbf{X}' \hat{\mathbf{\Omega}}^{-1} \mathbf{X} = \text{plim } \frac{1}{n} \mathbf{X}' \mathbf{\Omega}^{-1} \mathbf{X}$$

and

$$\text{plim } n^{-1/2} \mathbf{X}' \hat{\mathbf{\Omega}}^{-1} \mathbf{u} = \text{plim } n^{-1/2} \mathbf{X}' \mathbf{\Omega}^{-1} \mathbf{u}$$

then  $\hat{\boldsymbol{\beta}}_{FGLS}$  and  $\hat{\boldsymbol{\beta}}_{GLS}$  share the same asymptotic properties.

The consistency of  $\hat{\boldsymbol{\beta}}_{FGLS}$  and  $\hat{\boldsymbol{\beta}}_{GLS}$  hinges on consistent residual estimates which in turn require consistent OLS estimates. If  $\mathbf{\Omega}$  is not diagonal the consistency of  $\hat{\boldsymbol{\beta}}_{GLS}$  is lost if any element of  $\mathbf{X}_t$  is a lagged dependent variable. The lagged dependent variables are still predetermined to current  $u_t$  but due to serial correlation of the errors they are not uncorrelated to  $u_t$ .

### 6.3 Heteroscedasticity

As seen above, if we assume for the main diagonal of  $\mathbf{\Omega}$

$$\omega_t^2 = \exp(\delta + \mathbf{Z}_t \boldsymbol{\gamma})$$

or more general

$$\omega_t^2 = h(\delta + \mathbf{Z}_t\boldsymbol{\gamma})$$

and off-main diagonal elements equal to zero we obtain as a typical element of  $\boldsymbol{\Omega}^{-1}$

$$\omega_t^{-2} = \frac{1}{h(\delta + \mathbf{Z}_t\boldsymbol{\gamma})}$$

and therefore a typical element of  $\boldsymbol{\Psi}$  as

$$\omega_t^{-1} = \left( \frac{1}{h(\delta + \mathbf{Z}_t\boldsymbol{\gamma})} \right)^{1/2}$$

The transformed model then looks like

$$\omega_t^{-1}y_t = \omega_t^{-1}\mathbf{X}_t\boldsymbol{\beta} + \omega_t^{-1}u_t$$

This is the reason why it is sometimes called *weighted least squares* (WLS). Note that if  $\mathbf{X}_t$  contained a constant, after transformation it becomes to  $1/\omega_t$ .

### 6.3.1 Testing for Heteroscedasticity

Assuming  $E(u_t^2|\mathbf{Z}_t) = h(\delta + \mathbf{Z}_t\boldsymbol{\gamma})$  we can write this as

$$u_t^2 = h(\delta + \mathbf{Z}_t\boldsymbol{\gamma}) + \text{residual}$$

This is a non-linear regression for which we get

$$\begin{aligned} h'(\hat{\delta} + \mathbf{Z}_t\hat{\boldsymbol{\gamma}}) \left( u_t^2 - h(\hat{\delta} + \mathbf{Z}_t\hat{\boldsymbol{\gamma}}) \right) &= 0 \\ h'(\hat{\delta} + \mathbf{Z}_t\hat{\boldsymbol{\gamma}})\mathbf{Z}_t' \left( u_t^2 - h(\hat{\delta} + \mathbf{Z}_t\hat{\boldsymbol{\gamma}}) \right) &= \mathbf{0} \end{aligned}$$

Thus an artificial regression of  $u_t^2 - h(\hat{\delta} + \mathbf{Z}_t\hat{\boldsymbol{\gamma}})$  on  $h'(\hat{\delta} + \mathbf{Z}_t\hat{\boldsymbol{\gamma}})$  and  $h'(\hat{\delta} + \mathbf{Z}_t\hat{\boldsymbol{\gamma}})\mathbf{Z}_t$  should result in coefficients equal to zero:

$$u_t^2 - h(\hat{\delta} + \mathbf{Z}_t\hat{\boldsymbol{\gamma}}) = h'(\hat{\delta} + \mathbf{Z}_t\hat{\boldsymbol{\gamma}})b_\delta + h'(\hat{\delta} + \mathbf{Z}_t\hat{\boldsymbol{\gamma}})\mathbf{Z}_t\mathbf{b}_\gamma + \text{residual}$$

Under the null hypothesis of no heteroscedasticity, i.e.  $H_0 : \boldsymbol{\gamma} = \mathbf{0}$ , this regression reduces to

$$u_t^2 - h(\delta) = h'(\delta)b_\delta + h'(\delta)\mathbf{Z}_t\mathbf{b}_\gamma + \text{residual}$$

The constant factor  $h'(\delta)$  changes the coefficients, they do not change the explanatory power (also SSR) of the model. Thus a test for heteroscedasticity then simply is running the regression

$$\hat{u}_t^2 = b_\delta + \mathbf{Z}_t \mathbf{b}_\gamma + \text{residual}$$

and performing an F-test for  $\mathbf{b}_\gamma = \mathbf{0}$ . Note that it does not depend on the form of  $h(\cdot)$ .

### 6.3.2 Heteroscedasticity-Consistent Covariance Matrix (HCCM)

Using  $\hat{\boldsymbol{\beta}}_{OLS}$  we saw above that for  $\mathbf{X}$  considered as fix and given  $\boldsymbol{\Omega}$

$$V(\hat{\boldsymbol{\beta}}_{OLS}) = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\boldsymbol{\Omega}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$$

Looking at  $\sqrt{n}(\hat{\boldsymbol{\beta}}_{OLS} - \boldsymbol{\beta})$  we find that the asymptotic covariance matrix to be

$$\left( \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}'\mathbf{X} \right)^{-1} \left( \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}'\boldsymbol{\Omega}\mathbf{X} \right) \left( \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}'\mathbf{X} \right)^{-1} \quad (10)$$

In one of the probably most famous papers in econometrics White (1980) showed that the middle term can be consistently estimated by

$$\frac{1}{n} \sum_{t=1}^n \hat{u}_t^2 \mathbf{X}'_t \mathbf{X}_t$$

Thus an estimate for the asymptotic HCCM in (10) is

$$\widehat{V(\hat{\boldsymbol{\beta}}_{OLS})} = (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{t=1}^n \hat{u}_t^2 \mathbf{X}'_t \mathbf{X}_t \right) (\mathbf{X}'\mathbf{X})^{-1}$$

## 6.4 Autoregressive Processes

In regression models for time series data it is an often encountered problem that due to various reasons error terms might be correlated.

### 6.4.1 AR(1) Process

The simplest autoregressive structure is the first order process written as

$$u_t = \rho u_{t-1} + \varepsilon_t \quad \varepsilon_t \sim IID(0, \sigma^2) \quad |\rho| < 1 \quad (11)$$

The requirement  $|\rho| < 1$  is called **stationarity condition**. Weak stationarity of process  $u_t$  requires

- $E(u_t) = \mu$
- $\text{Var}(u_t) = \sigma_u^2$
- $\text{Cov}(u_t, u_{t-s}) = \gamma_s$

All these moments, of course they should exist, are independent of  $t$ . The covariance only depends on the number of lags considered but is also independent of  $t$ .

Substituting recursively  $u_{t-1} = \rho u_{t-2} + \varepsilon_{t-1}$  in (14) again and again yields

$$u_t = \varepsilon_t + \rho \varepsilon_{t-1} + \rho^2 \varepsilon_{t-2} + \rho^3 \varepsilon_{t-3} + \rho^4 \varepsilon_{t-4} + \dots$$

Since the  $\varepsilon_i$  are all independent it is easy to compute the variance of error term  $u_t$ :

$$\sigma_u^2 \equiv \text{Var}(u_t) = \sum_{i=0}^{\infty} (\rho^i)^2 \sigma_\varepsilon^2 = \sigma_\varepsilon^2 \sum_{i=0}^{\infty} (\rho^2)^i = \frac{\sigma_\varepsilon^2}{1 - \rho^2}$$

The requirement for the geometric series to converges is  $|\rho^2| < 1$  which implies  $|\rho| < 1$ . The covariance of  $u_t$  and  $u_{t-1}$  is

$$\text{Cov}(u_t, u_{t-1}) = E(u_t \cdot u_{t-1}) = E((\rho u_{t-1} + \varepsilon_t) \cdot u_{t-1}) = E(\rho u_{t-1}^2) = \rho \sigma_u^2$$

The covariance of  $u_t$  and  $u_{t-1}$  is

$$\text{Cov}(u_t, u_{t-2}) = E(u_t \cdot u_{t-2}) = E((\rho^2 u_{t-2} + \varepsilon_t + \rho \varepsilon_{t-1}) \cdot u_{t-2}) = E(\rho^2 u_{t-2}^2) = \rho^2 \sigma_u^2$$

Or more general:  $\gamma_s = \text{Cov}(u_t, u_{t-s}) = \rho^s \sigma_u^2$

As a result, the covariance matrix  $\mathbf{\Omega}$  of  $n$ -vector  $\mathbf{u}$  is

$$\mathbf{\Omega}(\rho) = \frac{\sigma_\varepsilon^2}{1 - \rho^2} \begin{pmatrix} 1 & \rho & \rho^2 & \dots & \rho^{n-1} \\ \rho & 1 & \rho & \dots & \rho^{n-2} \\ \vdots & \vdots & \vdots & & \vdots \\ \rho^{n-1} & \rho^{n-2} & \rho^{n-3} & \dots & 1 \end{pmatrix} \quad (12)$$

### 6.4.2 Higher-Order Autoregressive Processes

We consider an AR(p) process

$$u_t = \alpha_1 u_{t-1} + \alpha_2 u_{t-2} + \dots + \alpha_p u_{t-p} + \varepsilon_t \quad \varepsilon_t \sim \text{IID}(0, \sigma_\varepsilon^2)$$

For these higher order processes it is useful to introduce the lag-operator  $L$  which has the properties  $Lu_t = u_{t-1}$  and  $LLu_t = L^2 u_t = u_{t-2}$ . The AR(p) process can now be written as

$$(1 - \alpha_1 L - \alpha_2 L^2 - \dots - \alpha_p L^p)u_t = \varepsilon_t$$

Note, that the AR(1) process with  $(1 - \rho L)u_t = \varepsilon_t$  can be **inverted** to

$$u_t = (1 - \rho L)^{-1} \varepsilon_t$$

Thus the term  $(1 - \rho L)^{-1}$  represents a geometric series of  $(\rho L)$ .

The notation using the lag-operator is helpful since it closely relates to the **characteristic equation**

$$1 - \alpha_1 z - \alpha_2 z^2 - \dots - \alpha_p z^p = 0$$

If the roots of this characteristic equation, some of them might be complex, all lie outside the unit circle then the AR(p) process is stationary. A complex number  $z = a + bi$  lies outside the unit circle if  $a^2 + b^2 > 1$ .

### 6.4.3 Testing for Serial Correlation

The most common test statistic to test for serial correlation is the Durbin-Watson DW-statistic:

$$d = \frac{\sum_{t=2}^n (\hat{u}_t - \hat{u}_{t-1})^2}{\sum_{t=1}^n \hat{u}_t^2} = \frac{\sum_{t=2}^n \hat{u}_t^2}{\sum_{t=1}^n \hat{u}_t^2} - 2 \cdot \frac{\sum_{t=2}^n \hat{u}_t \hat{u}_{t-1}}{\underbrace{\sum_{t=2}^n \hat{u}_t^2}_{\hat{\rho}}} + \frac{\sum_{t=2}^n \hat{u}_{t-1}^2}{\sum_{t=1}^n \hat{u}_t^2}$$

If we assume  $\hat{u}_0 = 0$  the middle term is the sample correlation coefficient  $\hat{\rho}$ . Thus for  $n \rightarrow \infty$  the first and the last term of the right hand side tend to 1 so that  $d$  is asymptotically equal to  $2 - 2\hat{\rho}$ .

As we have seen in section 2 the residual vector  $\hat{\mathbf{u}}$  of a correctly specified model is equal to  $\mathbf{M}_X \mathbf{u}$ . It implies that in finite samples even if the components in  $\mathbf{u}$  are serially uncorrelated the residuals  $\hat{\mathbf{u}}$  display a certain amount of serial correlation depending on matrix  $\mathbf{X}$ .

A critical value  $d$  for a given significance level thus depends on  $\mathbf{X}$ . But it has been shown that there are lower  $d_l$  and upper  $d_u$  bounds for it only depending on  $n$  and  $k$ . In practice, it is still quite common to reject  $H_0 : \rho = 0$  if  $d < d_l$ , not reject if  $d > d_u$ , and be indecisive if  $d_l < d < d_u$  although everyone could perform a Monte Carlo based test to overcome this problem.

### 6.4.4 Estimation

If we consider a model with AR(1) error term

$$y_t = \mathbf{X}_t \boldsymbol{\beta} + u_t \quad u_t = \rho u_{t-1} + \varepsilon_t \quad \varepsilon_t \sim IID(0, \sigma^2) \quad |\rho| < 1$$

we can replace  $u_{t-1} = y_{t-1} - \mathbf{X}_{t-1} \boldsymbol{\beta}$  and obtain

$$y_t = \rho y_{t-1} + \mathbf{X}_t \boldsymbol{\beta} - \rho \mathbf{X}_{t-1} \boldsymbol{\beta} + \varepsilon_t$$

If we ignore the first observation this equation can be estimated using non-linear least squares.

An alternative approach is the feasible GLS estimator. We just have to find a decomposition of covariance matrix (12) fulfilling  $\mathbf{\Omega}^{-1} = \mathbf{\Psi}\mathbf{\Psi}'$  in a way that  $\mathbf{\Psi}'\mathbf{u} = \boldsymbol{\varepsilon}$ . For observations  $t = 2, \dots, n$  we know that

$$\varepsilon_t = u_t - \rho u_{t-1}$$

Thus the  $t^{th}$  row of  $\mathbf{\Psi}'$  must have a 1 in the  $t^{th}$  position,  $-\rho$  in the  $(t-1)^{th}$  position, and 0 elsewhere. For the first row of  $\mathbf{\Psi}'\mathbf{u}$  we use the unconditional variance of  $u_1$  which from above is  $\text{Var}(u_1) = \sigma_\varepsilon^2/(1 - \rho^2)$ . Therefore, if we define

$$\varepsilon_1 = (1 - \rho^2)^{1/2} u_1$$

it is obvious that  $\varepsilon_1$  is independent of all  $\varepsilon_2, \varepsilon_3$  etc. and additionally we obtain  $\text{Var}(\varepsilon_1) = (1 - \rho^2)/(1 - \rho^2) = 1$ . As a result, we get

$$\mathbf{\Psi}' = \begin{pmatrix} (1 - \rho^2)^{1/2} & 0 & \dots & & 0 \\ -\rho & 1 & 0 & \dots & 0 \\ 0 & -\rho & 1 & 0 & \dots & 0 \\ \vdots & & & & & \vdots \\ 0 & \dots & 0 & -\rho & 1 & 0 \\ 0 & \dots & & 0 & -\rho & 1 \end{pmatrix}$$

The next step is calculating the transformed observations  $\mathbf{\Psi}'\mathbf{y}$  and  $\mathbf{\Psi}'\mathbf{X}$  and running an OLS regression for these transformed variables.

Whenever  $\rho$  is known this GLS estimator can be applied immediately. If it is unknown it can be estimated using  $\hat{\mathbf{u}}$  from a preceding OLS estimation.

The FGLS residual  $(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{FGLS})$  can be used to re-estimate  $\rho$  and reiterate the GLS step. By succeedingly repeating these steps we come quite close to the non-linear least squares estimates.

A third alternative estimator is MLE if one is willing to assume  $\varepsilon_t \sim IIN(0, \sigma_\varepsilon^2)$ .



## 6.5 Panel Data

Panel data have two dimensions, a cross-section dimension and a time dimension, and can be seen as repeated cross-sections. Our model then becomes to

$$y_{it} = \mathbf{X}_{it}\boldsymbol{\beta} + u_{it} \quad i = 1, \dots, m \quad t = 1, \dots, T$$

where  $\mathbf{X}_{it}$  is a  $1 \times k$  vector of explanatory variables. There are now  $n = mT$  observations.

- If  $E(u_{it}|\mathbf{X}_{it}) = 0$  OLS yields unbiased but probably inefficient estimates if  $\text{Cov}(u_{it}, u_{jt}) \neq 0$  or  $\text{Cov}(u_{it}, u_{i,t+s}) \neq 0$ .
- If  $\mathbf{X}_{it}$  contains lagged dependent variables, OLS estimates are inconsistent if  $u_{it}$  are serially correlated.

### Error Components

It is quite common to assume

$$u_{it} = e_t + \nu_i + \varepsilon_{it}$$

with  $e_t$  being independent across  $t$ ,  $\nu_i$  being independent across  $i$ , and  $\varepsilon_{it}$  being independent across  $i$  and  $t$ .

We will simplify by assuming  $e_t=0$ .

### Fixed-Effects

If we assume  $\nu_i$  to be non-stochastic or fixed but unknown to the researcher they can be treated as parameters to be estimated. In such a case we can write the model as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{D}\boldsymbol{\eta} + \boldsymbol{\varepsilon}$$

The covariance matrix for error  $\boldsymbol{\varepsilon}$  then is  $E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}') = \sigma_\varepsilon^2 \mathbf{I}_n$ . Vector  $\mathbf{y}$  has  $y_{it}$  as typical element and  $\mathbf{D}$  is an  $n \times m$  matrix of dummy variables. A 1 in column  $j$  indicates

that this specific observation in that row belongs to cross-sectional unit  $j$ . There are  $T$  ones in each column of  $\mathbf{D}$ . The  $m$ -dimensional vector  $\boldsymbol{\eta}$  contains the fixed effects  $\nu_i$ .

The **fixed-effects** estimator of  $\boldsymbol{\beta}$ , sometimes called least squares dummy variables (LSDV) estimator, then is

$$\hat{\boldsymbol{\beta}}_{FE} = (\mathbf{X}'\mathbf{M}_D\mathbf{X})^{-1}\mathbf{X}'\mathbf{M}_D\mathbf{y}$$

$\mathbf{D}'\mathbf{y}$  yields an  $m$ -dimensional vector with typical element  $\sum_t y_{it}$  and therefore

$$(\mathbf{D}'\mathbf{D})^{-1}\mathbf{D}'\mathbf{y} = \begin{pmatrix} T^{-1}\sum_t y_{1t} \\ \vdots \\ T^{-1}\sum_t y_{mt} \end{pmatrix} = \begin{pmatrix} \bar{y}_1 \\ \vdots \\ \bar{y}_m \end{pmatrix}$$

Thus  $\mathbf{M}_D\mathbf{y}$  has typical element  $y_{it} - \bar{y}_i$ . The  $\hat{\boldsymbol{\beta}}_{FE}$  estimator makes only use of a variable's variation around the mean for each of the  $m$  groups and, therefore, it is called **within-group estimator**.

The estimator is unbiased since  $\mathbf{X}$  and  $\mathbf{D}$  are exogenous and also efficient. However, if some explanatory variables for an individual  $i$  do not vary over time then columns of  $\mathbf{X}$  and  $\mathbf{D}$  are collinear.

### Random-Effects

If we treat  $\nu_i$  to be stochastic but independent of  $\varepsilon_{it}$  we need to restrict that  $\nu_i$  should be independent of  $\mathbf{X}$ . In that case, OLS estimates  $\boldsymbol{\beta}$  are unbiased but inefficient. This is due to the fact that

$$\begin{aligned} \text{Var} &= \sigma_\nu^2 + \sigma_\varepsilon^2 \\ \text{Cov}(u_{it}, u_{is}) &= \sigma_\nu^2 \\ \text{Cov}(u_{it}, u_{js}) &= 0 \quad \forall i \neq j \end{aligned}$$

Ordering the  $n = mT$  observations by cross-sectional units, we can write the covariance Matrix  $E(\mathbf{u}\mathbf{u}') = \mathbf{\Omega}$  as a block-diagonal matrix

$$\mathbf{\Omega} = \begin{pmatrix} \mathbf{\Sigma} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{\Sigma} & \cdots & \mathbf{0} \\ \vdots & \vdots & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{\Sigma} \end{pmatrix}$$

with  $\mathbf{\Sigma} = \sigma_\varepsilon^2 \mathbf{I}_T + \sigma_\nu^2 \mathbf{u}\mathbf{u}'$ .

GLS yields an efficient estimator if  $\mathbf{\Sigma}$  is known. Problems with the consistency of FGLS arise since  $n \rightarrow \infty$  can be achieved by fixed  $m$  and increasing  $T$  or the reverse.

The counterpart to the within-group estimator is the **between-group estimator** stemming from the fact that  $\mathbf{y} = (\mathbf{M}_D + \mathbf{P}_D)\mathbf{y} = \mathbf{M}_D\mathbf{y} + \mathbf{P}_D\mathbf{y}$ :

$$\hat{\boldsymbol{\beta}}_{BG} = (\mathbf{X}'\mathbf{P}_D\mathbf{X})^{-1}\mathbf{X}'\mathbf{P}_D\mathbf{y}$$

As we have seen before, vector  $\mathbf{P}_D\mathbf{y}$  contains the  $m$  group sample means, however, replicated  $T$  times. The regression

$$\mathbf{P}_D\mathbf{y} = \mathbf{P}_D\mathbf{X}\boldsymbol{\beta} + \text{residual}$$

is a basis for estimating  $\sigma_\nu^2$  whereas the fixed effects model is appropriate to estimate  $\sigma_\varepsilon^2$ .

It is insightful to note that the OLS estimator mentioned at the beginning of this section is a weighted average of  $\hat{\boldsymbol{\beta}}_{BG}$  and  $\hat{\boldsymbol{\beta}}_{FE}$ :

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{M}_D\mathbf{y} + \mathbf{P}_D\mathbf{y}) \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{M}_D\mathbf{X}\hat{\boldsymbol{\beta}}_{FE} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{P}_D\mathbf{X}\hat{\boldsymbol{\beta}}_{BG} \end{aligned}$$

### 6.5.1 Instrumental Variables Estimator

Writing the fixed-effects estimator alternatively as

$$\begin{aligned}\hat{\beta}_{FE} &= \left( \sum_{i=1}^m \sum_{t=1}^T (\mathbf{X}_{it} - \bar{\mathbf{X}}_i)' (\mathbf{X}_{it} - \bar{\mathbf{X}}_i) \right)^{-1} \sum_{i=1}^m \sum_{t=1}^T (\mathbf{X}_{it} - \bar{\mathbf{X}}_i)' (y_{it} - \bar{y}_i) \\ &= \left( \sum_{i=1}^m \sum_{t=1}^T (\mathbf{X}_{it} - \bar{\mathbf{X}}_i)' \mathbf{X}_{it} \right)^{-1} \sum_{i=1}^m \sum_{t=1}^T (\mathbf{X}_{it} - \bar{\mathbf{X}}_i)' y_{it}\end{aligned}$$

the estimator can be interpreted as an instrumental variables estimator for the model

$$y_{it} = \mathbf{X}_{it} \boldsymbol{\beta} + (\nu_i + \varepsilon_{it})$$

where  $\mathbf{X}_{it}$  is instrumented by  $\mathbf{X}_{it} - \bar{\mathbf{X}}_i$ . Since  $E((\mathbf{X}_{it} - \bar{\mathbf{X}}_i)' \nu_i) = \mathbf{0}$  by construction, consistency of this estimator is given provided  $E((\mathbf{X}_{it} - \bar{\mathbf{X}}_i)' \varepsilon_{it}) = \mathbf{0}$ . If we have components in  $\mathbf{X}_{it}$  that are obviously uncorrelated with  $\nu_i$  there is no need for instrumenting those. This allows incorporation of certain time-invariant variables.

An extension to that is a generalized version by noting that not only  $\mathbf{X}_{it} - \bar{\mathbf{X}}_i$  is a valid instrument for  $\mathbf{X}_{it}$  but also  $\mathbf{X}_{i1} - \bar{\mathbf{X}}_i$  up to  $\mathbf{X}_{i,t-1} - \bar{\mathbf{X}}_i$  resulting in more instruments than unknown parameters.

### 6.5.2 Dynamic Linear Models

Extending our linear panel model by a lagged dependent variable leads to dynamic model

$$y_{it} = \mathbf{X}_{it} \boldsymbol{\beta} + \gamma y_{i,t-1} + (\nu_i + \varepsilon_{it})$$

First we study a simpler version of that autoregressive panel model that has no exogenous variables

$$y_{it} = \gamma y_{i,t-1} + (\nu_i + \varepsilon_{it})$$

with  $|\gamma| < 1$ . The fixed effects estimator for  $\gamma$  is

$$\hat{\gamma}_{FE} = \frac{\sum_{i=1}^m \sum_{t=2}^T (y_{it} - \bar{y}_i)(y_{i,t-1} - \bar{y}_{i,-1})}{\sum_{i=1}^m \sum_{t=2}^T (y_{i,t-1} - \bar{y}_{i,-1})^2}$$

where  $\bar{y}_i = (T - 1)^{-1} \sum_{t=2}^T y_{it}$  and  $\bar{y}_{i,-1} = (T - 1)^{-1} \sum_{t=2}^T y_{i,t-1}$ . For **fixed**  $T$  and increasing  $m$  it is an inconsistent estimate since

$$\begin{aligned} \varepsilon_{it} = y_{it} - \gamma y_{i,t-1} - \nu_i &\Rightarrow \sum_{t=2}^T \varepsilon_{it} = \sum_{t=2}^T y_{it} - \gamma \sum_{t=2}^T y_{i,t-1} - \sum_{t=2}^T \nu_i \\ \Rightarrow \bar{\varepsilon}_{it} = \bar{y}_{it} - \gamma \bar{y}_{i,-1} - \nu_i &\Rightarrow \varepsilon_{it} - \bar{\varepsilon}_i = (y_{it} - \bar{y}_i) - \gamma(y_{i,t-1} - \bar{y}_{i,-1}) \\ &\Rightarrow \text{plim}_{m \rightarrow \infty} (\varepsilon_{it} - \bar{\varepsilon}_i)(y_{i,t-1} - \bar{y}_{i,-1}) \neq 0 \end{aligned}$$

One possible solution is a different data transformation to cancel out  $\nu_i$ :

$$y_{it} - y_{i,t-1} = \gamma(y_{i,t-1} - y_{i,t-2}) + (\varepsilon_{it} - \varepsilon_{i,t-1})$$

OLS estimation of this transformation is obviously inconsistent since  $y_{i,t-1}$  is correlated with  $\varepsilon_{i,t-1}$ . However, e.g.  $y_{i,t-2}$  is uncorrelated to both error terms and is therefore a valid instrument if  $\varepsilon_{it}$  is serially uncorrelated. This suggests an IV estimator

$$\gamma_{IV} = \frac{\sum_{i=1}^m \sum_{t=3}^T y_{i,t-2}(y_{it} - y_{i,t-1})}{\sum_{i=1}^m \sum_{t=3}^T y_{i,t-2}(y_{i,t-1} - y_{i,t-2})}$$

Of course,  $(y_{i,t-2} - y_{i,t-3})$  could be used as an instrument as well but suffers from losing one sample period. Again, if  $\varepsilon_{it}$  is serially uncorrelated this estimator is also consistent since

$$E((y_{i,t-2} - y_{i,t-3})(\varepsilon_{it} - \varepsilon_{i,t-1})) = 0$$

This is a moment condition just like  $E((y_{i,t-2})(\varepsilon_{it} - \varepsilon_{i,t-1})) = 0$  so that a GMM approach combining both should be more efficient. Considering a specific  $t$  we have  $t - 2$  moment conditions which can be exploited.

## 6.6 Systems of Regression Equations

In this section we focus on **multivariate models** in which more than one regression equation are considered. There are multiple endogenous variables to be explained simultaneously.

### 6.6.1 Seemingly Unrelated Regressions

We consider  $g$  endogenous variables, indexed by  $i$ , each of them having a linear regression function:

$$\mathbf{y}_i = \underset{(n \times k_i)}{\mathbf{X}_i} \boldsymbol{\beta}_i + \mathbf{u}_i \quad E(\mathbf{u}_i \mathbf{u}_i') = \sigma_{ii} \mathbf{I}$$

In this respect the  $g$  regressions are unrelated leading to the common name *seemingly unrelated regressions* (SUR). However, contemporaneous covariances across these equations are allowed:

$$E(u_{ti} u_{tj}) = \sigma_{ij} \quad \text{for all } t, \quad E(u_{ti} u_{sj}) = 0 \quad \text{for all } t \neq s$$

Defining an  $1 \times g$  vector  $\mathbf{U}_t$  with typical element  $u_{ti}$  and with this an  $n \times g$  matrix of error terms  $\mathbf{U}$  (note that  $\mathbf{U} = [\mathbf{u}_1 : \mathbf{u}_2 : \dots : \mathbf{u}_g]$ ) we can denote

$$E(\mathbf{U}_t' \mathbf{U}_t) = \frac{1}{n} E(\mathbf{U}' \mathbf{U}) = \boldsymbol{\Sigma}$$

In this system we possibly have  $\sum_{i=1}^g k_i = k$  regressors which are not all distinct (e.g.  $\iota$  is most likely included in all  $\mathbf{X}_i$ ). We combine all non-redundant regressors into (full column rank) matrix  $\mathbf{X}$  having  $l < k$  columns. We then can either assume strict exogeneity  $E(\mathbf{U} \mid \mathbf{X}) = \mathbf{O}$  or the weaker condition  $E(\mathbf{U}_t \mid \mathbf{X}_t) = \mathbf{0}$ . Note that  $\mathbf{X}_t$  now denotes the  $t^{\text{th}}$  row of  $\mathbf{X}$  whereas before  $\mathbf{X}_i$  denotes an  $n \times k_i$  matrix.

The entire SUR system can be written compactly as

$$\mathbf{y}_\bullet = \mathbf{X}_\bullet \boldsymbol{\beta}_\bullet + \mathbf{u}_\bullet$$

where  $\mathbf{y}_\bullet$  is a  $gn$  vector resulting from stacking  $\mathbf{y}_1$  through  $\mathbf{y}_g$ . Matrix  $\mathbf{X}_\bullet$  is block-diagonal of dimension  $gn \times k$ :

$$\mathbf{X}_\bullet = \begin{pmatrix} \mathbf{X}_1 & \mathbf{O} & \cdots & \mathbf{O} \\ \mathbf{O} & \mathbf{X}_2 & \cdots & \mathbf{O} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{O} & \mathbf{O} & \cdots & \mathbf{X}_g \end{pmatrix}$$

The covariance matrix of  $\mathbf{u}_\bullet$ , in which the columns of  $\mathbf{U}$  are vertically stacked, is

$$E(\mathbf{u}_\bullet \mathbf{u}_\bullet') = \begin{pmatrix} E(\mathbf{u}_1 \mathbf{u}_1') & \cdots & E(\mathbf{u}_1 \mathbf{u}_g') \\ \vdots & \ddots & \vdots \\ E(\mathbf{u}_g \mathbf{u}_1') & \cdots & E(\mathbf{u}_g \mathbf{u}_g') \end{pmatrix} = \begin{pmatrix} \sigma_{11} \mathbf{I} & \cdots & \sigma_{1g} \mathbf{I} \\ \vdots & \ddots & \vdots \\ \sigma_{g1} \mathbf{I} & \cdots & \sigma_{gg} \mathbf{I} \end{pmatrix} = \boldsymbol{\Sigma}_\bullet$$

This matrix can be written as a *Kronecker product*  $\boldsymbol{\Sigma}_\bullet = \boldsymbol{\Sigma} \otimes \mathbf{I}$ . For  $p \times q$  matrix  $\mathbf{A}$  and an  $r \times s$  matrix  $\mathbf{B}$  the resulting Kronecker product  $\mathbf{A} \otimes \mathbf{B}$  is of dimension  $pr \times qs$  and has as typical elements  $r \times s$  blocks  $a_{ij} \mathbf{B}$ . The general properties are (conformable matrices required):

$$\begin{aligned} (\mathbf{A} \otimes \mathbf{B})' &= \mathbf{A}' \otimes \mathbf{B}' \\ (\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) &= (\mathbf{AC}) \otimes (\mathbf{BD}) \\ (\mathbf{A} \otimes \mathbf{B})^{-1} &= \mathbf{A}^{-1} \otimes \mathbf{B}^{-1} \end{aligned}$$

### GLS Estimation

Using this notation the GLS estimator for a given matrix  $\boldsymbol{\Sigma}$  then is:

$$\begin{aligned} \hat{\boldsymbol{\beta}}_\bullet^{GLS} &= (\mathbf{X}_\bullet' \boldsymbol{\Sigma}_\bullet^{-1} \mathbf{X}_\bullet)^{-1} \mathbf{X}_\bullet' \boldsymbol{\Sigma}_\bullet^{-1} \mathbf{y}_\bullet \\ &= (\mathbf{X}_\bullet' (\boldsymbol{\Sigma}^{-1} \otimes \mathbf{I}) \mathbf{X}_\bullet)^{-1} \mathbf{X}_\bullet' (\boldsymbol{\Sigma}^{-1} \otimes \mathbf{I}) \mathbf{y}_\bullet \end{aligned}$$

The covariance matrix of that estimator thus is

$$V(\hat{\boldsymbol{\beta}}_\bullet^{GLS}) = (\mathbf{X}_\bullet' (\boldsymbol{\Sigma}^{-1} \otimes \mathbf{I}) \mathbf{X}_\bullet)^{-1}$$

**Feasible GLS** first uses the inefficient but at least consistent OLS estimator  $\hat{\boldsymbol{\beta}}_\bullet^{OLS} = (\mathbf{X}_\bullet' \mathbf{X}_\bullet)^{-1} \mathbf{X}_\bullet' \mathbf{y}_\bullet$  to compute the residuals  $\hat{\mathbf{U}}$  and in turn estimating  $\boldsymbol{\Sigma}$  by

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \hat{\mathbf{U}}' \hat{\mathbf{U}}$$

This estimator replaces  $\Sigma$  in the above formulas resulting in  $\hat{\beta}_{\bullet}^{FGLS}$

### GMM Estimation

Considering a single equation isolated we can interpret  $E(\mathbf{X}_i \mathbf{u}_i) = \mathbf{0}$  as exactly identifying moment conditions resulting in the OLS estimator for that isolated equation. For all the other regressors in the other equations not included in  $\mathbf{X}_i$  the moment condition with  $\mathbf{u}_i$  should hold as well. Therefore, we can state  $gl$  theoretical moment conditions

$$E(\mathbf{X}'(\mathbf{y}_i - \mathbf{X}_i \beta_i)) = \mathbf{0} \quad i = 1, \dots, g$$

What we can learn from this is that if all the  $g$  equations use exactly the same regressors  $\mathbf{X} = \mathbf{X}_i$  the empirical counterpart of these moment conditions is

$$\mathbf{X}'(\mathbf{y}_i - \mathbf{X} \beta_i) = \mathbf{0} \quad i = 1, \dots, g$$

and thus OLS estimation for each equation separately is efficient.

For the general case we can write  $\mathbf{X}'(\mathbf{y}_i - \mathbf{X} \beta_i) = \mathbf{0}$  for  $i = 1, \dots, g$  more compactly as

$$(\mathbf{I}_g \otimes \mathbf{X})'(\mathbf{y}_{\bullet} - \mathbf{X}_{\bullet} \beta_{\bullet}) = \mathbf{0}$$

By the same line of argument as in section 6.1 the over-identifying moments in  $\mathbf{X}$  need to be weighted to obtain an efficient GMM estimator. The optimally weighted instruments are thus again  $(\mathbf{I}_g \otimes \mathbf{X})'(\Sigma^{-1} \otimes \mathbf{I}_n) = (\Sigma^{-1} \otimes \mathbf{X}')$  leading to optimally weighted moments

$$(\Sigma^{-1} \otimes \mathbf{X}') \mathbf{u}_{\bullet} = (\Sigma^{-1} \otimes \mathbf{X}')(\mathbf{y}_{\bullet} - \mathbf{X}_{\bullet} \beta_{\bullet})$$

having covariance matrix

$$(\Sigma^{-1} \otimes \mathbf{X}')(\Sigma \otimes \mathbf{I}_n)(\Sigma^{-1} \otimes \mathbf{X}) = \Sigma^{-1} \otimes \mathbf{X}' \mathbf{X}$$

implying the criterion function of the fully efficient GMM:

$$(\mathbf{y}_{\bullet} - \mathbf{X}_{\bullet} \beta_{\bullet})'(\Sigma^{-1} \otimes \mathbf{X})(\Sigma \otimes (\mathbf{X}' \mathbf{X})^{-1})(\Sigma^{-1} \otimes \mathbf{X}')(\mathbf{y}_{\bullet} - \mathbf{X}_{\bullet} \beta_{\bullet})$$

This reduces to  $(\mathbf{y}_{\bullet} - \mathbf{X}_{\bullet} \beta_{\bullet})'(\Sigma^{-1} \otimes \mathbf{P}_X)(\mathbf{y}_{\bullet} - \mathbf{X}_{\bullet} \beta_{\bullet})$  leading to

$$\hat{\beta}_{\bullet}^{GMM} = (\mathbf{X}'_{\bullet}(\Sigma^{-1} \otimes \mathbf{P}_X)\mathbf{X}_{\bullet})^{-1} \mathbf{X}'_{\bullet}(\Sigma^{-1} \otimes \mathbf{P}_X)\mathbf{y}_{\bullet}$$

Now we can see that the resulting estimator is equivalent to  $\hat{\beta}_{\bullet}^{GLS}$ .



### 6.6.2 Linear Simultaneous Equation Models

The  $g$  equations are now extended by allowing some of the regressor variables to be endogenous variables from other equations:

$$\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta}_i + \mathbf{u}_i = \mathbf{Z}_i\boldsymbol{\beta}_{1i} + \mathbf{Y}_i\boldsymbol{\beta}_{2i} + \mathbf{u}_i$$

The  $\mathbf{Z}_i$  are the exogenous or predetermined variables. All the distinct columns of these matrices  $\mathbf{Z}_1, \dots, \mathbf{Z}_g$  are again summarized in the  $n \times l$  matrix  $\mathbf{Z}$ . It must be possible to solve for the endogenous variables as functions of the predetermined variables  $\mathbf{Z}$  which must be linear functions

$$\mathbf{y}_i = \mathbf{Z}\boldsymbol{\pi}_i + \text{residuals}$$

Thus  $\mathbf{Z}$  serve as instruments for  $\mathbf{y}_i$  by using  $\mathbf{P}_Z$  and of course for  $\mathbf{Z}_i$ . We can use

$$\hat{\mathbf{X}}_i = [\mathbf{Z}_i : \mathbf{P}_Z\mathbf{Y}_i] = \mathbf{P}_Z[\mathbf{Z}_i : \mathbf{Y}_i] = \mathbf{P}_Z\mathbf{X}_i$$

as conditional expectations of the regressors. By analogy to the previous subsection, the efficient GMM estimator is

$$\hat{\boldsymbol{\beta}}_{\bullet}^{GMM} = (\mathbf{X}'_{\bullet}(\boldsymbol{\Sigma}^{-1} \otimes \mathbf{P}_Z)\mathbf{X}_{\bullet})^{-1} \mathbf{X}'_{\bullet}(\boldsymbol{\Sigma}^{-1} \otimes \mathbf{P}_Z)\mathbf{y}_{\bullet}$$

#### Feasible GMM or 3SLS

If the matrix  $\boldsymbol{\Sigma}$  is unknown it can be replaced by a consistent estimate. Thus we first estimate  $\boldsymbol{\beta}_{\bullet}$  by an instrumental variable approach using the moment conditions

$$\mathbf{X}'_{\bullet}(\mathbf{I}_g \otimes \mathbf{P}_Z)'(\mathbf{y}_{\bullet} - \mathbf{X}_{\bullet}\boldsymbol{\beta}_{\bullet}) = \mathbf{O}$$

leading to a consistent but inefficient estimator which is usually termed as 2SLS

$$\hat{\boldsymbol{\beta}}_{\bullet}^{2SLS} = (\mathbf{X}'_{\bullet}(\mathbf{I}_g \otimes \mathbf{P}_Z)\mathbf{X}_{\bullet})^{-1} \mathbf{X}'_{\bullet}(\boldsymbol{\Sigma}^{-1} \otimes \mathbf{P}_Z)\mathbf{y}_{\bullet}$$

Using this estimator we can compute residuals and as above estimate the contemporaneous covariance matrix which we denote as  $\hat{\boldsymbol{\Sigma}}_{2SLS}$ . In the next stage we compute the above efficient GMM estimator which is usually called 3SLS:

$$\hat{\boldsymbol{\beta}}_{\bullet}^{3SLS} = \left( \mathbf{X}'_{\bullet}(\hat{\boldsymbol{\Sigma}}_{2SLS}^{-1} \otimes \mathbf{P}_Z)\mathbf{X}_{\bullet} \right)^{-1} \mathbf{X}'_{\bullet}(\hat{\boldsymbol{\Sigma}}_{2SLS}^{-1} \otimes \mathbf{P}_Z)\mathbf{y}_{\bullet}$$

## 7 Discrete and Limited Dependent Variables

### 7.1 Binary Response Models

Binary response models are characterized by the fact that the dependent value  $y_t$  can only take on two values usually coded 0 and 1. As usual in econometrics the conditional mean of  $y_t$  given  $\mathbf{X}_t$  is the main focus.

A binomially distributed random variable  $Z \sim B(1, p)$  has the moments

$$E(Z) = p \quad \text{Var}(Z) = p(1 - p)$$

Thus we consider  $p_t \equiv P(y_t = 1 | \mathbf{X}_t) = E(y_t | \mathbf{X}_t)$  which only can vary between 0 and 1 requiring a mapping of  $\mathbf{X}_t$  and  $\boldsymbol{\beta}$  on to the unit interval. Therefore, the usual specification is

$$p_t \equiv P(y_t = 1 | \mathbf{X}_t) = E(y_t | \mathbf{X}_t) = F(\mathbf{X}_t \boldsymbol{\beta})$$

CDF's fulfill the requirement for such a function  $F(\cdot)$ . Common choices are the cdf of the normal distribution  $F(\mathbf{X}_t \boldsymbol{\beta}) = \Phi(\mathbf{X}_t \boldsymbol{\beta})$  (probit model) or the logistic distribution (logit model). Note, that since cdf's are non-linear the impact of a change in one of the regressors is not constant anymore

$$\frac{\partial p_t}{\partial x_{tj}} = \frac{\partial E(y_t | \mathbf{X}_t)}{\partial x_{tj}} = \frac{\partial F(\mathbf{X}_t \boldsymbol{\beta})}{\partial x_{tj}} = f(\mathbf{X}_t \boldsymbol{\beta}) \cdot \beta_j$$

### Maximum Likelihood Estimation

Using the indicator  $y_t$  we can write the probability function compactly as

$$P(y_t) = p_t^{y_t} \cdot (1 - p_t)^{1 - y_t}$$

Using the expression for  $p_t$  the likelihood function then is

$$f(\boldsymbol{\beta}; \mathbf{y}, \mathbf{X}) = \prod_{t=1}^n F(\mathbf{X}_t \boldsymbol{\beta})^{y_t} \cdot (1 - F(\mathbf{X}_t \boldsymbol{\beta}))^{1 - y_t}$$

from which we obtain the loglikelihood function

$$\ell(\boldsymbol{\beta}; \mathbf{y}, \mathbf{X}) = \sum_{t=1}^n y_t \log F(\mathbf{X}_t \boldsymbol{\beta}) + (1 - y_t) \log(1 - F(\mathbf{X}_t \boldsymbol{\beta}))$$

Taking first derivatives

$$\begin{aligned}\frac{\partial \ell}{\partial \boldsymbol{\beta}} &= \sum_{t=1}^n \frac{y_t}{F(\mathbf{X}_t \boldsymbol{\beta})} f(\mathbf{X}_t \boldsymbol{\beta}) \mathbf{X}_t' + \frac{1-y_t}{1-F(\mathbf{X}_t \boldsymbol{\beta})} (-f(\mathbf{X}_t \boldsymbol{\beta}) \mathbf{X}_t') \\ &= \sum_{t=1}^n \frac{(y_t - F(\mathbf{X}_t \boldsymbol{\beta})) f(\mathbf{X}_t \boldsymbol{\beta})}{F(\mathbf{X}_t \boldsymbol{\beta})(1-F(\mathbf{X}_t \boldsymbol{\beta}))} \mathbf{X}_t'\end{aligned}\quad (13)$$

and second derivatives (using short-hand  $F_t = F(\mathbf{X}_t \boldsymbol{\beta})$  and  $f_t = f(\mathbf{X}_t \boldsymbol{\beta})$ )

$$\frac{\partial^2 \ell}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} = \sum_{t=1}^n \frac{(-f_t^2 + (y_t - F_t) \cdot df_t/dz) F_t(1 - F_t) - (y_t - F_t) f_t (f_t - 2F_t f_t)}{(F_t(1 - F_t))^2} \mathbf{X}_t' \mathbf{X}_t$$

Taking expectations of that Hessian and multiply it with  $-1$  yields the information matrix

$$\mathbf{I}(\boldsymbol{\beta}) = \sum_{t=1}^n \frac{f_t^2}{F_t(1 - F_t)} \mathbf{X}_t' \mathbf{X}_t$$

## Probit Model

The density function of a standard normal distribution is

$$\phi(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}z^2\right)$$

It is easy to see that the first derivative is  $d\phi(z)/dz = \phi(z) \cdot (-z)$ .

The conditional probability or conditional mean in the probit model  $P(y_t = 1 | \mathbf{X}_t) = \Phi(\mathbf{X}_t \boldsymbol{\beta})$  can be derived alternatively assuming a latent variable  $y_t^*$  and its generating process

$$y_t^* = \mathbf{X}_t \boldsymbol{\beta} + u_t \quad u_t \sim IIN(0, 1)$$

We only observe whether the latent variable is larger than a given threshold (0 in this specific case):

$$y_t = \begin{cases} 1 & \text{if } y_t^* > 0 \\ 0 & \text{if } y_t^* \leq 0 \end{cases}$$

From this we derive

$$P(y_t = 1) = P(y_t^* > 0) = P(\mathbf{X}_t \boldsymbol{\beta} + u_t > 0) = 1 - P(u_t < -\mathbf{X}_t \boldsymbol{\beta}) = \Phi(\mathbf{X}_t \boldsymbol{\beta})$$

Replacing  $F(\mathbf{X}_t \boldsymbol{\beta})$  by  $\Phi(\mathbf{X}_t \boldsymbol{\beta})$  and  $f(\mathbf{X}_t \boldsymbol{\beta})$  by  $\phi(\mathbf{X}_t \boldsymbol{\beta})$  in (13) and setting these equations equal to  $\mathbf{0}$  defines the MLE.

## Logit Model

The cdf for the logistic distribution is given by

$$\Lambda(z) = \frac{1}{1 + e^{-z}} = \frac{e^z}{1 + e^z}$$

We see that this distribution is symmetric around 0 since  $\Lambda(-z) = 1 - \Lambda(z)$ . The pdf has the form

$$\lambda(z) = \frac{e^z}{(1 + e^z)^2} = \Lambda(z)\Lambda(-z) = \Lambda(z)(1 - \Lambda(z))$$

Using these expressions replacing  $F$  and  $f$ , respectively, in (13) we obtain

$$\frac{\partial \ell}{\partial \boldsymbol{\beta}} = \sum_{t=1}^n (y_t - \Lambda(\mathbf{X}_t \boldsymbol{\beta})) \mathbf{X}'_t$$

Thus the Hessian simplifies to

$$\frac{\partial^2 \ell}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} = \sum_{t=1}^n -\lambda(\mathbf{X}_t \boldsymbol{\beta}) \mathbf{X}'_t \mathbf{X}_t = -\sum_{t=1}^n \Lambda(\mathbf{X}_t \boldsymbol{\beta})(1 - \Lambda(\mathbf{X}_t \boldsymbol{\beta})) \mathbf{X}'_t \mathbf{X}_t$$

In this case the Hessian does not depend on  $y_t$  so that the negative of it is already the information matrix. If we define a diagonal matrix  $\boldsymbol{\Psi}$  with typical diagonal element  $p_t(1 - p_t) = \Lambda(\mathbf{X}_t \boldsymbol{\beta})(1 - \Lambda(\mathbf{X}_t \boldsymbol{\beta}))$  the Hessian can be written as

$$\frac{\partial^2 \ell}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} = -\mathbf{X}' \boldsymbol{\Psi} \mathbf{X}$$

from which we can deduce negative definiteness if  $\Lambda(\mathbf{X}_t \boldsymbol{\beta})$  shows variation across individuals. The loglikelihood function is therefore globally concave with a unique maximum.

It is interesting to note that  $\Lambda(z)/(1 - \Lambda(z)) = e^z$  and thus

$$\log \left( \frac{p_t}{1 - p_t} \right) = \mathbf{X}_t \boldsymbol{\beta}$$

## 7.2 Ordered Discrete Response Models

We extend the latent variable model

$$y_t^* = \mathbf{X}_t \boldsymbol{\beta} + u_t \quad u_t \sim IIN(0, \sigma^2)$$

by allowing more than two discrete responses, say  $r$  categories, so that the measurement equation becomes to

$$y_t = j \quad \Leftrightarrow \quad \gamma_{t,j-1} < y_t^* \leq \gamma_{t,j} \quad \text{with} \quad \gamma_0 = -\infty \quad \text{and} \quad \gamma_r = \infty$$

The probability for  $y_t$  being observed in category  $j$ :

$$\begin{aligned} P(y_t = j | \mathbf{X}_t) &= \int_{\gamma_{j-1}}^{\gamma_j} f(y_t^*) dy_t^* = \int_{-\infty}^{\gamma_j} f(y_t^*) dy_t^* - \int_{-\infty}^{\gamma_{j-1}} f(y_t^*) dy_t^* \\ &= P(y_t^* < \gamma_j) - P(y_t^* < \gamma_{j-1}) \\ &= F\left(\frac{\gamma_j - \mathbf{X}_t \boldsymbol{\beta}}{\sigma}\right) - F\left(\frac{\gamma_{j-1} - \mathbf{X}_t \boldsymbol{\beta}}{\sigma}\right) \end{aligned}$$

As implicitly done in the last section, some identifying restrictions are required:  $\sigma = 1$  and in case  $\mathbf{X}_t$  contains a constant e.g.  $\gamma_1 = 0$ .

If we define a set of dummy variables  $y_{ij}$  with

$$y_{ij} = \begin{cases} 1 & \text{if } y_i = j \\ 0 & \text{else} \end{cases}$$

we can write probabilities compactly as

$$P(y_i = k) = \prod_{j=1}^r P(y_i = j)^{y_{ij}}$$

resulting in a loglikelihood function

$$\ell(\boldsymbol{\beta}, \boldsymbol{\gamma}) = \sum_{t=1}^n \sum_{j=1}^r y_{tj} \log P(y_t = j | \mathbf{X}_t)$$

### 7.3 Unordered Discrete Response Models

The discrete response categories in this class of models do not show any natural ordering scheme. These models are also referred to as *discrete choice* models.

#### Multinomial Logit Model

A natural extension for the multinomial case is by generalizing the logit case:

$$P(y_t = l) = \frac{\exp(\mathbf{W}_{tl} \boldsymbol{\beta}^l)}{\sum_{j=1}^r \exp(\mathbf{W}_{tj} \boldsymbol{\beta}^j)} \quad \text{for } l = 1, \dots, r$$

The (row) vector of regressors  $\mathbf{W}_{tl}$  now contains individual specific as well as category specific variables. If there are only individual specific regressors,  $\mathbf{W}_{tl}$  is replaced by  $\mathbf{X}_t$ . This shows that there must be category specific parameters  $\beta^l$ . It also shows that not all parameter vectors are identified:

$$\begin{aligned} P(y_t = l) &= \frac{\exp(\mathbf{X}_t \beta^l)}{\sum_{j=1}^r \exp(\mathbf{X}_t \beta^j)} \cdot \frac{\exp(-\mathbf{X}_t \beta^1)}{\exp(-\mathbf{X}_t \beta^1)} = \frac{\exp(\mathbf{X}_t (\beta^l - \beta^1))}{\sum_{j=1}^r \exp(\mathbf{X}_t (\beta^j - \beta^1))} \\ &= \frac{\exp(\mathbf{X}_t (\beta^l - \beta^1))}{1 + \sum_{j=2}^r \exp(\mathbf{X}_t (\beta^j - \beta^1))} \end{aligned}$$

Thus we can restrict e.g.  $\beta^1 = \mathbf{0}$  which must be taken into account interpreting estimated parameters.

Another special case is the *conditional logit model* which uses individual specific choice characteristics  $\mathbf{W}_{tl}$  and fixed  $\beta$

$$P(y_t = l) = \frac{\exp(\mathbf{W}_{tl} \beta)}{\sum_{j=1}^r \exp(\mathbf{W}_{tj} \beta)} \quad \text{for } l = 1, \dots, r$$

One important property of the multinomial logit model is the so-called *independence of irrelevant alternatives* (IIA) which can be seen from

$$\frac{P(y_t = l)}{P(y_t = j)} = \frac{\exp(\mathbf{W}_{tl} \beta^l)}{\exp(\mathbf{W}_{tl} \beta^j)}$$

The ratio of these probabilities does not depend on the regressors or parameters of other response categories. To avoid that sometimes unrealistic property this model is extended to the *nested logit model* where the  $r$  choices are partitioned into  $m$  disjoint subsets  $A_i$  so that

$$P(y_t = l | y_t \in A_i) = \frac{\exp(\mathbf{W}_{tl} \beta^l / \theta_i)}{\sum_{j \in A_i} \exp(\mathbf{W}_{tj} \beta^j / \theta_i)}$$

where  $\theta_i$  is a scale parameter equal for all  $l \in A_i$ . These parameters determines the probability of choosing an element of  $A_i$ .

## Multinomial Pobit Model

A somewhat different foundation is by assuming a latent model

$$y_{tj}^* = \mathbf{W}_{tj} \beta^j + u_{tj} \quad \mathbf{u}_t \sim IIN(\mathbf{0}, \Omega)$$

where category  $j$  is observed and thus  $y_{tj} = 1$  if  $y_{tj}^* - y_{ti}^* \geq 0$  for all  $i = 1, \dots, r$ . The latent model is often interpreted as describing the utility level  $y_{tj}^*$ . The choice with the highest utility is observed. Again, not all the parameter vectors are identified leading to the often used restriction  $\beta^1 = \mathbf{0}$ .

The advantage of the multinomial probit model is the flexible covariance structure  $\Omega$  avoiding the IIA property if the categories are not modelled to be uncorrelated. However, the joint probabilities of an  $r$  dimensional normal distribution (actually only  $r - 1$  dimensions are required) are computationally burdensome since it involves approximation of  $r - 1$  fold integrals.

## 7.4 Count Data Models

For count data the most common distribution used is the Poisson distribution having the probability function

$$P(Y = y) = \frac{e^{-\lambda} \lambda^y}{y!} \quad y = 0, 1, 2, \dots,$$

It follows that

$$E(Y) = \lambda \quad \text{Var}(Y) = \lambda$$

Modelling the mean of  $Y$  as a function of regressors requires its non-negativity. A common choice is

$$\lambda_t = \exp(\mathbf{X}_t \boldsymbol{\beta})$$

implying

$$P(Y_t = y) = \frac{\exp(-\exp(\mathbf{X}_t \boldsymbol{\beta})) \exp(\mathbf{X}_t \boldsymbol{\beta})^y}{y!} \quad y = 0, 1, 2, \dots,$$

The loglikelihood function is then obtained as

$$\ell(\mathbf{y}, \boldsymbol{\beta}) = \sum_{t=1}^n (-\exp(\mathbf{X}_t \boldsymbol{\beta}) + y_t \mathbf{X}_t \boldsymbol{\beta} - \log y_t!)$$

resulting in the FOC

$$\frac{\partial \ell(\mathbf{y}, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \sum_{t=1}^n (y_t - \exp(\mathbf{X}_t \boldsymbol{\beta})) \mathbf{X}_t' = \mathbf{0}$$

and the Hessian

$$\mathbf{H}(\boldsymbol{\beta}) = - \sum_{t=1}^n \exp(\mathbf{X}_t \boldsymbol{\beta}) \mathbf{X}_t' \mathbf{X}_t = -\mathbf{X}' \boldsymbol{\Psi} \mathbf{X}$$

with diagonal matrix  $\boldsymbol{\Psi}$  having typical diagonal element  $\exp(\mathbf{X}_t \boldsymbol{\beta})$ . The Hessian is negative definite if there is variation in the rows of  $\mathbf{X}$  implying global concavity of  $\ell(\mathbf{y}, \boldsymbol{\beta})$  and thus a unique maximum. Note the resemblance of FOC and Hessian to the binary logit model.

## 7.5 Truncated and Censored Models

Both situations can be modelled using the latent model

$$y_t^* = \mathbf{X}_t \boldsymbol{\beta} + u_t \quad u_t \sim IIN(0, \sigma^2)$$

and the measurement equation

$$y_t = \begin{cases} y_t^* & \text{if } y_t^* > 0 \\ 0 & \text{else} \end{cases}$$

### 7.5.1 Truncated Regression Models

If only the observations  $y_t^* > 0$  are observed the sample is truncated. Thus we get

$$\begin{aligned} P(y_t^* > 0) &= P(\mathbf{X}_t \boldsymbol{\beta} + u_t > 0) \\ &= 1 - P(u_t < -\mathbf{X}_t \boldsymbol{\beta}) = 1 - P(u_t/\sigma < -\mathbf{X}_t \boldsymbol{\beta}/\sigma) \\ &= 1 - \Phi(-\mathbf{X}_t \boldsymbol{\beta}/\sigma) = \Phi(\mathbf{X}_t \boldsymbol{\beta}/\sigma) \end{aligned}$$

If all the latent values could be observed directly their density is equal to  $\sigma^{-1} \phi((y_t^* - \mathbf{X}_t \boldsymbol{\beta})/\sigma)$ . Truncation requires that it has to be multiplied by a factor so that it integrates to one:

$$\frac{\sigma^{-1} \phi((y_t^* - \mathbf{X}_t \boldsymbol{\beta})/\sigma)}{\Phi(\mathbf{X}_t \boldsymbol{\beta}/\sigma)}$$



Assume we have  $n_1$  truncation points so that we have  $n_2 = n - n_1$  observations. The loglikelihood function then looks like

$$\ell(\mathbf{y}, \boldsymbol{\beta}, \sigma^2) = -\frac{n_2}{2} \log(2\pi) - n_2 \log \sigma - \frac{1}{2\sigma^2} \sum_{t=1}^{n_2} (y_t^* - \mathbf{X}_t \boldsymbol{\beta})^2 - \sum_{t=1}^{n_2} \log \Phi(\mathbf{X}_t \boldsymbol{\beta} / \sigma)$$

Without the last term it is just the loglikelihood function of a classical normal regression. It shows that OLS estimation with the truncated sample yields inconsistent estimates since it neglects the last term in the loglikelihood function.

### 7.5.2 Censored Regression Models

If we observe all the  $y_t$  values including the  $n_1$  zeros we call that the *Tobit model*. Note that if no observation is censored the loglikelihood function contribution of an observation

$$\ell_t(\mathbf{y}, \boldsymbol{\beta}, \sigma^2) = \sigma^{-1} \phi((y_t^* - \mathbf{X}_t \boldsymbol{\beta}) / \sigma)$$

would result. On the other hand, since

$$P(y_t = 0) = P(y_t^* \leq 0) = P(u_t \leq \mathbf{X}_t \boldsymbol{\beta}) = \Phi(-\mathbf{X}_t \boldsymbol{\beta} / \sigma)$$

we have loglikelihood contributions of the censored observations

$$\ell_t(\mathbf{y}, \boldsymbol{\beta}, \sigma^2) = \log \Phi(-\mathbf{X}_t \boldsymbol{\beta} / \sigma)$$

Combining both yields the loglikelihood function of the Tobit model

$$\ell(\mathbf{y}, \boldsymbol{\beta}, \sigma^2) = \sum_{y_t=0} \log \Phi(-\mathbf{X}_t \boldsymbol{\beta} / \sigma) + \sum_{y_t>0} \log (\sigma^{-1} \phi((y_t^* - \mathbf{X}_t \boldsymbol{\beta}) / \sigma))$$

An OLS regression using the  $y_t$  information thus is inconsistent as well since it is equivalent to a loglikelihood function

$$\sum_{y_t=0} \log (\sigma^{-1} \phi((0 - \mathbf{X}_t \boldsymbol{\beta}) / \sigma)) + \sum_{y_t>0} \log (\sigma^{-1} \phi((y_t^* - \mathbf{X}_t \boldsymbol{\beta}) / \sigma))$$

It is interesting to note that we can write the loglikelihood function as

$$\begin{aligned} \ell(\mathbf{y}, \boldsymbol{\beta}, \sigma^2) &= \sum_{y_t=0} \log \Phi(-\mathbf{X}_t \boldsymbol{\beta} / \sigma) + \sum_{y_t>0} \log \Phi(\mathbf{X}_t \boldsymbol{\beta} / \sigma) \\ &\quad + \sum_{y_t>0} \log (\sigma^{-1} \phi((y_t^* - \mathbf{X}_t \boldsymbol{\beta}) / \sigma)) - \sum_{y_t>0} \log \Phi(\mathbf{X}_t \boldsymbol{\beta} / \sigma) \end{aligned}$$

The first two terms constitute the loglikelihood function of a binary probit just using the censored/not censored information while the last two terms are the loglikelihood function of the truncated sample

## 7.6 Sample Selectivity Models

Closely related to the tobit model is the sample selectivity model:

$$\begin{pmatrix} y_t^* \\ z_t^* \end{pmatrix} = \begin{pmatrix} \mathbf{X}_t \boldsymbol{\beta} \\ \mathbf{W}_t \boldsymbol{\gamma} \end{pmatrix} + \begin{pmatrix} u_t \\ v_t \end{pmatrix} \quad \begin{pmatrix} u_t \\ v_t \end{pmatrix} \sim IIN \left( \mathbf{0}, \begin{pmatrix} \sigma^2 & \rho\sigma \\ \rho\sigma & 1 \end{pmatrix} \right)$$

Measurement equations

$$y_t = \begin{cases} y_t^* & \text{if } z_t^* > 0 \\ \text{unobserved} & \text{otherwise} \end{cases}$$

$$z_t = \begin{cases} 1 & \text{if } z_t^* > 0 \\ 0 & \text{otherwise} \end{cases}$$

A loglikelihood contribution can be written as

$$(1 - z_t)P(z_t = 0) + z_t P(z_t = 1) f(y_t^* | z_t = 1)$$

A joint density  $f(y, z)$  can either be factorized as  $f(y|z)f(z)$  or as  $f(z|y)f(y)$  which is also true for our context. Thus we are concerned in finding an expression for  $P(z_t = 1|y_t^*)$ . In doing so we apply this factorization also to the joint distribution of  $u_t$  and  $v_t$ . The conditional distribution of a normal distribution again is normal. Therefore, we can write

$$v_t = \frac{\rho}{\sigma} u_t + \varepsilon_t \quad \varepsilon_t \sim IIN(0, 1 - \rho^2)$$

To see that this is true we calculate the unconditional first two moments of  $v_t$ :

$$E(v_t) = \frac{\rho}{\sigma} E(u_t) + E(\varepsilon_t) = 0 \quad \text{Var}(v_t) = \frac{\rho^2}{\sigma^2} \text{Var}(u_t) + (1 - \rho^2) = 1$$

Now we can calculate  $P(z_t = 1|y_t^*)$ :

$$P(z_t = 1|y_t^*) = \Phi \left( \frac{\mathbf{W}_t \boldsymbol{\gamma} + \rho(y_t^* - \mathbf{X}_t \boldsymbol{\beta})/\sigma}{(1 - \rho^2)^{1/2}} \right)$$

Combining everything yields the loglikelihood function

$$\begin{aligned} \ell(\boldsymbol{\beta}, \boldsymbol{\gamma}, \rho, \sigma^2) = & \sum_{z_t=0} \log \Phi(-\mathbf{W}_t \boldsymbol{\gamma}) + \sum_{z_t=1} \log \left( \frac{1}{\sigma} \phi((y_t - \mathbf{X}_t \boldsymbol{\beta})/\sigma) \right) \\ & + \sum_{z_t=1} \log \Phi \left( \frac{\mathbf{W}_t \boldsymbol{\gamma} + \rho(y_t - \mathbf{X}_t \boldsymbol{\beta})/\sigma}{(1 - \rho^2)^{1/2}} \right) \end{aligned}$$

In practice, the maximum likelihood estimation is often replaced by *Heckman's two-step estimator*.

## 7.7 Duration Models

Non-negative continuous variables are often modelled using gamma distribution, log-normal distribution, or Weibull distribution.

### Gamma Distribution

The density function of a gamma distribution is

$$f(t) = \frac{\kappa^\theta}{\Gamma(\theta)} t^{\theta-1} \cdot \exp(-\kappa \cdot t) \quad \kappa > 0, \theta > 0, t > 0$$

We obtain  $E(t) = \theta/\kappa$  and  $\text{Var}(t) = \theta/\kappa^2$ . Note, if  $\theta = 1$  we obtain the exponential distributions. The mean depends on two parameters so that it possible to model either of these as a function of explanatory variables. For instance, if we model  $\kappa(\mathbf{X}_i) = \exp(\mathbf{X}_i \boldsymbol{\beta})$  we have

$$\begin{aligned} E(t_i | \mathbf{X}_i) &= \frac{\theta}{\exp(\mathbf{X}_i \boldsymbol{\beta})} = \theta \cdot \exp(-\mathbf{X}_i \boldsymbol{\beta}) \\ \Rightarrow \frac{\partial E(t_i | \mathbf{X}_i)}{\partial x_{ij}} &= \underbrace{\theta \cdot \exp(-\mathbf{X}_i \boldsymbol{\beta})}_{>0} \cdot (-\beta_j) \end{aligned}$$

The loglikelihood function then is

$$\ell(\mathbf{t}, \boldsymbol{\beta}, \theta) = \sum_{i=1}^n \theta \cdot \mathbf{X}_i \boldsymbol{\beta} - \log \Gamma(\theta) + (\theta - 1) \log t_i - t_i \cdot \exp(\mathbf{X}_i \boldsymbol{\beta})$$

## Log-Normal Distribution

Density function:

$$f(t; \mu, \sigma^2) = \frac{1}{t \cdot \sigma \sqrt{2\pi}} \cdot \exp\left(-\frac{1}{2} \left(\frac{\log(t) - \mu}{\sigma}\right)^2\right)$$

The mean value of the random variable  $t$

$$\Rightarrow E(t) = \exp\left(\mu + \frac{\sigma^2}{2}\right)$$

again depends on two parameters allowing to specify either  $\mu_i = \mu(\mathbf{X}_i)$  or  $\sigma_i^2 = \sigma^2(\mathbf{X}_i)$ . Usually the first alternative is chosen. With that the loglikelihood function looks like

$$\ell(\mathbf{t}; \mu, \sigma^2) = \sum_{i=1}^n -\log(t_i \sqrt{2\pi} \sigma) - \frac{1}{2\sigma^2} (\log(t_i) - \mathbf{X}_i \boldsymbol{\beta})^2$$

Considering the F.O.C.

$$\begin{aligned} \frac{\partial \ell(\mathbf{t}; \boldsymbol{\beta}, \sigma^2)}{\partial \boldsymbol{\beta}} &= -\frac{1}{2\sigma^2} \sum_{i=1}^n 2(\log(t_i) - \mathbf{X}_i \boldsymbol{\beta}) \cdot (-\mathbf{X}_i)' \stackrel{!}{=} 0 \\ \frac{\partial \ell(\mathbf{t}; \boldsymbol{\beta}, \sigma^2)}{\partial \sigma} &= \sum_{i=1}^n -\frac{1}{\sigma} - \frac{(-2)}{2\sigma^3} (\log(t_i) - \mathbf{X}_i \boldsymbol{\beta})^2 \stackrel{!}{=} 0 \end{aligned}$$

we see from the first set of equations that this is simply an OLS regression of  $\log(t_i)$  on  $\mathbf{X}_i$  and thus  $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\tilde{\mathbf{t}}$  with  $\tilde{t}_i = \log(t_i)$ . The second set yields

$$\hat{\sigma}_{ML}^2 = \frac{1}{n} \sum_{i=1}^n (\ln(t_i) - \mathbf{X}_i \hat{\boldsymbol{\beta}})^2$$

## Weibull Distribution

Density function:

$$f(t; \alpha, \theta) = \alpha \theta^\alpha t^{\alpha-1} \exp(-(\theta t)^\alpha) \quad \alpha > 0, \theta > 0.$$

The mean value results

$$E(t) = \frac{1}{\theta} \cdot \Gamma\left(1 - \frac{1}{\alpha}\right)$$

for which commonly  $\theta_i = \theta(\mathbf{X}_i) = \exp(\mathbf{X}_i \boldsymbol{\beta})$  is specified.

## Hazard Functions

As an alternative to modelling the mean value of a distribution of random variable  $T$  we can specify the hazard rate as a function of explanatory variables. For a distribution with density  $f(t)$ , cdf  $F(t)$ , and *survivor function*  $S(t) = 1 - F(t)$  the hazard rate is defined

$$h(t) = \frac{f(t)}{S(t)}$$

For interpretation it is quite interesting to note that

$$\begin{aligned} h(t) &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \cdot P(t \leq T \leq t + \Delta t | T \geq t) \\ &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \cdot \frac{P(t \leq T \leq t + \Delta t \wedge T \geq t)}{P(T \geq t)} \\ &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \cdot \frac{F(t + \Delta t) - F(t)}{S(t)} = \frac{F'(t)}{S(t)} \end{aligned}$$

For the loglikelihood function we use the fact that  $f(t) = h(t)S(t)$

$$\ell(\mathbf{t}, \boldsymbol{\theta}) = \sum_{i=1}^n \log f(t_i | \mathbf{X}_i, \boldsymbol{\theta}) = \sum_{i=1}^n \log h(t_i | \mathbf{X}_i, \boldsymbol{\theta}) + \sum_{i=1}^n \log S(t_i | \mathbf{X}_i, \boldsymbol{\theta})$$

For instance, the Weibull distribution has cdf  $F(t) = 1 - \exp(-(\theta t)^\alpha)$  and therefore a hazard function  $h(t) = \alpha \theta^\alpha t^{\alpha-1}$ . We see that if we model  $\theta_i = \exp(\mathbf{X}_i \boldsymbol{\beta})$  it affects the mean value as well as the hazard function.

## Right Censoring

At time  $t'$  not all the spells might have ended so that we have censored observation. In section 7.5.2 with left censoring (censoring of  $y_t^* < 0$ ) we saw that the loglikelihood function has the form

$$\sum_{i \in C} \log F + \sum_{i \in U} \log f$$

Right censoring is therefore characterized by

$$\sum_{i \in U} \log f + \sum_{i \in C} \log(1 - F(t')) = \sum_{i \in U} \log f + \sum_{i \in C} \log S(t')$$

Thus the loglikelihood function for a sample of right censored observation then is

$$\begin{aligned}\ell(\mathbf{t}, \boldsymbol{\theta}) &= \sum_{i \in U} \log f(t_i | \mathbf{X}_i, \boldsymbol{\theta}) + \sum_{i \in C} \log S(t') \\ &= \sum_{i \in U} \log h(t_i | \mathbf{X}_i, \boldsymbol{\theta}) + \sum_{i \in U} \log S(t_i | \mathbf{X}_i, \boldsymbol{\theta}) + \sum_{i \in C} \log S(t' | \mathbf{X}_i, \boldsymbol{\theta}) \\ &= \sum_{i \in U} \log h(t_i | \mathbf{X}_i, \boldsymbol{\theta}) + \sum_{i=1}^n \log S(t_i | \mathbf{X}_i, \boldsymbol{\theta})\end{aligned}$$

## 8 Time Series Analysis

In section 6.4 we already introduced AR(p) processes to model time series behavior. We now use AR and Moving-Average (MA) processes to model a univariate time series of data  $\mathbf{y} = (y_1, y_2, \dots, y_T)'$ .

### 8.1 Stationary Univariate Processes

We already introduced the notion of weak stationarity and also the lag operator  $Ly_t = y_{t-1}$ . Besides AR(p) we also consider MA(q) which is given by

$$y_t = \varepsilon_t - \alpha_1 \varepsilon_{t-1} - \alpha_2 \varepsilon_{t-2} - \dots - \alpha_q \varepsilon_{t-q} \quad \varepsilon_t \sim IID(0, \sigma_\varepsilon^2)$$

Thus, the error term is assumed to be a stationary process called **White Noise**. It is easy to see that an MA(q) is always stationary since

$$\text{Var}(y_t) = (1 + \alpha_1^2 + \alpha_2^2 + \dots + \alpha_q^2) \sigma_\varepsilon^2$$

We start with a AR(1) process

$$y_t = \theta y_{t-1} + \varepsilon_t \quad \varepsilon_t \sim IID(0, \sigma_\varepsilon^2)$$

Repeated substitution yields:

$$\begin{aligned} y_t &= \theta y_{t-1} + \varepsilon_t = \theta(\theta y_{t-2} + \varepsilon_{t-1}) + \varepsilon_t = \theta^2 y_{t-2} + \varepsilon_t + \theta \varepsilon_{t-1} \\ &= \theta^2(\theta y_{t-3} + \varepsilon_{t-2}) + \varepsilon_t + \theta \varepsilon_{t-1} = \theta^3 y_{t-3} + \varepsilon_t + \theta \varepsilon_{t-1} + \theta^2 \varepsilon_{t-2} \\ &\vdots \\ &= \theta^k y_{t-k} + \varepsilon_t + \theta \varepsilon_{t-1} + \theta^2 \varepsilon_{t-2} + \dots + \theta^{k-1} \varepsilon_{t-k+1} \end{aligned}$$

Provided that  $|\theta| < 1$  the first term vanishes if  $k \rightarrow \infty$  and we obtain an MA( $\infty$ ) as a representation of the AR(1) process having

$$\text{Var}(y_t) = (1 + \alpha_1^2 + \alpha_2^2 + \dots) \sigma_\varepsilon^2 = \sigma_\varepsilon^2 \sum_{j=0}^{\infty} (\alpha_j^j)^2 = \sigma_\varepsilon^2 \sum_{j=0}^{\infty} (\alpha_j^2)^j = \frac{\sigma_\varepsilon^2}{1 - \alpha_j^2}$$

This is called **inversion** of an AR(1) into an MA( $\infty$ ) which can shortly be written by using the lag operator ( $y_t = \theta L y_t + \varepsilon_t$ )

$$(1 - \theta L)y_t = \varepsilon_t \quad \Leftrightarrow \quad y_t = (1 - \theta L)^{-1}\varepsilon_t = (1 + \theta L + \theta^2 L^2 + \theta^3 L^3 + \dots)\varepsilon_t$$

Repeated substitution can also be done for an MA(1)

$$\begin{aligned} \varepsilon_t &= y_t + \alpha \varepsilon_{t-1} = y_t + \alpha(y_{t-1} + \alpha \varepsilon_{t-2}) = y_t + \alpha y_{t-1} + \alpha^2 \varepsilon_{t-2} \\ &= y_t + \alpha y_{t-1} + \alpha^2(y_{t-2} + \alpha \varepsilon_{t-3}) = y_t + \alpha y_{t-1} + \alpha^2 y_{t-2} + \alpha^3 \varepsilon_{t-3} \\ &\vdots \\ &= y_t + \alpha y_{t-1} + \alpha^2 y_{t-2} + \dots + \alpha^k y_{t-k} + \alpha^{k+1} \varepsilon_{t-k-1} \end{aligned}$$

Given that  $|\alpha| < 1$  the MA(1) is represented by an AR( $\infty$ )

$$y_t = (1 - \alpha L)\varepsilon_t \quad \Leftrightarrow \quad \varepsilon_t = (1 - \alpha L)^{-1}y_t = (1 + \alpha L + \alpha^2 L^2 + \alpha^3 L^3 + \dots)y_t$$

An AR(2) model

$$y_t = \theta_1 y_{t-1} + \theta_2 y_{t-2} + \varepsilon_t = \theta_1 L y_t + \theta_2 L^2 y_t + \varepsilon_t \quad \Leftrightarrow \quad (1 - \theta_1 L - \theta_2 L^2)y_t = \varepsilon_t$$

is stationary if the (possibly complex) roots of the characteristic equation

$$\theta(z) = 1 - \theta_1 z - \theta_2 z^2 = 0$$

all lie outside the unit circle. Thus the process  $y_t = 1.4y_{t-1} - 0.9y_{t-2} + \varepsilon_t$  is stationary since the roots are 1.25 and 4 which is also true for  $y_t = 1.05y_{t-1} - 0.2y_{t-2} + \varepsilon_t$  having complex roots outside the unit circle.

The combination of AR(p) and MA(q) yields the ARMA(p,q) model

$$y_t = \theta_1 y_{t-1} + \theta_2 y_{t-2} + \dots + \theta_p y_{t-p} + \varepsilon_t - \alpha_1 \varepsilon_{t-1} - \alpha_2 \varepsilon_{t-2} - \dots - \alpha_q \varepsilon_{t-q}$$

which can be written as  $\theta(L)y_t = \alpha(L)\varepsilon_t$  and thus (if invertability is given)  $y_t = \theta(L)^{-1}\alpha(L)\varepsilon_t$  or  $\alpha(L)^{-1}\theta(L)y_t = \varepsilon_t$ .

An AR(p) process can be estimated using OLS regressing  $y_t$  on its lagged observations. As stated in chapter 2 the assumption of  $E(\varepsilon_t | y_{t-1}, \dots, y_{t-p}) = 0$  is, however, not sufficient to ensure unbiasedness of the estimator but it ensures consistency (see section 3.2).



## 8.2 Stationary Multivariate Processes, VAR

Now we consider a vector of  $m$  time series  $\mathbf{y}_t = (y_{1,t}, y_{2,t}, \dots, y_{m,t})'$  and model them simultaneously using a stationary autoregressive model

$$\mathbf{y}_t = \Theta_1 \mathbf{y}_{t-1} + \boldsymbol{\varepsilon}_t \quad \boldsymbol{\varepsilon}_t \sim IID(\mathbf{0}, \Sigma)$$

As an example, a bivariate vector autoregressive model reads like

$$\begin{aligned} y_{1,t} &= \theta_{11}y_{1,t-1} + \theta_{12}y_{2,t-1} + \varepsilon_{1,t} \\ y_{2,t} &= \theta_{21}y_{1,t-1} + \theta_{22}y_{2,t-1} + \varepsilon_{2,t} \\ \Leftrightarrow \begin{pmatrix} y_{1,t} \\ y_{2,t} \end{pmatrix} &= \begin{pmatrix} \theta_{11} & \theta_{12} \\ \theta_{21} & \theta_{22} \end{pmatrix} \begin{pmatrix} y_{1,t-1} \\ y_{2,t-1} \end{pmatrix} + \begin{pmatrix} \varepsilon_{1,t} \\ \varepsilon_{2,t} \end{pmatrix} \end{aligned}$$

If e.g.  $\theta_{22} = 0$  and we insert the second equation into the first, the system reduce to a univariate AR(2) process of  $y_{1,t}$ . The corresponding parameters are  $\theta_1 = \theta_{11}$  and  $\theta_2 = \theta_{12}\theta_{21}$ , respectively.

To determine the joint stationarity of vector  $\mathbf{y}_t$  we have to solve the characteristic equation obtained from  $|\mathbf{I} - \Theta_1 z| = 0$ :

$$\begin{vmatrix} 1 - \theta_{11}z & -\theta_{12}z \\ -\theta_{21}z & 1 - \theta_{22}z \end{vmatrix} = (1 - \theta_{11}z)(1 - \theta_{22}z) - \theta_{12}\theta_{21}z^2 = 0$$

If all the roots of this equation lie outside the unit circle the multivariate process is jointly stationary. For example, for

$$\Theta_1 = \begin{pmatrix} .8 & -.6 \\ .7 & .2 \end{pmatrix}$$

the complex roots of the characteristic function lie outside the unit circle so that the joint process is stationary. Another aspect in this context is **Granger causality**.

In this simple bivariate context  $y_2$  does not Granger cause  $y_1$  if  $\theta_{12}$  is equal to zero.

The generalization is the VAR(p) where vector  $\mathbf{y}_t$  is explained by lagged vectors up to order p:

$$\begin{aligned} \mathbf{y}_t &= \Theta_1 \mathbf{y}_{t-1} + \Theta_2 \mathbf{y}_{t-2} + \dots + \Theta_p \mathbf{y}_{t-p} + \boldsymbol{\varepsilon}_t \\ \Leftrightarrow (\mathbf{I} - \Theta_1 L - \Theta_2 L^2 - \dots - \Theta_p L^p) \mathbf{y}_t &= \boldsymbol{\varepsilon}_t \quad \Leftrightarrow \Theta(L) \mathbf{y}_t = \boldsymbol{\varepsilon}_t \end{aligned}$$

The process is jointly stationary if the roots of the characteristic equation  $|\Theta(z)| = 0$  all lie outside the unit circle.

Estimation of a VAR(p) can be performed with a SUR. Since each equation depends on lags of all the  $m$  variables the regressor matrix is the same for each variable. Therefore, the VAR system can be estimated equation by equation using OLS to obtain asymptotically efficient estimators.

If the VAR(p) is invertible, i.e.  $|\Theta(1)| = |(\mathbf{I} - \Theta_1 - \Theta_2 - \dots - \Theta_p)| \neq 0$ , it can be written as a VMA( $\infty$ )

$$\mathbf{y}_t = \Theta(L)^{-1} \boldsymbol{\varepsilon}_t = \boldsymbol{\varepsilon}_t + \mathbf{A}_1 \boldsymbol{\varepsilon}_{t-1} + \mathbf{A}_2 \boldsymbol{\varepsilon}_{t-2} + \dots$$

from which we derive

$$\frac{\partial \mathbf{y}_{t+s}}{\partial \boldsymbol{\varepsilon}_t'} = \mathbf{A}_s \quad s > 0$$

Thus a change of the  $k$ -th component in  $\boldsymbol{\varepsilon}_t$  by one unit we find the effects on  $\mathbf{y}_{t+s}$  in the  $k$ -th column of  $\mathbf{A}_s$ . Therefore, we derive the dynamic effects of such a change on the  $j$ -th component of vector  $\mathbf{y}$  by plotting the  $j, k$ -elements of matrices  $\mathbf{I}, \mathbf{A}_1, \mathbf{A}_2, \dots$  which is called the **impulse-response function**.

## 8.3 Nonstationary Time Series

### 8.3.1 Deterministic and Stochastic Trends

Economic time series like GDP, consumption, investments, etc. often show a positive trend. With a slightly modified AR(p)

$$y_t = \mu + \theta_1 y_{t-1} + \theta_2 y_{t-2} + \dots + \theta_p y_{t-p} + \varepsilon_t \quad (14)$$

we obtain for a stationary process an unconditional mean of

$$E(y_t) = \frac{\mu}{1 - \sum_{j=1}^p \theta_j}$$

which is obviously not a good description of the real data. The same is true for a modified MA(q)  $y_t = \mu + \alpha(L)\varepsilon_t$  having an unconditional mean of  $\mu$ . One way

to obtain a nonstationary series with an increasing unconditional mean is a **deterministic trend** model

$$y_t = \mu + \delta \cdot t + \alpha(L)\varepsilon_t$$

Due to the MA structure it is stationary around the trend. Another common model is the *random walk with drift*

$$y_t = \delta + y_{t-1} + \varepsilon_t \quad (15)$$

Note that  $\Delta y_t = (1 - L)y_t = y_t - y_{t-1} = \delta + \varepsilon_t$ . From section 8.1 it is clear that the process with an assumed  $y_0 = 0$  can be written as

$$y_t = \delta \cdot t + \sum_{i=1}^t \varepsilon_i$$

Comparing an AR(1) from (14) and (15) we find that the constant term has a different meaning once  $\theta_1$  becomes zero which also holds true for the variances. Due to that the process is said to have a **stochastic trend**. Processes like (15) are called *integrated* of order one, or  $I(1)$ , since the first differences are stationary,  $I(0)$ . It is interesting to consider the effects of differencing a stationary AR(1) process

$$\begin{aligned} y_t = \mu + \theta_1 y_{t-1} + \varepsilon_t &\Leftrightarrow y_t - y_{t-1} = \mu + \theta_1 y_{t-1} - (\mu + \theta_1 y_{t-2} + \varepsilon_{t-1}) + \varepsilon_t \\ &\Leftrightarrow \Delta y_t = \theta_1 \Delta y_{t-1} - (1 - L)\varepsilon_t \end{aligned}$$

Thus the resulting ARMA(1,1) process contains a noninvertible MA(1) (unit root in the MA part) which is sometimes called an *over-differenced* process.

In the case of an integrated random walk with drift *and* trend

$$y_t = \delta + \gamma t + y_{t-1} + \varepsilon_t$$

we get by setting  $y_0 = 0$

$$\begin{aligned} y_1 &= \delta + \gamma \cdot 1 + \varepsilon_1 \\ y_2 &= 2\delta + \gamma \cdot (1 + 2) + \varepsilon_1 + \varepsilon_2 \\ &\vdots \\ y_t &= \delta \cdot t + \gamma \cdot \frac{t(t+1)}{2} + \sum_{i=1}^t \varepsilon_i \end{aligned}$$

The series shows a quadratic trend and an increasing variance whereas the differenced series is white noise plus a constant:  $\Delta y_t = \gamma + \varepsilon_t$ . Moreover, the differenced series of a process

$$y_t = \delta + \gamma t + \lambda t^2 + y_{t-1} + \varepsilon_t$$

has the form

$$\Delta y_t = (\gamma - \kappa) + (2\kappa) \cdot t + \varepsilon_t$$

### 8.3.2 Estimation of AR(1) Models

Estimation of a stationary AR(1) model which has no constant using OLS obtains

$$\hat{\theta} = \frac{\sum_{t=2}^T y_{t-1} y_t}{\sum_{t=2}^T y_{t-1}^2} = \frac{T^{-1} \sum_{t=2}^T y_{t-1} (\theta y_{t-1} + \varepsilon_t)}{T^{-1} \sum_{t=2}^T y_{t-1}^2} = \theta + \frac{T^{-1} \sum_{t=2}^T y_{t-1} \varepsilon_t}{T^{-1} \sum_{t=2}^T y_{t-1}^2}$$

Rearranging the first equation yields  $\sum_{t=2}^T y_{t-1} y_t = \hat{\theta} \sum_{t=2}^T y_{t-1}^2$ . The estimator's variance is

$$\text{Var}(\hat{\theta}) = \frac{\sigma_\varepsilon^2}{\sum_{t=2}^T y_{t-1}^2}$$

so that via the estimate of  $\sigma_\varepsilon^2$

$$\begin{aligned} \hat{\sigma}_\varepsilon^2 &= T^{-1} \sum_{t=2}^T (y_t - \hat{\theta} y_{t-1})^2 = T^{-1} \left( \sum_{t=2}^T y_t^2 - 2\hat{\theta} \sum_{t=2}^T y_t y_{t-1} + \hat{\theta}^2 \sum_{t=2}^T y_{t-1}^2 \right) \\ &= T^{-1} \left( \sum_{t=2}^T y_t^2 - \hat{\theta}^2 \sum_{t=2}^T y_{t-1}^2 \right) = T^{-1} \left( \sum_{t=3}^{T+1} y_{t-1}^2 - \hat{\theta}^2 \sum_{t=2}^T y_{t-1}^2 \right) \\ &\approx T^{-1} (1 - \hat{\theta}^2) \sum_{t=2}^T y_{t-1}^2 \end{aligned}$$

we get

$$\widehat{\text{Var}}(\hat{\theta}) \approx T^{-1} (1 - \hat{\theta}^2)$$

and thus using asymptotic reasoning

$$\sqrt{T}(\hat{\theta} - \theta) \xrightarrow{a} N(0, 1 - \theta^2)$$

Thus the t-statistic is asymptotically standard normally distributed. It can be shown that for the nonstationary AR(1) with  $\theta = 1$  the convergence rate is  $T$  and not  $\sqrt{T}$

so that under  $H_0 : \theta = 1$  it holds

$$T(\hat{\theta} - 1) \xrightarrow{a} N(0, V)$$

There is no closed analytical expression for variance  $V$  but it can be obtained using Monte Carlo methods. As a consequence the usual t-statistic

$$t = \frac{\sqrt{T}(\hat{\theta} - 1)}{\sqrt{1 - \hat{\theta}^2}} \quad (16)$$

is not asymptotically standard normal and needs to be simulated as well.

### 8.3.3 Unit Roots Tests

Several papers by Dickey and Fuller (1979) give critical values for the t-statistic in (16) of the one sided test of  $H_0 : \theta = 1$  for different  $T$  and different specifications. In order to obtain the test statistic immediately from a standard regression output the regression model is specified as

$$\Delta y_t = y_t - y_{t-1} = (\theta - 1)y_{t-1} + \varepsilon_t$$

To allow for a trend stationary AR(1) as the alternative we could specify a constant and a time trend in the regression specification. Note that in the presence of a unit root it implies a deterministic linear and quadratic time trend as shown in section 8.3.1. To be precise, in such a situation the the null hypothesis should read like  $H_0 : \theta = 1, \text{ no constant, no trend}$  which could be tested jointly. Due to practical reasons, however, solely  $\theta - 1 = 0$  is usually tested.

Critical values for  $\alpha = 5\%$  and different specifications are:

	No constant	Constant	Constant
T	No trend	No trend	Trend
25	-1.95	-3.00	-3.60
50	-1.95	-2.93	-3.50
500	-1.95	-2.87	-3.42
⋮			
∞	-1.95	-2.86	-3.41

## Higher Order AR Processes

An AR(2) model can be rewritten as

$$\begin{aligned} y_t &= \theta_1 y_{t-1} + \theta_2 y_{t-2} + \varepsilon_t \\ \Leftrightarrow y_t - y_{t-1} &= \theta_1 y_{t-1} - y_{t-1} + \theta_2 y_{t-1} - \theta_2 y_{t-2} + \theta_2 y_{t-2} + \varepsilon_t \\ \Leftrightarrow \Delta y_t &= (\theta_1 + \theta_2 - 1)y_{t-1} - \theta_2 \Delta y_{t-1} + \varepsilon_t \end{aligned}$$

From the characteristic equation of this AR model  $1 - \theta_1 z - \theta_2 z^2 = 0$  we see that if the roots are real values then the process is nonstationary if  $z = 1$  and thus  $1 - \theta_1 - \theta_2 = 0$ . Fortunately, under the null  $H_0 : \theta_1 + \theta_2 - 1 = 0$  the critical values of the t-statistic  $(\hat{\theta}_1 + \hat{\theta}_2 - 1)/\sqrt{\text{Var}(\hat{\theta}_1 + \hat{\theta}_2)}$  are the same as for the Dickey-Fuller test given above. A generalization is called the **Augmented Dickey-Fuller test** (ADF) which has the form

$$\Delta y_t = \kappa y_{t-1} + \vartheta_1 \Delta y_{t-1} + \dots + \vartheta_{p-1} \Delta y_{t-p+1} + \varepsilon_t$$

and we test  $H_0 : \kappa = 0$  which is equivalent to  $(\sum_{i=1}^p \theta_i - 1)$  against  $H_1 : \kappa < 0$ .

## 8.4 Cointegration

### 8.4.1 Error Correction model

A distributed lag model

$$Y_t = \mu + \sum_{i=0}^m \beta_i X_{t-i} + u_t$$

with  $\beta_i$  being proportional to probabilities of a geometric distribution with parameter  $\lambda$

$$\beta_i = \beta \cdot P(t = i) = \beta \cdot \lambda^i \cdot (1 - \lambda)$$

after the Koyck-Transformation can be written as

$$Y_t = \mu(1 - \lambda) + \lambda Y_{t-1} + \beta(1 - \lambda)X_t + u_t - \lambda u_{t-1}$$

A more general form of these distributed lag models can be defined as

$$Y_t = \delta + \theta Y_{t-1} + \phi_0 X_t + \phi_1 X_{t-1} + \varepsilon_t \quad (17)$$

which usually called autoregressive distributed lag model ADL(1,1). Such a model with stationary  $Y_t$  and  $X_t$  shows an immediate effect of a change in  $X_t$  expressed by

$$\frac{\partial Y_t}{\partial X_t} = \phi_0$$

and an effect in next period of

$$\frac{\partial Y_{t+1}}{\partial X_t} = \theta \frac{\partial Y_t}{\partial X_t} + \phi_1 = \theta\phi_0 + \phi_1$$

and after two periods

$$\frac{\partial Y_{t+2}}{\partial X_t} = \theta \frac{\partial Y_{t+1}}{\partial X_t} = \theta(\theta\phi_0 + \phi_1)$$

and so on. If we impose  $|\theta| < 1$  the long-run effect of a one unit change in  $X_t$  is

$$\begin{aligned} & \phi_0 + (\theta\phi_0 + \phi_1) + \theta(\theta\phi_0 + \phi_1) + \theta^2(\theta\phi_0 + \phi_1) + \dots \\ &= \phi_0 + (\theta\phi_0 + \phi_1)(1 + \theta + \theta^2 + \dots) = \frac{\phi_0(1 - \theta)}{1 - \theta} + \frac{\theta\phi_0 + \phi_1}{1 - \theta} \\ &= \frac{\phi_0 + \phi_1}{1 - \theta} \end{aligned}$$

The long-run equilibrium of (17) can be written as

$$E(Y_t) = \frac{\delta}{1 - \theta} + \frac{\phi_0 + \phi_1}{1 - \theta} E(X_t) = \alpha + \beta E(X_t)$$

If we subtract  $Y_{t-1}$  on both sides of (17) and expand with  $\phi_0 X_{t-1} - \phi_0 X_{t-1}$  we obtain

$$\Delta Y_t = \delta - (1 - \theta)Y_{t-1} + \phi_0 \Delta X_t + (\phi_1 - \phi_0)X_{t-1} + \varepsilon_t$$

or

$$\Delta Y_t = \phi_0 \Delta X_t - (1 - \theta)(Y_{t-1} - \alpha - \beta X_{t-1}) + \varepsilon_t$$

which is called the **error correction** formulation of the ADL(1,1). If  $Y_{t-1}$  is above its long-run equilibrium value  $\alpha + \beta X_{t-1}$  a negative adjustment on  $Y_t$  is generated with  $(1 - \theta)$  being the adjustment speed.

## 8.4.2 Spurious Regression

Let us assume a deterministic trend in the two variables  $Y_t$  and  $X_t$

$$\begin{aligned} Y_t &= a + b \cdot t + \varepsilon_{1,t} & \varepsilon_{1,t} &\sim IID(0, \sigma_{\varepsilon_1}^2) \\ X_t &= c + d \cdot t + \varepsilon_{2,t} & \varepsilon_{2,t} &\sim IID(0, \sigma_{\varepsilon_2}^2) \end{aligned}$$

with independent error terms. A regression of  $Y_t$  on  $X_t$  will result in usually highly significant regression coefficients and high  $R^2$ . In large samples the slope coefficient converges to  $b/d$  and the constant term to  $a - c \cdot b/d$  due to the omitted variable  $t$  which is a **spurious regression**. However, regressing  $Y_t$  on  $X_t$  and  $t$  the coefficient of  $X_t$  is mostly insignificant.

Considering two variables with stochastic trends like independent random walks

$$\begin{aligned} Y_t &= Y_{t-1} + \varepsilon_{1,t} & \varepsilon_{1,t} &\sim IID(0, \sigma_{\varepsilon_1}^2) \\ X_t &= X_{t-1} + \varepsilon_{2,t} & \varepsilon_{2,t} &\sim IID(0, \sigma_{\varepsilon_2}^2) \end{aligned}$$

we also obtain spurious results of high t- and F-statistics, a high  $R^2$ , and a Durbin-Watson statistic close to zero. The problem is that since  $Y_t$  and  $X_t$  are both  $I(1)$ , the error term  $\varepsilon_t$  is also  $I(1)$  and therefore the t- and F-test statistics have no well-defined asymptotic distributions. However, a specification like (17) will obtain consistent estimates of  $\theta = 1$  and  $\phi_0 = \phi_1 = 0$  due to the residual being  $I(0)$ .

### 8.4.3 Cointegration of two Variables

If we have two random walk processes as described above where the two variables are each  $I(1)$  and hence their respective first differences  $I(0)$  but additionally there exists a value  $\beta$  such that  $Y_t - \beta X_t$  is  $I(0)$  then we describe those variables to be **cointegrated**. We call the vector  $\beta = (1; -\beta)'$  the cointegration vector for the variable vector  $(Y_t; X_t)'$  since  $\beta'(Y_t; X_t)'$  is  $I(0)$ . The OLS estimates of the specification

$$Y_t = a + bX_t + e_t \tag{18}$$

are called *super consistent*. The intuition behind that is, that estimates  $\hat{b} \neq \beta$  produces residuals that are not  $I(0)$  and thus the sample variance of  $\hat{e}_t$  is much larger than in the case of  $\hat{b} = \beta$  in which the residuals are  $I(0)$  and thus have smaller sample variance. Thus the estimates are forced to be closer to the true value than in a situation where the alternative is also  $I(0)$ . In fact,  $\hat{b}$  is not  $\sqrt{T}$  consistent for  $\beta$  but has a faster convergence rate of  $T$ .



**Testing** for cointegration requires in a first step separate unit root tests for  $Y_t$  and  $X_t$ , respectively. In a second step OLS estimates of (18) are computed and a Dickey-Fuller test for a unit root in the residuals is performed

$$\Delta \hat{e}_t = \varphi \hat{e}_t + \nu_t$$

If the null hypothesis  $H_0 : \varphi = 0$  can be rejected the integrated variables  $Y_t$  and  $X_t$  are cointegrated. However, the Dickey-Fuller critical values need to be adjusted for the  $\hat{e}_t$  are not a given time series but are estimated. Again simulation techniques can be used to determine those critical values for this case. For instance, for  $T = 500$  the critical value for the usual test statistic of  $\hat{\varphi}$  is -3.34 ( $\alpha = 5\%$ ) instead of -2.87 for, say,  $Y_t$ .

Assuming a partial adjustment model for  $Y_t$  where the long-run target is  $Y_t^* = \beta X_t$

$$\begin{aligned} \Delta Y_t &= \lambda_1 \Delta Y_t^* + \lambda_2 (Y_{t-1} - Y_{t-1}^*) + \varepsilon_t \\ &= \lambda_1 \beta \Delta X_t + \lambda_2 (Y_{t-1} - \beta X_{t-1}) + \varepsilon_t \end{aligned}$$

where we get an error correction formulation where all the terms are  $I(0)$ . If we allow for a constant in the ECM it implies an additional deterministic trend in both variables.

#### 8.4.4 Cointegration in a VAR

Extending the idea of cointegration to a  $k$ -vector of variables with more than two components each being  $I(1)$  then there might be a  $(k \times r)$  cointegration matrix  $\beta$  with rank  $r \leq k - 1$  such that  $\beta' \mathbf{y}_t$  describing various  $I(0)$  equilibrium relations. We restrict our analysis to a VAR(3) process

$$\begin{aligned} \mathbf{y}_t &= \Theta_1 \mathbf{y}_{t-1} + \Theta_2 \mathbf{y}_{t-2} + \Theta_3 \mathbf{y}_{t-3} + \varepsilon_t \\ \Leftrightarrow (\mathbf{I} - \Theta_1 L - \Theta_2 L^2 - \Theta_3 L^3) \mathbf{y}_t &= \varepsilon_t \quad \Leftrightarrow \Theta(L) \mathbf{y}_t = \varepsilon_t \end{aligned}$$

We rewrite this in the manner of higher order univariate AR in 8.3.3.

$$\begin{aligned}
\Delta \mathbf{y}_t &= (\Theta_1 + \Theta_2 - \mathbf{I})\mathbf{y}_{t-1} - \Theta_2\Delta \mathbf{y}_{t-1} + \Theta_3\mathbf{y}_{t-3} + \varepsilon_t \\
&= (\Theta_1 + \Theta_2 + \Theta_3 - \mathbf{I})\mathbf{y}_{t-1} - \Theta_2\Delta \mathbf{y}_{t-1} - \Theta_3(\Delta \mathbf{y}_{t-1} + \Delta \mathbf{y}_{t-2}) + \varepsilon_t \\
&= (\Theta_1 + \Theta_2 + \Theta_3 - \mathbf{I})\mathbf{y}_{t-1} + \Gamma_1\Delta \mathbf{y}_{t-1} + \Gamma_2\Delta \mathbf{y}_{t-2} + \varepsilon_t
\end{aligned}$$

with  $\Gamma_1 = -\Theta_2 - \Theta_3$  and  $\Gamma_2 = -\Theta_3$ . For a VAR(p) we obtain in the same fashion

$$\Delta \mathbf{y}_t = \Gamma_1\Delta \mathbf{y}_{t-1} + \dots + \Gamma_{p-1}\Delta \mathbf{y}_{t-p+1} + \Pi\mathbf{y}_{t-1} + \varepsilon_t$$

with  $\Pi = -(\mathbf{I} - \Theta_1 - \dots - \Theta_p) = -\Theta(1)$ . Now we can consider three cases:

- $\Pi = \mathbf{0}$  which implies non-stationarity for the  $\mathbf{y}_t$  components since  $\Theta(1) = \mathbf{0}$  and thus a stationary VAR(p-1) for  $\Delta \mathbf{y}_t$ .
- If all  $\mathbf{y}_t$  components are stationary  $\Pi = -\Theta(1)$  must have full rank and can be inverted. Thus a VMA representation  $\mathbf{y}_t = \Theta(L)^{-1}\varepsilon_t$  is obtained
- $\Pi$  is not zero and not invertible (having rank  $r$ ) and can be written as  $\Pi = \gamma\beta'$  where the latter two matrices both have rank  $r$ . Thus we have  $r$  cointegrating relationships so that we obtain  $\mathbf{z}_t = \beta'\mathbf{y}_t$  which represent  $r$  stationary series. Inserting it in the last equation we can interpret it as a **vector error correction model** (VECM)

$$\Delta \mathbf{y}_t = \Gamma_1\Delta \mathbf{y}_{t-1} + \dots + \Gamma_{p-1}\Delta \mathbf{y}_{t-p+1} + \gamma\beta'\mathbf{y}_{t-1} + \varepsilon_t$$