

# Modern Basic Econometrics ⊖ Lectures

## INTRODUCTION TO MATRIX APPROACH TO ECONOMETRICS

Esben Høg

Department of Mathematical Sciences  
Aalborg University

18. februar 2016

# List of Slides

1	Introduction	2
2	The Multiple Regression Model	3
3	The Multivariate Normal Distribution	6
4	Conditioning a multivariate normal distribution	8
5	The Bivariate Normal Distribution	9
6	Least Squares and differentiation w.r.t. a vector	18
	6.1 Some formulas of Matrix Differentiation	21
	6.2 The distribution of $\beta$	29
	6.3 The Gauss-Markov theorem	32
7	The decomposition of sums of squares	33
8	Geometric interpretation of OLS	37
9	Linear Restrictions	41
	9.1 Estimation	42
	9.2 Tests of linear restrictions	47
	9.3 The Chow Forecast test	49
	9.4 Chow's Breakpoint Test	57
10	Heteroskedasticity defined	59
11	OLS with Heteroskedastic Errors	60
12	Heteroskedasticity-Robust Inference	63
13	GLS estimation	65
14	FGLS estimation	71
	14.1 Whites HCSE	72

# 1. Introduction

- In the following I will treat a number of subjects closely related to the standard multiple regression models in Econometrics.
- These subjects are more or less thoroughly examined in other courses. The major difference is that I almost exclusively use the notation of matrix algebra here. However the first half of the present slide set is perhaps more or less well known to most of you.
- There are several reasons to make yourself well acquainted with this:
  - ◆ Linear models are easier to work with in matrix notation.
  - ◆ A lot of scientific papers use this notation.
  - ◆ Several programming languages are matrix-oriented, which ease the implementation of calculations.

## 2. The Multiple Regression Model

- The general linear multiple regression model with  $k$  **explanatory** variables and  $n$  observations is fundamentally a system of  $n$  equations:

$$y_1 = \beta_1 + \beta_2 x_{12} + \dots + \beta_k x_{1k} + u_1,$$

$$y_2 = \beta_1 + \beta_2 x_{22} + \dots + \beta_k x_{2k} + u_2,$$

$$y_3 = \beta_1 + \beta_2 x_{32} + \dots + \beta_k x_{3k} + u_3,$$

$$\vdots$$

$$y_n = \beta_1 + \beta_2 x_{n2} + \dots + \beta_k x_{nk} + u_n.$$

- with  $u_1, \dots, u_n \sim \text{iid} - \mathcal{N}(0, \sigma^2)$ .

- If we construct the following vectors and matrices

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} \quad \text{and} \quad \mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix},$$

- and

$$\mathbf{X} = \begin{bmatrix} 1 & x_{12} & \cdots & x_{1k} \\ 1 & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n2} & \cdots & x_{nk} \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix},$$

- the regression model can be written as

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}$$

$$\Leftrightarrow \underset{n \times 1}{\mathbf{y}} = \underset{n \times k}{\mathbf{X}} \underset{k \times 1}{\boldsymbol{\beta}} + \underset{n \times 1}{\mathbf{u}}, \quad \mathbf{u} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n).$$

- This is perhaps known from the theory of linear normal models.

### 3. The Multivariate Normal Distribution

We introduce some concepts from multivariate normal (Gaussian) distributions.

$\mathbf{y}$  is Multivariate Normal Distributed with mean vector  $\boldsymbol{\mu}$  and non-singular covariance matrix  $\boldsymbol{\Sigma}$ .



$$\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}).$$



Multivariate density for  $\mathbf{y}$  is

$$f(\mathbf{y}) = \frac{1}{(2\pi)^{n/2} \sqrt{|\boldsymbol{\Sigma}|}} \exp \left( -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right).$$

Mean vector:

Covariance matrix:



## 4. Conditioning a multivariate normal distribution

Partitioned vectors and covariance matrix

$$\mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}.$$

If  $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  then the distribution of  $\mathbf{y}_2$  conditional on  $\mathbf{y}_1$  is multivariate normal with

$$\mathbf{y}_2 \mid \mathbf{y}_1 \sim \mathcal{N}(\boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}(\mathbf{y}_1 - \boldsymbol{\mu}_1), \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12}).$$

In particular **this is important if  $\mathbf{y}$  is a time series**, and for example  $\mathbf{y}_2$  corresponds to a future period and  $\mathbf{y}_1$  corresponds to the present. Then under normality we have a conditional normal distribution for the future observations given the present observations.

## 5. The Bivariate Normal Distribution

Special case: The bivariate case, that is  $\mathbf{y}_1 = y_1$  and  $\mathbf{y}_2 = y_2$  are both scalars,

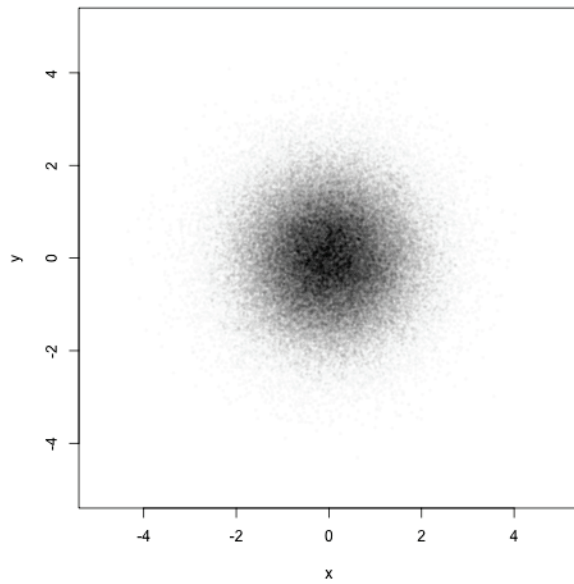
$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}.$$

Under normality  $y_2$  is conditionally normal given  $y_1$ :

$$y_2 | y_1 \sim \mathcal{N}\left(\mu_2 + \frac{\sigma_2}{\sigma_1}\rho(y_1 - \mu_1), (1 - \rho^2)\sigma_2^2\right),$$

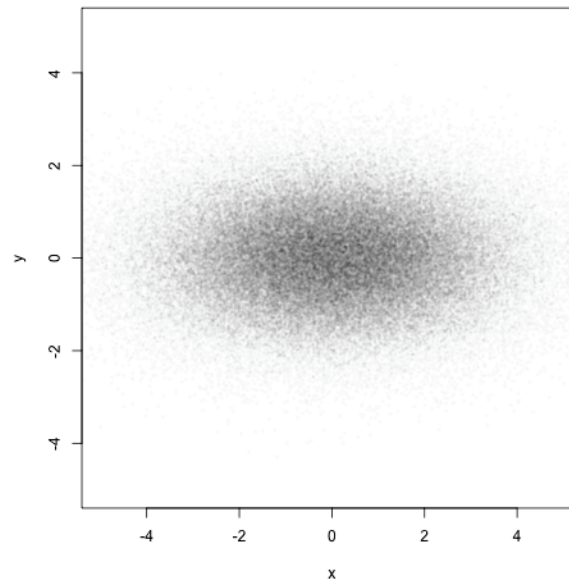
where  $\rho$  is the correlation coefficient between  $y_1$  and  $y_2$ .

► Examples for  $(x, y)$  bivariate normal distributed (no. 1).



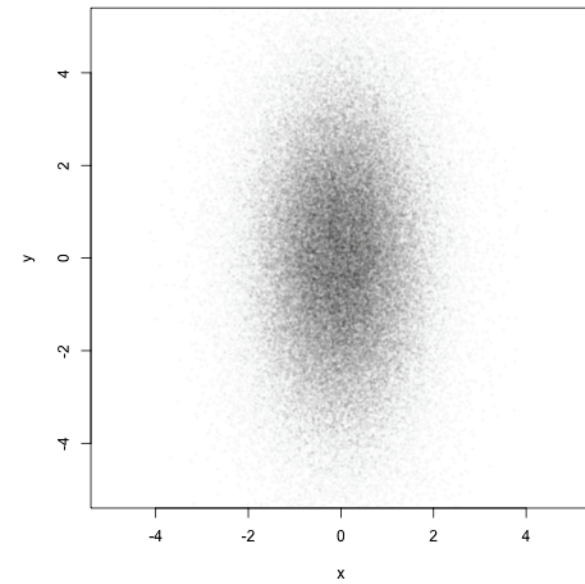
$$X \sim \mathcal{N}(0, 1), Y \sim \mathcal{N}(0, 1)$$

$$\rho = 0$$



$$X \sim \mathcal{N}(0, 2), Y \sim \mathcal{N}(0, 1)$$

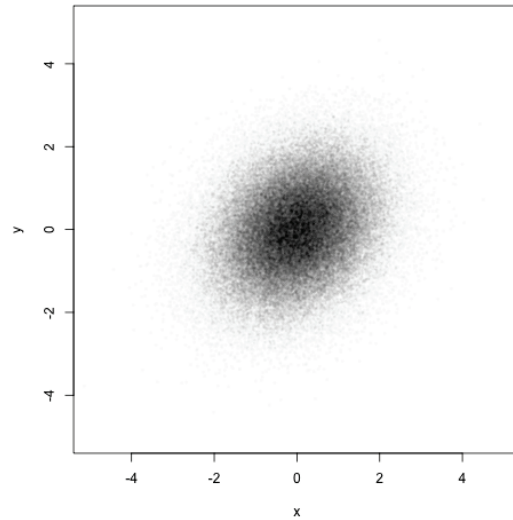
$$\rho = 0$$



$$X \sim \mathcal{N}(0, 1), Y \sim \mathcal{N}(0, 2)$$

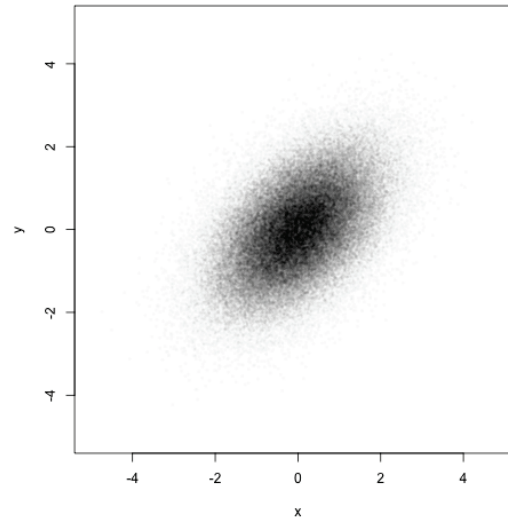
$$\rho = 0$$

► Examples for  $(x, y)$  bivariate normal distributed (no. 2).



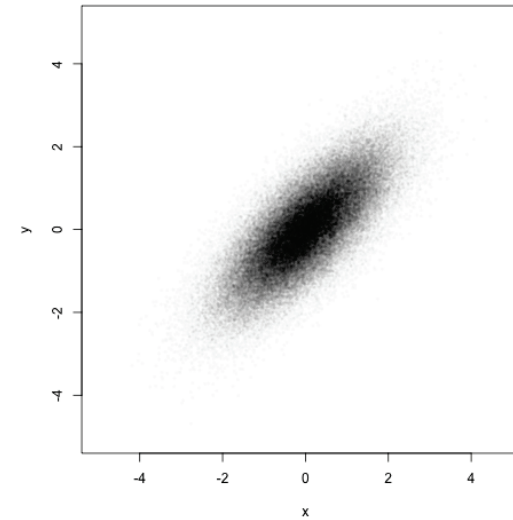
$$X \sim \mathcal{N}(0, 1), Y \sim \mathcal{N}(0, 1)$$

$$\rho = 0.25$$



$$X \sim \mathcal{N}(0, 1), Y \sim \mathcal{N}(0, 1)$$

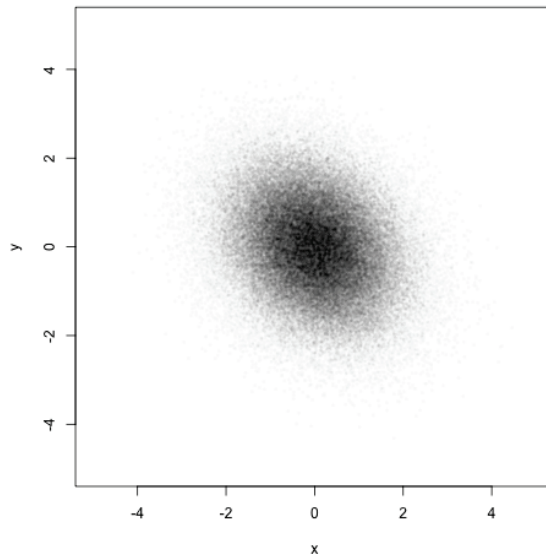
$$\rho = 0.5$$



$$X \sim \mathcal{N}(0, 1), Y \sim \mathcal{N}(0, 1)$$

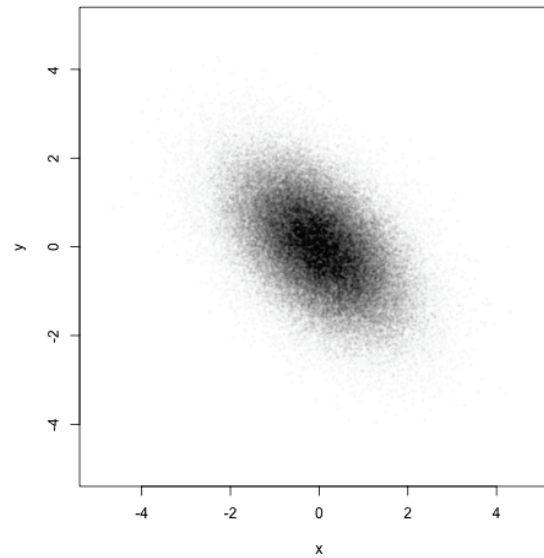
$$\rho = 0.75$$

► Examples for  $(x, y)$  bivariate normal distributed (no. 3).



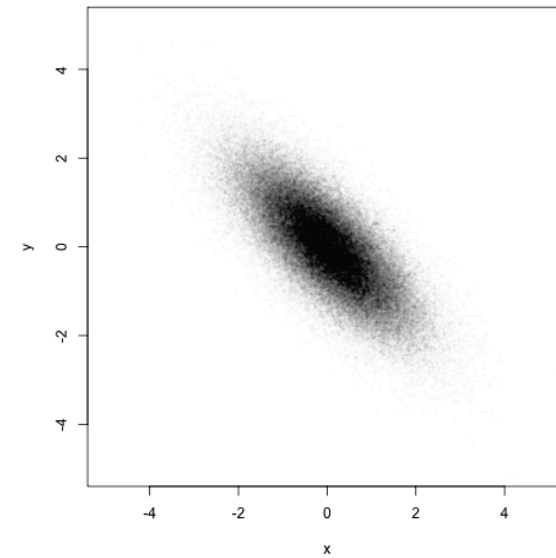
$$X \sim \mathcal{N}(0, 1), Y \sim \mathcal{N}(0, 1)$$

$$\rho = -0.25$$



$$X \sim \mathcal{N}(0, 1), Y \sim \mathcal{N}(0, 1)$$

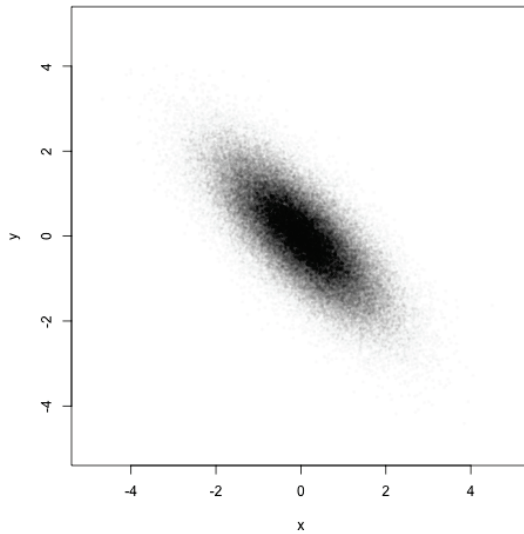
$$\rho = -0.5$$



$$X \sim \mathcal{N}(0, 1), Y \sim \mathcal{N}(0, 1)$$

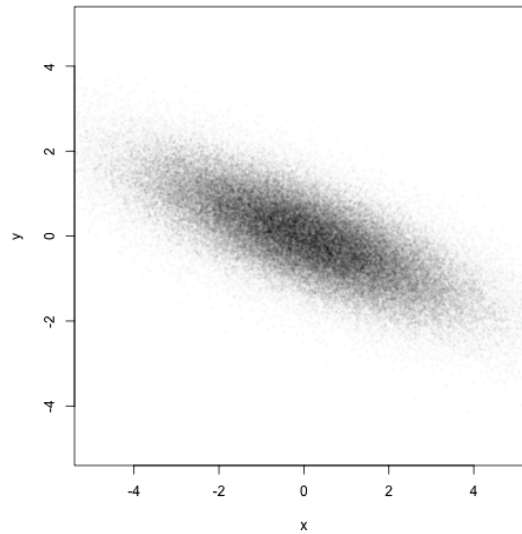
$$\rho = -0.75$$

► Examples for  $(x, y)$  bivariate normal distributed (no. 4).



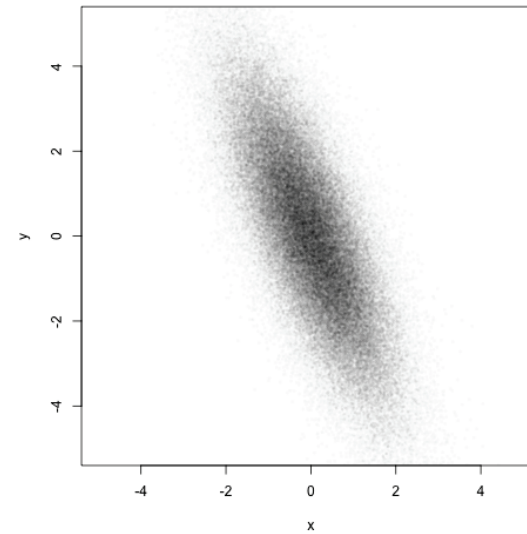
$$X \sim \mathcal{N}(0, 1), Y \sim \mathcal{N}(0, 1)$$

$$\rho = -0.75$$



$$X \sim \mathcal{N}(0, 2), Y \sim \mathcal{N}(0, 1)$$

$$\rho = -0.75$$



$$X \sim \mathcal{N}(0, 1), Y \sim \mathcal{N}(0, 2)$$

$$\rho = -0.75$$

# Classical Assumptions in Econometrics:

**A1.** zero conditional mean:  $E(\mathbf{u}|\mathbf{X}) = \mathbf{0}$

**A2.** No perfect collinearity:  $\text{rank}(\mathbf{X}) = k \Leftrightarrow |\mathbf{X}^\top \mathbf{X}| \neq 0$

**A3** Homoskedasticity and  
no serial correlation  $\text{var}(\mathbf{u}|\mathbf{X}) = \sigma^2 \mathbf{I}_n$

**A4.** Normality of errors:  $\mathbf{u} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$

where:

$$\sigma^2 \mathbf{I}_n = \sigma^2 \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}.$$

## The case of a General Covariance Matrix

- The general case where we may have **heteroscedasticity** and/or autocorrelation corresponds to cases where the covariance matrix of  $\mathbf{u}$  is no longer  $\sigma^2 \mathbf{I}_n$ . Instead the covariance matrix of  $\mathbf{u}$  has the following general form:

$$\begin{aligned} \text{var}(\mathbf{u}) = E(\mathbf{u}\mathbf{u}^\top) &= \begin{bmatrix} \text{var}(u_1) & \text{cov}(u_1 u_2) & \dots & \text{cov}(u_1 u_n) \\ \text{cov}(u_2 u_1) & \text{var}(u_2) & \dots & \text{cov}(u_2 u_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(u_n u_1) & \text{cov}(u_n u_2) & \dots & \text{var}(u_n) \end{bmatrix} \\ &= \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1n} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_n^2 \end{bmatrix} = \Sigma. \end{aligned}$$



## Example: Heteroscedasticity

- For so-called heteroscedasticity in the form:

$$\text{var}(u_i) = \sigma_i^2 = \sigma^2 z_i^2 \quad \text{for } i = 1, \dots, n,$$

and where  $z_1, \dots, z_n$  are some known values,

- the “General Covariance Matrix” then has the form

$$\text{var}(\mathbf{u}) = \mathbf{\Sigma} = \sigma^2 \begin{bmatrix} z_1 & 0 & \dots & 0 \\ 0 & z_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & z_n \end{bmatrix}.$$

## Example: Autocorrelation

- For 1st order autocorrelation in the form:

$$u_t = \rho u_{t-1} + \epsilon_t, \quad \epsilon_t \sim \text{nid}(0, \sigma_\epsilon^2) \quad \text{for } t = 1, \dots, T,$$

- it turns out that the “General Covariance Matrix” has the form

$$\text{var}(\mathbf{u}) = \mathbf{\Sigma} = \frac{\sigma_\epsilon^2}{1 - \rho^2} \begin{bmatrix} 1 & \rho & \dots & \rho^{T-2} & \rho^{T-1} \\ \rho & 1 & \dots & \rho^{T-3} & \rho^{T-2} \\ \vdots & \vdots & & \vdots & \vdots \\ \rho^{T-1} & \rho^{T-2} & \dots & \rho & 1 \end{bmatrix}.$$

- Exercise: Show that  $\mathbf{\Sigma}$  indeed has this form.

## 6. Least Squares and differentiation w.r.t. a vector

- The first task is to estimate the regression parameters:  $\beta_1, \dots, \beta_k$  by OLS.
- OLS minimizes the sum of squared residuals:

$$\min_{\beta_1, \beta_2, \dots, \beta_k} \sum_{i=1}^n u_i^2 = \min_{\beta_1, \beta_2, \dots, \beta_k} Q(\beta_1, \beta_2, \dots, \beta_k) = \min_{\boldsymbol{\beta}} Q(\boldsymbol{\beta}).$$

- We interpret the sum of squared errors as a function of the unknown parameter vector  $\boldsymbol{\beta}$ . By minimizing  $Q(\boldsymbol{\beta})$  wrt.  $\boldsymbol{\beta}$  we find the OLS estimate  $\hat{\boldsymbol{\beta}}$ .
- To do that we may differentiate  $Q(\boldsymbol{\beta})$  wrt.  $\boldsymbol{\beta}$ , put the 1st order derivatives equal to zero and solve the resulting  $k$  equations which we call the **normal equations**.
- The second order differentiation w.r.t.  $\boldsymbol{\beta}$  produces a square matrix (a Hessian) which should be positive definite in order for the normal equations to define a minimum.

- We use the notation:  $\mathbf{u}^\top = \begin{bmatrix} u_1 & u_2 & \cdots & u_n \end{bmatrix}$ , where  $\mathbf{u}$  is the  $n \times 1$  column vector of residuals.
- Obviously by using the definition of inner product

$$\begin{aligned} \mathbf{u}^\top \mathbf{u} &= u_1 u_1 + u_2 u_2 + \cdots + u_n u_n \\ &= \sum_{i=1}^n u_i^2 = Q(\boldsymbol{\beta}). \end{aligned}$$

- which motivates the need to calculate a derivative like

$$\frac{\partial \mathbf{u}^\top \mathbf{u}}{\partial \boldsymbol{\beta}}.$$

- By definition:  $\underset{n \times 1}{\mathbf{u}} = \underset{n \times 1}{\mathbf{y}} - \underset{n \times k}{\mathbf{X}} \underset{k \times 1}{\boldsymbol{\beta}}$ . If we insert this we obtain

$$\begin{aligned} Q(\boldsymbol{\beta}) &= \mathbf{u}^\top \mathbf{u} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &= \mathbf{y}^\top \mathbf{y} - \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} \\ &= \mathbf{y}^\top \mathbf{y} - 2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}, \end{aligned}$$

- and since  $\frac{\partial Q(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \frac{\partial \mathbf{y}^\top \mathbf{y}}{\partial \boldsymbol{\beta}} - \frac{\partial 2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y}}{\partial \boldsymbol{\beta}} + \frac{\partial \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}}{\partial \boldsymbol{\beta}}$ ,
- we obviously have to know how to calculate these three scalar functions of  $\boldsymbol{\beta}$ :

$$\frac{\partial \mathbf{y}^\top \mathbf{y}}{\partial \boldsymbol{\beta}}, \quad \frac{\partial 2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y}}{\partial \boldsymbol{\beta}} \quad \text{and} \quad \frac{\partial \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}}{\partial \boldsymbol{\beta}}.$$

## 6.1. Some formulas of Matrix Differentiation

- Consider the case where we want to differentiate a scalar function of a vector,  $f(\boldsymbol{\delta})$ , with respect to that very same vector  $\boldsymbol{\delta}$ .  
 $k \times 1$

- Then the rule is:
$$\frac{\partial f(\boldsymbol{\delta})}{\partial \boldsymbol{\delta}} = \begin{bmatrix} \frac{\partial f(\boldsymbol{\delta})}{\partial \delta_1} \\ \frac{\partial f(\boldsymbol{\delta})}{\partial \delta_2} \\ \vdots \\ \frac{\partial f(\boldsymbol{\delta})}{\partial \delta_k} \end{bmatrix}.$$

- From the previous we have three “relevant” forms of scalar functions of  $\boldsymbol{\beta}$ :

$$f(\boldsymbol{\beta}) = \mathbf{y}^\top \mathbf{y}, \quad f(\boldsymbol{\beta}) = \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y} = \mathbf{y}^\top \mathbf{X} \boldsymbol{\beta} \quad \text{and} \quad f(\boldsymbol{\beta}) = \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}.$$

- I.e. a constant, a linear function, and a quadratic form.

1. The constant function is

$$f(\boldsymbol{\beta}) = \mathbf{y}^\top \mathbf{y}.$$

Therefore we have

$$\frac{\partial \mathbf{y}^\top \mathbf{y}}{\partial \boldsymbol{\beta}} = \begin{bmatrix} \frac{\partial(\sum_{i=1}^n y_i^2)}{\partial \beta_1} \\ \frac{\partial(\sum_{i=1}^n y_i^2)}{\partial \beta_2} \\ \vdots \\ \frac{\partial(\sum_{i=1}^n y_i^2)}{\partial \beta_k} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \mathbf{0}_{k \times 1}.$$

2. The linear function is

$$f(\boldsymbol{\beta}) = \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y} = \mathbf{y}^\top \mathbf{X} \boldsymbol{\beta}.$$

This is because  $\mathbf{X}^\top \mathbf{y}$  is a  $k \times 1$  vector (call it  $\mathbf{a}$  for example),

we have with  $\mathbf{a} = \begin{bmatrix} a_1 & a_2 & \dots & a_n \end{bmatrix}^\top$

$$f(\boldsymbol{\beta}) = \boldsymbol{\beta}^\top \mathbf{a} = \beta_1 a_1 + \beta_2 a_2 + \dots + \beta_k a_k = \mathbf{a}^\top \boldsymbol{\beta}.$$

So

$$\frac{\partial(\boldsymbol{\beta}^\top \mathbf{a})}{\partial \boldsymbol{\beta}} = \frac{\partial(\mathbf{a}^\top \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \begin{bmatrix} \frac{\partial(\beta_1 a_1 + \beta_2 a_2 + \dots + \beta_k a_k)}{\partial \beta_1} \\ \frac{\partial(\beta_1 a_1 + \beta_2 a_2 + \dots + \beta_k a_k)}{\partial \beta_2} \\ \vdots \\ \frac{\partial(\beta_1 a_1 + \beta_2 a_2 + \dots + \beta_k a_k)}{\partial \beta_k} \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_k \end{bmatrix} = \mathbf{a} = \mathbf{X}^\top \mathbf{y}.$$



### 3. The quadratic form is

$$f(\boldsymbol{\beta}) = \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}.$$

$\mathbf{X}^\top \mathbf{X}$  is a symmetric  $k \times k$  matrix. Call it  $\mathbf{A}$ . Then we have

$$\begin{aligned} f(\boldsymbol{\beta}) &= \underset{1 \times k}{\boldsymbol{\beta}^\top} \underset{k \times k}{\mathbf{A}} \underset{k \times 1}{\boldsymbol{\beta}} \\ &= \begin{bmatrix} \beta_1 & \beta_2 & \cdots & \beta_k \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1k} \\ a_{21} & a_{22} & \cdots & a_{2k} \\ \vdots & \vdots & \ddots & \\ a_{1k} & a_{2k} & \cdots & a_{kk} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix}. \end{aligned}$$

This derivative in this case is best understood via a simple illustration.

Therefore, let us calculate

$$\begin{aligned}
 f(\boldsymbol{\beta}) &= \underset{1 \times 3}{\boldsymbol{\beta}^\top} \underset{3 \times 3}{\mathbf{A}} \underset{3 \times 1}{\boldsymbol{\beta}} = \begin{bmatrix} \beta_1 & \beta_2 & \beta_3 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{12} & a_{22} & a_{23} \\ a_{13} & a_{23} & a_{33} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = \\
 &\begin{bmatrix} a_{11}\beta_1 + a_{12}\beta_2 + a_{13}\beta_3 & a_{12}\beta_1 + a_{22}\beta_2 + a_{23}\beta_3 & a_{13}\beta_1 + a_{23}\beta_2 + a_{33}\beta_3 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} \\
 &= a_{11}\beta_1^2 + a_{12}\beta_1\beta_2 + a_{13}\beta_1\beta_3 + a_{12}\beta_2\beta_1 + a_{22}\beta_2^2 + a_{23}\beta_2\beta_3 + a_{13}\beta_3\beta_1 + a_{23}\beta_3\beta_2 + a_{33}\beta_3^2 \\
 &= a_{11}\beta_1^2 + a_{22}\beta_2^2 + a_{33}\beta_3^2 + 2a_{12}\beta_1\beta_2 + 2a_{13}\beta_1\beta_3 + 2a_{23}\beta_3\beta_2.
 \end{aligned}$$

So,

$$\begin{aligned}
 \frac{\partial \boldsymbol{\beta}^\top \mathbf{A} \boldsymbol{\beta}}{\partial \boldsymbol{\beta}} &= \begin{bmatrix} \frac{\partial f(\boldsymbol{\beta})}{\partial \beta_1} \\ \frac{\partial f(\boldsymbol{\beta})}{\partial \beta_2} \\ \frac{\partial f(\boldsymbol{\beta})}{\partial \beta_3} \end{bmatrix} = \begin{bmatrix} 2(a_{11}\beta_1 + a_{12}\beta_2 + a_{13}\beta_3) \\ 2(a_{22}\beta_2 + a_{12}\beta_1 + a_{23}\beta_3) \\ 2(a_{33}\beta_3 + a_{13}\beta_1 + a_{23}\beta_2) \end{bmatrix} \\
 &= \begin{bmatrix} 2(a_{11}\beta_1 + a_{12}\beta_2 + a_{13}\beta_3) \\ 2(a_{12}\beta_1 + a_{22}\beta_2 + a_{23}\beta_3) \\ 2(a_{13}\beta_1 + a_{23}\beta_2 + a_{33}\beta_3) \end{bmatrix} \\
 &= 2 \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{12} & a_{22} & a_{23} \\ a_{13} & a_{23} & a_{33} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = 2\mathbf{A}\boldsymbol{\beta}.
 \end{aligned}$$

Hence  $\frac{\partial(\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = 2\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}.$

- So much tedious algebra to convince ourselves that

$$\frac{\partial \mathbf{y}^\top \mathbf{y}}{\partial \boldsymbol{\beta}} = \mathbf{0}, \quad \frac{\partial \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y}}{\partial \boldsymbol{\beta}} = \mathbf{X}^\top \mathbf{y}, \quad \text{and} \quad \frac{\partial \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}}{\partial \boldsymbol{\beta}} = 2\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}.$$

- Hence we must have:

$$\frac{\partial Q(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \frac{\partial \mathbf{y}^\top \mathbf{y}}{\partial \boldsymbol{\beta}} - 2 \frac{\partial \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y}}{\partial \boldsymbol{\beta}} + \frac{\partial \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}}{\partial \boldsymbol{\beta}} = -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}.$$

- And therefore the normal equations are

$$\frac{\partial Q(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \mathbf{0} \Leftrightarrow \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^\top \mathbf{y}.$$

- There exist three **equivalent** necessary and sufficient conditions that ensure that  $\mathbf{X}^\top \mathbf{X}$  is invertible:
  1. The columns in  $\mathbf{X}$  are linearly independent.
  2.  $\mathbf{X}^\top \mathbf{X}$  has full rank, that is  $\rho(\mathbf{X}^\top \mathbf{X}) = k$ .
  3.  $|\mathbf{X}^\top \mathbf{X}| \neq 0$ .
- If these conditions are not met we are faced with so-called perfect collinearity.
- If however the three conditions are fulfilled  $(\mathbf{X}^\top \mathbf{X})^{-1}$  is well defined and the OLS estimate of  $\boldsymbol{\beta}$  is

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

## 6.2. The distribution of $\hat{\beta}$

- Note that  $\hat{\beta}$  is also the MLE of  $\beta$ .
- Given  $X$ ,  $\hat{\beta}$  is a linear function of  $y$ . From the theory of linear normal models we then know that  $\hat{\beta}$  is normally distributed.
- It is easily shown that

$E[\hat{\beta}|X] = \beta$ , therefore  $\hat{\beta}$  is an unbiased estimator of  $\beta$ .

- Furthermore it is easily shown that the variance matrix of  $\hat{\beta}$  is

$$\text{var}(\hat{\beta}|\mathbf{X}) = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}.$$

- Hence  $\hat{\boldsymbol{\beta}}$  is multivariate normal distributed according to

$$\hat{\boldsymbol{\beta}} \sim N_k(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}).$$

- This fact is used to test linear hypotheses on  $\boldsymbol{\beta}$ .
- We estimate  $\sigma^2$  by the unbiased estimator  $s^2$ :

$$s^2 = \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{n - k} = \frac{\hat{\mathbf{u}}^\top \hat{\mathbf{u}}}{n - k}.$$

- Then

$$\frac{\hat{\beta}_j - \beta_j}{s\sqrt{x^{jj}}} \sim t(n - k),$$

where  $x^{jj}$  is the  $j$ th diagonal element of  $(\mathbf{X}^\top \mathbf{X})^{-1}$ .



## 6.3. The Gauss-Markov theorem

- This theorem states that the OLS estimate  $\hat{\beta}$  is BLUE, (*Best Linear Unbiased Estimator*).

**BLUE:**

1. Unbiased:  $E[\hat{\beta}] = \beta$ .
2. Efficient (min. variance):  $\text{var}(\hat{\beta})$  smallest possible.
3. Linear function of the  $y$ 's.

2. Says that the estimate has the smallest<sup>a</sup> variance among all possible unbiased linear estimators.

---

<sup>a</sup> A square matrix  $A$  is here defined to be smaller than another square matrix  $B$  if  $B - A$  is positive definite.

## 7. The decomposition of sums of squares

- We can construct the usual decomposition of the variation of  $\mathbf{y}$ , since

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

where  $\hat{\mathbf{u}}$  is the residual vector. Then

$$\begin{aligned}\mathbf{y}^\top \mathbf{y} &= (\mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\mathbf{u}})^\top (\mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\mathbf{u}}) \\ &= \hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} + 2\hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \hat{\mathbf{u}} + \hat{\mathbf{u}}^\top \hat{\mathbf{u}}.\end{aligned}$$

► Since

$$\hat{\beta}^\top X^\top \hat{u} = \hat{\beta}^\top X^\top (y - X\hat{\beta})$$

$$= ((X^\top X)^{-1} X^\top y)^\top X^\top (y - X(X^\top X)^{-1} X^\top y)$$

$$= y^\top X(X^\top X)^{-1} X^\top (y - X(X^\top X)^{-1} X^\top y)$$

$$= y^\top X(X^\top X)^{-1} X^\top y - y^\top X(X^\top X)^{-1} X^\top X(X^\top X)^{-1} X^\top y$$

$$= y^\top X(X^\top X)^{-1} X^\top y - y^\top X(X^\top X)^{-1} X^\top y = 0,$$

► we have

$$y^\top y = \hat{\beta}^\top X^\top X \hat{\beta} + \hat{u}^\top \hat{u}.$$

- This only provides us with the sum of squares, but what we want is the squared deviations from averages. This is obtained by subtracting  $n\bar{y}^2 = \frac{1}{n}(\mathbf{1}^\top \mathbf{y})^2$  on both sides:

$$\mathbf{y}^\top \mathbf{y} - \frac{1}{n}(\mathbf{1}^\top \mathbf{y})^2 = \hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} - \frac{1}{n}(\mathbf{1}^\top \mathbf{y})^2 + \hat{\mathbf{u}}^\top \hat{\mathbf{u}}$$

Total variation	=	Explained variation	+	Unexplained variation
<b>SS Total</b>		<b>SS Explained</b>		<b>SS Residuals</b>

Here:

**SSE**: The part of the variation in  $y$ , which is explained by  $\mathbf{X}$ . **SSR**: The part of the variation in  $y$ , which is not explained by the model.

- The coefficient of determination is now defined as:

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS}$$

- and the adjusted coefficient of determination

$$R_{\text{adj}}^2 = 1 - \frac{n-1}{n-k} (1 - R^2)$$

- Note the intuition in the definition of  $R_{\text{adj}}^2$ : Mean sums of squares are used instead of sums of squares.

## 8. Geometric interpretation of OLS

- Given the OLS estimators  $\hat{\beta}$ , the *predicted* value of  $\mathbf{y}$  is defined by

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \underbrace{\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}_{\mathbf{P}} \mathbf{y} = \mathbf{P}\mathbf{y}.$$

- The residuals are defined as

$$\hat{\mathbf{u}} = \mathbf{y} - \mathbf{X}\hat{\beta} = \mathbf{y} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \underbrace{[\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top]}_{\mathbf{M}} \mathbf{y} = \mathbf{M}\mathbf{y} = \mathbf{M}\mathbf{u},$$

- Where

$$\mathbf{P} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \quad \text{and} \quad \mathbf{M} = \mathbf{I} - \mathbf{P}$$

are the so-called fundamental OLS matrices.

- They are **idempotent** (and symmetric and hence by definition projection matrices) and satisfy the equations  $\mathbf{P}\mathbf{X} = \mathbf{X}$  and  $\mathbf{M}\mathbf{X} = \mathbf{0}$ , and

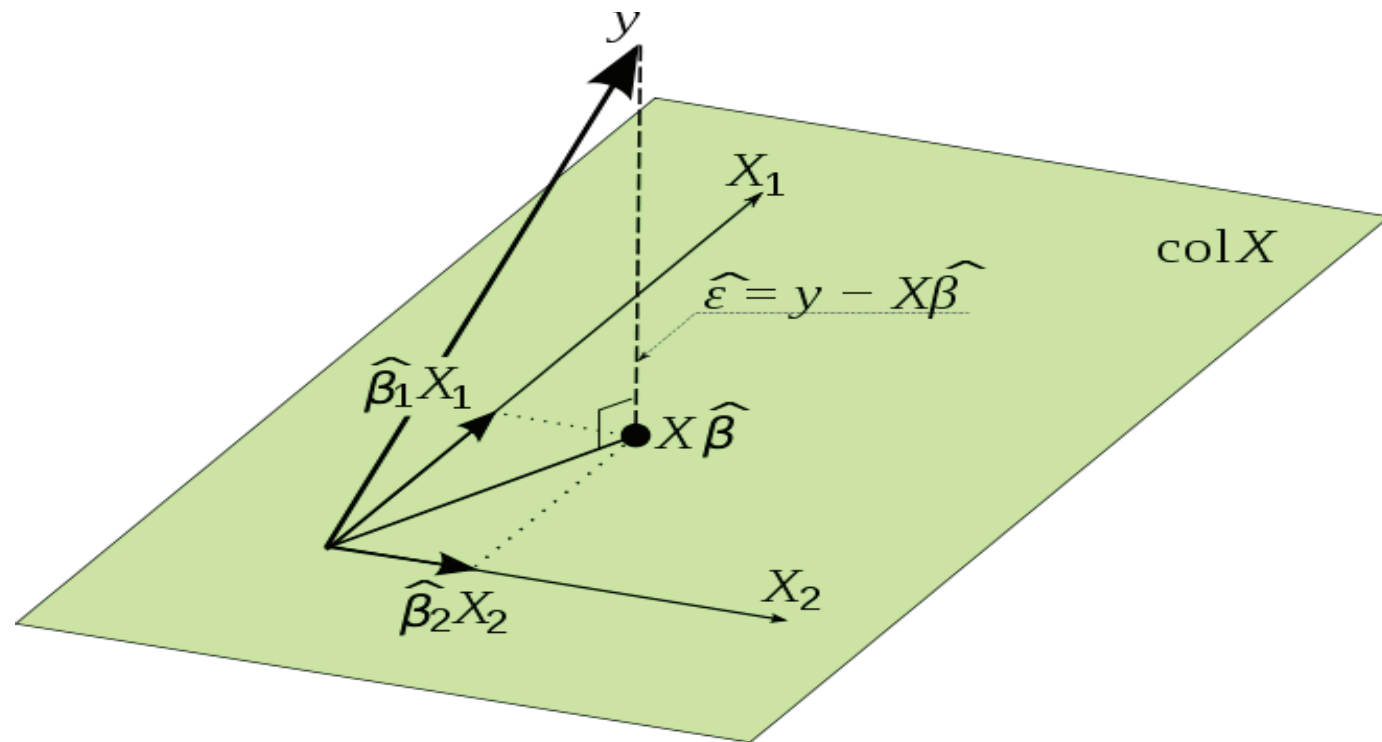
$$\begin{aligned}\hat{\mathbf{y}}^\top \hat{\mathbf{u}} &= (\mathbf{P}\mathbf{y})^\top [\mathbf{M}\mathbf{y}] = \mathbf{y}^\top \mathbf{P}^\top (\mathbf{I} - \mathbf{P})\mathbf{y} \\ &= \mathbf{y}^\top (\mathbf{P}^\top - \mathbf{P}^2)\mathbf{y} \\ &= \mathbf{y}^\top (\mathbf{P} - \mathbf{P})\mathbf{y} \\ &= 0,\end{aligned}$$

which shows that the residuals and the predicted responses are orthogonal.

- Further it can be seen that

$$\begin{aligned}\mathbf{y} &= \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\mathbf{u}} \\ &= \mathbf{P}\mathbf{y} + \mathbf{M}\mathbf{y}.\end{aligned}$$

► The projection related to OLS:





- The previous and the fact that both  $\mathbf{P}$  and  $\mathbf{M}$  are projection matrices with  $\mathbf{P} + \mathbf{M} = \mathbf{I}$ , show that OLS geometrically can be interpreted as a decomposition of  $\mathbf{y}$  in two orthogonal components.
  - ◆  $\mathbf{P}$  projects  $\mathbf{y}$  on the (*hyper*-)plane , which is spanned by the columns in  $\mathbf{X}$ , and
  - ◆  $\mathbf{M}$  projects  $\mathbf{y}$  on the (*hyper*-)plane, which is orthogonal to the first (*hyper*-)plane.
- Thus,  $\mathbf{y}$  may be seen as a sum of two orthogonal vectors, that constitute the sides in a right-angled triangle enclosing the right angle.

## 9. Linear Restrictions

- We want to apply some linear restrictions to the model:

$$r_{11}\beta_1 + r_{12}\beta_2 + \dots + r_{1k}\beta_k = r_1$$

$$r_{21}\beta_1 + r_{22}\beta_2 + \dots + r_{2k}\beta_k = r_2$$

$$\vdots$$

$$r_{q1}\beta_1 + r_{q2}\beta_2 + \dots + r_{qk}\beta_k = r_q.$$

- In matrix notation this is written  $\underset{q \times k}{\mathbf{R}} \underset{k \times 1}{\boldsymbol{\beta}} = \underset{q \times 1}{\mathbf{r}}'$ ,  
where the elements of  $\mathbf{R}$  gives the restrictions.

## 9.1. Estimation

- We want to estimate the model under such restrictions.
- Each row in  $\mathbf{R}$  is a linear restriction of the coefficient vector.
- Typically  $\mathbf{R}$  will have few rows and many zeros in each row.
- **Examples:**
  - One of the coefficients is 0,  $\beta_j = 0$ :

$$\mathbf{R} = \begin{bmatrix} 0 & 0 & \dots & 1 & 0 & \dots & 0 \end{bmatrix} \quad \text{and} \quad \mathbf{r} = 0.$$

- Two of the coefficients are equal,  $\beta_k = \beta_j$ :

$$\mathbf{R} = \begin{bmatrix} 0 & 0 & 1 & \dots & -1 & \dots & 0 \end{bmatrix} \quad \text{and} \quad \mathbf{r} = 0.$$

► **More examples:**

- Some of the coefficients add up to 1,  $\beta_2 + \beta_3 + \beta_4 = 1$ :

$$\mathbf{R} = \begin{bmatrix} 0 & 1 & 1 & 1 & 0 & \dots & 0 \end{bmatrix} \quad \text{and} \quad \mathbf{r} = 1.$$

- Some of the coefficients are 0,  $\beta_1 = \beta_2 = \beta_4 = 0$ :

$$\mathbf{R} = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 1 & \dots & 0 \end{bmatrix} \quad \text{and} \quad \mathbf{r} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

- To estimate the model under the restrictions we set out the Lagrange function to minimizing:

$$\begin{aligned} L(\lambda, \beta) &= (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \lambda^\top (\mathbf{R}\beta - \mathbf{r}) \\ &= \mathbf{y}^\top \mathbf{y} - 2\beta^\top \mathbf{X}^\top \mathbf{y} + \beta^\top \mathbf{X}^\top \mathbf{X}\beta + \lambda^\top (\mathbf{R}\beta - \mathbf{r}). \end{aligned}$$

- The first order conditions are

$$\frac{\partial L(\lambda, \beta)}{\partial \beta} = -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\beta + \mathbf{R}^\top \lambda = 0,$$

$$\frac{\partial L(\lambda, \beta)}{\partial \lambda} = \mathbf{R}\beta - \mathbf{r} = 0.$$

- Let  $\hat{\hat{\beta}}$  denote the estimate of  $\beta$  under the restrictions.

- From the first set of restrictions we get

$$\boldsymbol{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{y} - \frac{1}{2} \mathbf{R}^\top \boldsymbol{\lambda}).$$

- Inserting in the last set of restrictions reveals

$$\mathbf{r} = \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} - \frac{1}{2} \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \boldsymbol{\lambda},$$

- which again may be solved for  $\boldsymbol{\lambda}$ :

$$\boldsymbol{\lambda} = 2 \left[ \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \right]^{-1} \left[ \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} - \mathbf{r} \right].$$

- This can be used to find  $\hat{\hat{\beta}}$ :

$$\begin{aligned}\hat{\hat{\beta}} &= (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{y} - \frac{1}{2} \mathbf{R}^\top \underbrace{2 \left[ \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \right]^{-1} \left[ \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} - \mathbf{r} \right]}_{\lambda}) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \left[ \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \right]^{-1} \left[ \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} - \mathbf{r} \right] \\ &= \hat{\beta} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \left[ \mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top \right]^{-1} \left[ \mathbf{R} \hat{\beta} - \mathbf{r} \right].\end{aligned}$$

- It is worth noting that if  $\hat{\hat{\beta}}$  fulfills the restrictions then

$$\mathbf{R} \hat{\hat{\beta}} - \mathbf{r} = 0 \quad \text{which implies} \quad \hat{\hat{\beta}} = \hat{\beta}.$$

## 9.2. Test linear restrictions

- We will see how to test if the restrictions are valid.
- Thus we want to test the null hypothesis:  $H_0 : \mathbf{R}\boldsymbol{\beta} = \mathbf{r}$ .
- This is done under assumption A4 via the following F-statistic

$$F = \frac{(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})^\top [\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top]^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}) / q}{\hat{\mathbf{u}}^\top \hat{\mathbf{u}} / (n - k)} \sim F(q, n - k)$$

- If A4 is not fulfilled, then we must use the so-called Wald statistic:

$$W = \frac{(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})^\top [\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top]^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})}{\hat{\mathbf{u}}^\top \hat{\mathbf{u}} / (n - k)} \sim \chi^2(q) \text{ (approximately)}$$



- Note that the F-statistic is exactly equal to:

$$F = \frac{(\hat{\mathbf{u}}_*^\top \hat{\mathbf{u}}_* - \hat{\mathbf{u}}^\top \hat{\mathbf{u}}) / q}{\hat{\mathbf{u}}^\top \hat{\mathbf{u}} / (n - k)} \sim F(q, n - k),$$

- where

$\hat{\mathbf{u}}_*^\top \hat{\mathbf{u}}_*$  is the sum of the squared residuals from the restricted model, and

$\hat{\mathbf{u}}^\top \hat{\mathbf{u}}$  is the sum of squared residuals from the unrestricted model.

- A nice feature of using the matrix version on the previous slide is that you only need the estimates from one regression, namely the unrestricted one.

## 9.3. The Chow Forecast test

- The idea is that, if the parameter vector is constant, then we have a specific confidence that *out-of-sample* predictions will fall within specified bounds.
  - ◆ these bounds are the well-known prediction intervals computed from sample data.
- Hence, large prediction errors will be critical for the constancy hypothesis.

- To test this hypothesis we split the sample in two.
  - ✴ Use  $n_1$  observations for estimation, and
  - ✴ use  $n_2 = n - n_1$  for testing.
- 1. For time series we would usually take the first  $n_1$  for estimation.
- 2. In cross-sections we could e.g. split the sample according to a size variable.
- It will be appropriate to reserve 5-15 percent of the observations for testing.

- In splitting the sample, we use the following notation:

$$\begin{aligned} \mathbf{y} &= \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} \leftarrow \begin{matrix} n_1 \times 1 \\ n_2 \times 1 \end{matrix} & \mathbf{X} &= \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix} \leftarrow \begin{matrix} n_1 \times k \\ n_2 \times k \end{matrix} \\ \mathbf{u} &= \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{bmatrix} \leftarrow \begin{matrix} n_1 \times 1 \\ n_2 \times 1 \end{matrix} \end{aligned}$$

- The complete model can then be written as

$$\begin{aligned} \mathbf{y}_1 &= \mathbf{X}_1 \boldsymbol{\beta} + \mathbf{u}_1, \\ \mathbf{y}_2 &= \mathbf{X}_2 \boldsymbol{\alpha} + \mathbf{u}_2. \end{aligned}$$

- The null hypothesis is

$$H_0 : \boldsymbol{\alpha} = \boldsymbol{\beta}.$$

- To use a dummy variable approach write the second equation as

$$\begin{aligned}y_2 &= X_2\alpha + u_2 = X_2\alpha + u_2 + X_2\beta - X_2\beta \\&= X_2\beta + X_2(\alpha - \beta) + u_2 \\&= X_2\beta + \gamma + u_2.\end{aligned}$$

- Hence if  $\gamma = 0$  then  $\alpha = \beta$ .
- Note that  $X_2(\alpha - \beta)$  is a  $n_2 \times 1$  “parameter-vector”.
- Now the model may be written compactly as

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \underbrace{\begin{bmatrix} X_1 & 0 \\ X_2 & I_{n_2} \end{bmatrix}}_Z \begin{bmatrix} \beta \\ \gamma \end{bmatrix} + \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}.$$

► Now it is simple to estimate  $\begin{bmatrix} \boldsymbol{\beta} & \boldsymbol{\gamma} \end{bmatrix}^\top$ :

$$\begin{aligned} \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\gamma}} \end{bmatrix} &= (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{y} \\ &= \left( \begin{bmatrix} \mathbf{X}_1 & \mathbf{0} \\ \mathbf{X}_2 & \mathbf{I}_{n_2} \end{bmatrix}^\top \begin{bmatrix} \mathbf{X}_1 & \mathbf{0} \\ \mathbf{X}_2 & \mathbf{I}_{n_2} \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{X}_1 & \mathbf{0} \\ \mathbf{X}_2 & \mathbf{I}_{n_2} \end{bmatrix}^\top \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} \\ &= \begin{bmatrix} (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{y}_1 \\ \mathbf{y}_2 - \mathbf{X}_2 (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{y}_1 \end{bmatrix}. \end{aligned}$$

► Hence  $\hat{\boldsymbol{\beta}}$  estimate  $\boldsymbol{\beta}$  using  $n_1$  data points.

- Noting that

$$\mathbf{X}_2(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{y}_1 = \mathbf{X}_2 \hat{\boldsymbol{\beta}} = \hat{\mathbf{y}}_2,$$

we discover that  $\hat{\boldsymbol{\gamma}} = \mathbf{y}_2 - \hat{\mathbf{y}}_2$  is the prediction errors, when we are trying to predict  $\mathbf{y}_2$  using an estimate of  $\boldsymbol{\beta}$  based on  $n_1$  observations.

- Hence to test for parameter constancy we test

$$H_0 : \boldsymbol{\gamma} = \mathbf{0},$$

using our test for linear restrictions to the model with the dummy variable  $\begin{bmatrix} \mathbf{0} & \mathbf{I}_{n_2} \end{bmatrix}^\top$ .

- The restriction matrices are:

$$\mathbf{R} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{n_2} \end{bmatrix} \quad \text{and} \quad \mathbf{r} = \begin{bmatrix} \mathbf{0} \end{bmatrix}.$$

- Hence, the test is:

$$F = \frac{\left( \mathbf{R} \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\gamma}} \end{bmatrix} \right)^\top \left[ (\mathbf{R} \mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{R}^\top \right]^{-1} \left( \mathbf{R} \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\gamma}} \end{bmatrix} \right) / n_2}{\hat{\mathbf{u}}^\top \hat{\mathbf{u}} / \underbrace{(n_1 - k)}_{=n_1+n_2-(k+n_2)}} \sim F(n_2, n_1 - k).$$



► If A4 does not apply, then use:

$$W = \frac{\left( \mathbf{R} \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\gamma}} \end{bmatrix} \right)^\top [(\mathbf{R}\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{R}^\top]^{-1} \left( \mathbf{R} \begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\gamma}} \end{bmatrix} \right)}{\hat{\mathbf{u}}^\top \hat{\mathbf{u}} / (n_1 - k)} \sim \chi^2(n_2).$$

## 9.4. Chow's Breakpoint Test

- If the forecast subset is large enough it might be better to estimate two regression functions, one for each sub-sample, and then test for common parameters.
- Then the unrestricted model may be written:

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} X_1 & 0 \\ 0 & X_2 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} + \begin{bmatrix} u_1 \\ u_2 \end{bmatrix},$$

where  $\beta_1$  and  $\beta_2$  are the  $k$ -vectors of the two sub-samples.

- The null hypothesis of no structural break is then

$$H_0 : \beta_1 = \beta_2.$$

- Again, use the test for linear restrictions to this model. We just have to set up the restriction matrices:

$$\mathbf{R} = \begin{bmatrix} \mathbf{I}_k & -\mathbf{I}_k \end{bmatrix} \quad \text{and} \quad \mathbf{r} = \begin{bmatrix} \mathbf{0} \end{bmatrix}.$$

- Hence, the test is:

$$F = \frac{\left( \mathbf{R} \begin{bmatrix} \hat{\boldsymbol{\beta}}_1 \\ \hat{\boldsymbol{\beta}}_2 \end{bmatrix} \right)^\top [(\mathbf{R}\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{R}^\top]^{-1} \left( \mathbf{R} \begin{bmatrix} \hat{\boldsymbol{\beta}}_1 \\ \hat{\boldsymbol{\beta}}_2 \end{bmatrix} \right) / k}{\hat{\mathbf{u}}^\top \hat{\mathbf{u}} / (n - 2k)} \sim F(k, n - 2k).$$

- If A4 does not apply, then use the chi-squared alternative (the Wald test).
- It is not hard to modify this test to relax the restriction on the intercept term, or on the slope.

## 10. Heteroskedasticity Defined

- A general *Heteroskedastic* model violates A3.
- It can be presented as:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad \mathbf{u} \sim N(\mathbf{0}, \sigma^2 \boldsymbol{\Sigma}).$$

where

$$\begin{aligned} \sigma^2 \boldsymbol{\Sigma} = \text{var}(\mathbf{u}) &= \begin{bmatrix} \text{var}(u_1) & \text{cov}(u_1, u_2) & \dots & \text{cov}(u_1, u_n) \\ \text{cov}(u_2, u_1) & \text{var}(u_2) & \dots & \text{cov}(u_2, u_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(u_n, u_1) & \text{cov}(u_n, u_2) & \dots & \text{var}(u_n) \end{bmatrix} \\ &= E(\mathbf{u}\mathbf{u}^\top) \end{aligned}$$

## 11. OLS with Heteroskedastic Errors

➤ For a general covariance matrix the following results hold:

1. The OLS estimator  $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$  is unbiased and consistent.
2. The OLS estimator is **inefficient**.
3. The message is the following:
  - ◆ We can get a fair estimate of the parameter vector using OLS, even when A3 is not satisfied, but
  - ◆ We can not perform hypothesis testing on the parameter vector, because the standard errors estimated using OLS are **WRONG** (i.e. biased and inconsistent).

► The correct variance matrix of the OLS estimator is

$$\begin{aligned}\text{var}(\hat{\boldsymbol{\beta}}|\mathbf{X}) &= E \left[ (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top | \mathbf{X} \right] \\ &= E \left[ ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u}) ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u})^\top | \mathbf{X} \right] \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top E \left[ \mathbf{u} \mathbf{u}^\top | \mathbf{X} \right] \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \sigma^2 \boldsymbol{\Sigma} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Sigma} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}\end{aligned}$$

► Hence, tests based on  $\sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$  are **invalid**.

To sum up:

- The usual OLS t-statistics are **not** t-distributed under heteroskedasticity, and the problem can not be resolved by using large sample sizes. This is also the case for the F-test.
- One will usually underestimate  $SE(\hat{\beta}_{j\text{OLS}})$ , which implies that confidence bands get too short and t-statistics get too large.
- ➡ Thus, hypotheses which should be retained may be rejected.

## 12. Heteroskedasticity-Robust Inference

- In the **simple** regression model with heteroskedasticity:

$$y_i = \beta_0 + \beta_1 x_i + u_i, \quad \text{var}(u_i) = \sigma_i^2$$

we have an easy correction.

- If we don't know  $\sigma_i^2$ , then we may use:

$$\widehat{\text{var}(\hat{\beta}_{1\text{OLS}})}^{\text{HRSE}} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \hat{u}_i^2}{\left(\sum_{i=1}^n (x_i - \bar{x})^2\right)^2}.$$

- This estimator of the standard error is consistent.



- In the **multiple** regression model one should use

$$\widehat{\text{var}(\hat{\beta}_{j\text{OLS}})}^{\text{HRSE}} = \frac{\sum_{i=1}^n \hat{r}_{ij}^2 \hat{u}_i^2}{\left(\sum_{i=1}^n \hat{r}_{ij}^2\right)^2},$$

- where  $\hat{r}_{ij}$  are the residuals from the following regression:

$$x_{ji} = \delta_0 + \delta_1 x_{1i} + \dots + \delta_{j-1} x_{j-1,i} + \delta_{j+1} x_{j+1,i} + \dots + \delta_k x_{ki} + r_{ij},$$

- E.g. the residuals from the regression where  $x_j$  are regressed on the remaining explanatory variables.

- $\sqrt{\widehat{\text{var}(\hat{\beta}_{j\text{OLS}})}^{\text{HRSE}}}$  is known as the **heteroskedasticity-robust standard error** for  $\beta_j$ .

## 13. GLS Estimation

► We have something like:  $\sigma_i^2 = \sigma^2 z_i$  with  $z_i$  known.

► The trick is:

$$\frac{y_i}{\sqrt{z_i}} = \alpha \frac{1}{\sqrt{z_i}} + \beta \frac{x_i}{\sqrt{z_i}} + \frac{u_i}{\sqrt{z_i}}, \quad (1)$$

$$\Rightarrow \text{var} \left( \frac{u_i}{\sqrt{z_i}} \right) = \sigma^2 \quad \text{for all } i$$

► and now OLS is applicable on (1),  $\hat{\beta} = \hat{\beta}_{\text{GLS}}$  with:

$$\widehat{\text{var}(\hat{\beta}_{\text{GLS}})} = \frac{s^2}{(\sum^n x_i^2 / z_i) - (\sum^n x_i / z_i)^2 (\sum^n z_i^{-1})^{-1}}, \quad s^2 = \frac{SSE_{\text{GLS}}}{n - 2}.$$

► This method is also called **Weighted Least Squares WLS**.

- **Hence in general:** Least Squares estimation, when  $\Sigma$  is known, can be done by transforming the model such that the classical assumptions are fulfilled.
- As  $\Sigma$  is a positive definite symmetric matrix we can write it as

$$\Sigma = C\Lambda C^{\top}.$$

- By the calculus rules for eigen-values and -vectors we know

$$\begin{aligned}\Sigma^{-1} &= C\Lambda^{-1}C^{\top} = C\Lambda^{-1/2}\Lambda^{-1/2}C^{\top} = C\Lambda^{-1/2}(C\Lambda^{-1/2})^{\top} \\ &= P^{\top}P,\end{aligned}$$

and

$$\Sigma = (\Sigma^{-1})^{-1} = (P^{\top}P)^{-1} = P^{-1}(P^{\top})^{-1}.$$

## Important general Rule

$$\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \text{ and } \boldsymbol{\Sigma} \text{ is invertible}$$



$$\boldsymbol{\Sigma}^{-1/2} \mathbf{y} \sim \mathcal{N}(\boldsymbol{\Sigma}^{-1/2} \boldsymbol{\mu}, \mathbf{I}).$$

Here  $\boldsymbol{\Sigma}^{-1/2}$  is the matrix defined such that  $\boldsymbol{\Sigma}^{-1} = (\boldsymbol{\Sigma}^{-1/2})^\top \boldsymbol{\Sigma}^{-1/2}$ .  
Let  $\mathbf{P} = \boldsymbol{\Sigma}^{-1/2}$ . Then  $\boldsymbol{\Sigma}^{-1} = \mathbf{P}^\top \mathbf{P}$ . Exactly as defined on the previous slide.

- **Hence in general;** Least Squares estimation, **when  $\boldsymbol{\Sigma}$  is known,** can be done by transforming the model such that the classical assumptions (classical design criteria) on the covariance matrix are fulfilled. See next slide.

► Now use  $P$  to define:

$$y_* = Py, \quad X_* = PX, \quad u_* = Pu,$$

► and note that this implies that

$$\begin{aligned} \text{var}(u_* | X) &= E[u_* u_*^\top | X] = E[Pu(Pu)^\top | X] \\ &= E[Puu^\top P^\top | X] = PE[uu^\top | X]P^\top \\ &= P\sigma^2 \Sigma P^\top \\ &= \sigma^2 PP^{-1}(P^\top)^{-1}P^\top \\ &= \sigma^2 I. \end{aligned}$$

- Hence, the transformed errors satisfy the design criteria, and it is valid to apply OLS to the model

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

- Hence, the estimates of  $\boldsymbol{\beta}$  and  $\sigma^2$  are the usual ones:

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{\text{GLS}} &= (\mathbf{X}_*^\top \mathbf{X}_*)^{-1} \mathbf{X}_*^\top \mathbf{y}_* \\ &= ((\mathbf{P}\mathbf{X})^\top \mathbf{P}\mathbf{X})^{-1} (\mathbf{P}\mathbf{X})^\top \mathbf{P}\mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{P}^\top \mathbf{P}\mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}^\top \mathbf{P}\mathbf{y} \\ &= (\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{y}. \end{aligned}$$

$$\begin{aligned}s_{\text{GLS}}^2 &= \frac{(\mathbf{y}_* - \mathbf{X}_* \hat{\boldsymbol{\beta}}_{\text{GLS}})^\top (\mathbf{y}_* - \mathbf{X}_* \hat{\boldsymbol{\beta}}_{\text{GLS}})}{n - k} \\ &= \frac{(\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{GLS}})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{GLS}})}{n - k}\end{aligned}$$

➡ These are called **GLS estimators**.

➤ Variance of  $\hat{\boldsymbol{\beta}}_{\text{GLS}}$ :

$$\text{var}(\hat{\boldsymbol{\beta}}_{\text{GLS}}) = \sigma^2 (\mathbf{X}_*^\top \mathbf{X}_*)^{-1} = \sigma^2 (\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1}.$$

➤ and the estimated variance:

$$\widehat{\text{var}}(\hat{\boldsymbol{\beta}}_{\text{GLS}}) = s_{\text{GLS}}^2 (\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1}.$$

➤ In conclusion this shows that if  $\boldsymbol{\Sigma}$  is known, then it is easy to take care of a general covariance matrix for the disturbances.

## 14. FGLS estimation

- FGLS is short for **Feasible Generalized Least Squares**, and it is the method to use when  $\Sigma$  is **unknown**.
- In the general cases the situation is

$$y = X\beta + u, \quad u \sim N(0, V),$$

where  $V$  is the unknown covariance matrix of the disturbances.

**Note that the scaling factor  $\sigma^2$  is included in  $V$ .**

- Now if we can come up with a consistent estimate of  $V$ , call it  $\hat{V}$ , then we are home free:

$$\begin{aligned}\hat{\beta}_{\text{FGLS}} &= (X^\top \hat{V}^{-1} X)^{-1} X^\top \hat{V}^{-1} y, \\ \text{var}(\hat{\beta}_{\text{FGLS}}) &= (X^\top \hat{V}^{-1} X)^{-1}.\end{aligned}$$



## 14.1. Whites HCSE

- Halbert White suggested some simple heteroskedasticity consistent standard errors for the OLS estimator in the case where  $V$  has the following structure:

$$V = \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{bmatrix},$$

where  $\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$  are unknown.

- This estimator is exactly equal to the robust one we encountered previously.

- Remember that the correct variance matrix for the OLS estimator is

$$\text{var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{V} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}.$$

- Hence, to perform FGLS we need an estimate of

$$\begin{aligned} \mathbf{X}^\top \mathbf{V} \mathbf{X} &= \begin{bmatrix} \vdots & \vdots & & \vdots \\ \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_n \\ \vdots & \vdots & & \vdots \end{bmatrix} \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{bmatrix} \begin{bmatrix} \cdots & \mathbf{x}_1^\top & \cdots \\ \cdots & \mathbf{x}_2^\top & \cdots \\ \vdots & \vdots & \\ \cdots & \mathbf{x}_n^\top & \cdots \end{bmatrix} \\ &= \sum_{t=1}^n \sigma_t^2 \mathbf{x}_t \mathbf{x}_t^\top, \end{aligned}$$

where  $\mathbf{x}_t = (1, x_{2t}, \cdots, x_{kt})$  is the  $t$ 'th row of  $\mathbf{X}$ .

- The White estimator of  $\text{var}(\hat{\boldsymbol{\beta}})$ , then, is to replace  $\sigma_t^2$  with the squared  $t$ th residual  $\hat{u}_t^2$ , such that

$$\widehat{\text{var}(\hat{\boldsymbol{\beta}})}^{\text{HCSE}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \hat{\mathbf{V}} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}$$
$$= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \begin{bmatrix} \hat{u}_1^2 & 0 & \cdots & 0 \\ 0 & \hat{u}_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \hat{u}_n^2 \end{bmatrix} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}.$$

- In general, one should always report these standard errors with the regression output.

# Problem 1

- Verify the formulas on slide 9 from the notation on slides 6 and 8.
- Show that  $\Sigma$  on slide 17 has the form as postulated.
- Show that  $s^2$  from slide 31 is unbiased
- Show that the OLS estimator  $\hat{\beta}$  on slide 32 has the smallest variance among all possible unbiased linear estimators.

## Problem 2

- Consider the equation

$$y = X\beta + u.$$

Suppose that the rows of  $X$  are a random sample from some distribution with mean  $\mathbf{0}$  and covariance matrix  $\Sigma > 0$ , and suppose that  $\text{var}(u|X) = \sigma^2 I_n$ , and  $E(u|X) = X\theta$ , where  $\theta \neq \mathbf{0}$  does not depend on  $n$ .

- Show that  $E(u) = \mathbf{0}$  and  $\text{var}(u) = (\sigma^2 + \theta^\top \Sigma \theta) I_n$ , so that the random disturbances have mean  $\mathbf{0}$ , constant variance, and are uncorrelated with one another, but are correlated with  $X$ .
- Show that the OLS estimator of  $\beta$  satisfies  $\hat{\beta} = \beta + (X^\top X)^{-1} X^\top u$ .
- Show that  $E(\hat{\beta}|X) = \beta + \theta$  and  $\text{var}(\hat{\beta}|X) = \sigma^2 (X^\top X)^{-1}$ .

## Problem 3

- Let  $\hat{\boldsymbol{\beta}}$  be the  $(k \times 1)$  vector of OLS estimates.
1. Show that for any  $(k + 1) \times 1$  vector  $\mathbf{b}$ , we can write the sum of squared residuals as

$$SSR(\mathbf{b}) = \hat{\mathbf{u}}^\top \hat{\mathbf{u}} + (\hat{\boldsymbol{\beta}} - \mathbf{b})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \mathbf{b}).$$

(Hint: Write  $(\mathbf{y} - \mathbf{X}\mathbf{b})^\top (\mathbf{y} - \mathbf{X}\mathbf{b}) = [\hat{\mathbf{u}} + \mathbf{X}(\hat{\boldsymbol{\beta}} - \mathbf{b})]^\top [\hat{\mathbf{u}} + \mathbf{X}(\hat{\boldsymbol{\beta}} - \mathbf{b})]$  and use the fact that  $\mathbf{X}^\top \hat{\mathbf{u}} = \mathbf{0}$ .)

2. Explain how the expression for  $SSR(\mathbf{b})$  in part 1. above proves that  $\hat{\boldsymbol{\beta}}$  uniquely minimizes  $SSR(\mathbf{b})$  over all possible values of  $\mathbf{b}$ , assuming  $\mathbf{X}$  has rank  $k \times 1$ .

## Problem 4

Let  $\hat{\beta}$  be the OLS estimate from the regression of  $\mathbf{y}$  on  $\mathbf{X}$ . Let  $\mathbf{A}$  be a  $(k+1) \times (k+1)$  nonsingular matrix and define  $z_t \equiv \mathbf{x}_t \mathbf{A}$ ,  $t = 1, \dots, n$ , where  $\mathbf{x}_t$  is the  $t$ th row in  $\mathbf{X}$ . Therefore,  $z_t$  is  $1 \times (k+1)$  and is a nonsingular linear combination of  $\mathbf{x}_t$ . Let  $\mathbf{Z}$  be the  $n \times (k+1)$  matrix with rows  $z_t$ . Let  $\tilde{\beta}$  denote the OLS estimate from a regression of  $\mathbf{y}$  on  $\mathbf{Z}$ .

- Show that  $\tilde{\beta} = \mathbf{A}^{-1} \hat{\beta}$ .
- Let  $\hat{y}_t$  be the fitted values from the original regression and let  $\tilde{y}_t$  be the fitted values from regressing  $\mathbf{y}$  on  $\mathbf{Z}$ . Show that  $\hat{y}_t = \tilde{y}_t$ , for all  $t = 1, \dots, n$ . How do the residuals from the two regressions compare?
- Show that the estimated variance matrix for  $\tilde{\beta}$  is  $s^2 \mathbf{A}^{-1} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^{-1\top}$ , where  $s^2$  is the usual variance estimate from regressing  $\mathbf{y}$  on  $\mathbf{X}$ .

## Problem 5

Consider the usual linear model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$ , where  $\mathbf{y}$  is  $(n \times 1)$ ,  $\mathbf{X}$  is  $(n \times 5)$ ,  $\boldsymbol{\beta}$  is  $(5 \times 1)$ , and  $\mathbf{u}$  is  $(n \times 1)$ ,  $\mathbf{E}(\mathbf{u}) = \mathbf{0}$  and  $\text{var}(\mathbf{u}) = \sigma^2 \mathbf{I}$ . The sample size  $n$  is equal to 500, and  $\mathbf{X}$  is non-stochastic and has full rank. OLS is used to estimate the following coefficients and variance matrix of the coefficients

$$\hat{\boldsymbol{\beta}} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \end{pmatrix} = \begin{pmatrix} 2.5 \\ 1.1 \\ -3 \\ 0.5 \\ 0.4 \end{pmatrix}, \quad \widehat{\text{var}(\hat{\boldsymbol{\beta}})} = \begin{bmatrix} 0.5 & 0.2 & 0.1 & -0.2 & -0.5 \\ 0.2 & 0.2 & 0.3 & -0.1 & -0.7 \\ 0.1 & 0.3 & 0.1 & 0.1 & -0.2 \\ -0.2 & -0.1 & 0.1 & 1 & -0.5 \\ -0.5 & -0.7 & -0.2 & -0.5 & 0.1 \end{bmatrix}.$$



## Problem 5 (continued)

Test the following hypotheses. In each case set up the test for the hypothesis, write down the  $\mathbf{R}$  and  $\mathbf{r}$  matrices and do the F and Wald tests, cf. slide 47.

1.  $H_0 : \beta_1 = 1$
2.  $H_0 : -\beta_2 + \beta_4 + \beta_5 = 0$
3.  $H_0 : \beta_1 = 1, \quad -\beta_2 + \beta_4 + \beta_5 = 0.$

## HINTS to problems

- Problem 1: The Gauss-Markov Theorem, slide page 24. The idea of the proof is to choose another arbitrary linear estimator for  $\beta$ . Call this estimator  $\tilde{\beta}$ . It has the form  $\tilde{\beta} = Ay$ . Define a matrix  $D = A - (X^\top X)^{-1}X^\top$ , and write the variance-covariance matrix of  $\tilde{\beta}$  as something that depends on  $D$  and  $X$ , and recognize that  $\text{var}(\tilde{\beta}) - \text{var}(\hat{\beta})$  must be positive semi definite.
- Problem 2, second bullet on slide 76: Here you should use the **vector versions** of two general and important rules: (1) The law of iterated expectations:  $E[Y] = E[E[Y|Z]]$  for any pair of random variables. (2) The law of unconditional variance:  $\text{var}[Y] = E[\text{var}[Y|Z]] + \text{var}[E[Y|Z]]$  for any pair of random variables.
- Problem 3: Should be straightforward.
- Problem 4: Should be straightforward when you realize that  $Z = XA$ .
- Problem 5: Rewrite the F and Wald statistics on slide 47 using the fact that

$$s^2 = \frac{\hat{u}^\top \hat{u}}{n - k} \quad \text{and} \quad \widehat{\text{var}(\hat{\beta})} = s^2 (X^\top X)^{-1}.$$