

# Econometric Methods (Econometrics I)

## Lecture 1: Introduction and Statistical Prerequisites

Prof. Dr. Kai Carstensen

Kiel University

Winter Term 2023/2024

1. Introduction: on estimating causal effects
2. Some math: vector and matrix derivatives
3. Review of basic OLS: estimation
4. Review of basic OLS: inference
5. Conditional expectations
6. Linear projections

Reference: Wooldridge, Chapters 1-2.

# 1. Introduction

Economics suggests important relationships, often with policy implications, but virtually never suggests quantitative magnitudes of **causal effects**.

- How does another year of education change earnings?
- Does reducing class size cause an improvement in student performance?
- What is the effect of having children on labor supply?
- Does lowering the business property tax rate (in Germany, say, the *Gewerbesteuerhebesatz*) cause an increase in city economic activity?
- What is the effect on output growth of a 1 percentage point increase in interest rates by the ECB?
- What is the effect on GDP of a reduction in VAT?

Note: Econometrics is also good for forecasting when using time-series data. But in this course we concentrate on cross section data. The time series dimension is discussed in Econometrics II. And panel data in the respective specialization.

# How use data to measure causal effects?

- Ideally, we would like to have data from a **randomized controlled experiment**.
  - **Controlled** means that there are a control group that receives no treatment and a treatment group that receives a treatment.
  - **Randomized** means that the treatment is assigned randomly.
- But almost always we only have observational (nonexperimental) data:
  - returns to education
  - wage statistics
  - survey data
- Most of the course deals with difficulties arising from using observational data to estimate causal effects:
  - confounding effects (omitted factors)
  - simultaneous causality
  - correlation does not imply causation
  - nonlinearities
  - etc

# How to find a causal relationship?

- Use or develop a plausible theory that explains why a causing variable,  $w$ , should have an impact on another variable,  $y$ .
- Mimic an experiment: Make sure that “everything else” is either held constant or randomly distributed between those treated and those not treated.
- How can we hold other factors constant? We can condition on them. This means that effectively we compare treated and non-treated entities which are otherwise identical.
- Take into account that economic relationships are not deterministic. Hence, focus on the expected (average) relationship between  $y$  and  $x$  holding a vector of control variables,  $\mathbf{c}$ , constant:  $E(y|w, \mathbf{c})$ .
- What we are typically interested in, is (for continuous variables) the partial effect

$$\frac{\partial E(y|w, \mathbf{c})}{\partial w}$$

- This sounds simple, especially if  $E(y|w, \mathbf{c})$  is a linear function of  $w$  and  $\mathbf{c}$ , but there are many reasons why it may not. This is what this lecture is about.

In this course, we assume **random sampling**. This means that we

1. specify a **population model** and
2. draw an **independently and identically distributed** (iid) random sample (individuals, firms, etc.)  $i = 1, 2, \dots, N$  from this population.

## Population model:

- Most often a model of the conditional expectation. For example,

$$E(y|w, \mathbf{c}) = \beta_0 + \beta_1 w + \boldsymbol{\gamma}' \mathbf{c}.$$

- Think about the population. This can be very important. For example, when you are interested in how another year of education changes earnings, you must be specific: the active working population? including the currently jobless? only people without university degree? only people without vocational training? etc
- Note: You have to adapt sampling and measuring (what are the earnings of jobless people?) to your definition of the population.



## Random sampling:

- Often a sensible assumption for cross-sectional data.
- Implies that **all** variables (dependent variables and explanatory variables) are random.
- Hence, the assumption of fixed regressors sometimes still used in the baseline regression model does not apply.
- Moreover, the normality assumptions often used in basic econometrics need not hold: why should the variables be normally distributed?
- Consequence: finite-sample distributions of classical  $t$ -tests and  $F$ -tests are unknown.
- Thus, we resort to large-sample approximations using asymptotic analysis. This seems justified as most cross-sectional data sets are of moderate or even large size.

Note: the random sampling assumption might be more difficult to defend for panel and time series data. But this is delegated to the respective courses.

Wage offer function for women: Suppose that the natural log of the wage offer,  $wage^o$ , is determined as

$$\log(wage^o) = \beta_0 + \beta_1 educ + \beta_2 exper + \beta_3 married + u \quad (1)$$

where

- $educ$  is years of schooling,  $exper$  is years of labor market experience, and  $married$  is a binary variable indicating marital status,
- $u$ , called the error term or disturbance, contains unobserved factors that affect the wage offer,
- interest lies in the unknown parameters, the  $\beta_j$ ,
- random sampling depends on the population of interest (all working women? all women over age 18?)

For deriving the properties of estimators, it is often useful to write the population model for a generic draw from the population. Instead of

$$\log(\text{wage}^o) = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{married} + u$$

we then write

$$\log(\text{wage}_i^o) = \beta_0 + \beta_1 \text{educ}_i + \beta_2 \text{exper}_i + \beta_3 \text{married}_i + u_i, \quad (2)$$

where  $i$  indexes the individual.

Stating assumptions in terms of  $u_i$  and  $x_i \equiv (\text{educ}_i, \text{exper}_i, \text{married}_i)$  is the same as stating assumptions in terms of  $u$  and  $x$ .

## 2. Some Math

Any book on matrix algebra and calculus should have a section on matrix differentiation.

A brief but nice treatment can be found in: H. Lütkepohl (2007), New Introduction to Multiple Time Series Analysis, Springer-Verlag, Berlin, Appendix A.13.

Also feasible: W.H. Greene (2012), Econometric Analysis, 7th ed., Pearson, Appendix A.

**Rule 1:** Let  $f(\beta)$  be a scalar function of the  $(n \times 1)$  vector  $\beta = (\beta_1, \dots, \beta_n)'$ . The column vector of first partial derivatives is

$$\underbrace{\frac{\partial f}{\partial \beta}}_{n \times 1} := \begin{bmatrix} \partial f / \partial \beta_1 \\ \vdots \\ \partial f / \partial \beta_n \end{bmatrix}.$$

The row vector of first partial derivatives (often called gradient, symbol:  $\nabla_\beta$ ) is

$$\underbrace{\frac{\partial f}{\partial \beta'}}_{1 \times n} = \nabla_\beta f(\beta) := \left[ \frac{\partial f}{\partial \beta_1}, \dots, \frac{\partial f}{\partial \beta_n} \right].$$

The  $(n \times n)$  Hessian matrix of second partial derivatives is

$$\underbrace{\frac{\partial^2 f}{\partial \beta \partial \beta'}}_{n \times n} = \frac{\partial \frac{\partial f}{\partial \beta}}{\partial \beta'} = \nabla_\beta [\nabla_\beta f(\beta)]' := \begin{bmatrix} \frac{\partial^2 f}{\partial \beta_1 \partial \beta_1} & \cdots & \frac{\partial^2 f}{\partial \beta_1 \partial \beta_n} \\ \vdots & & \vdots \\ \frac{\partial^2 f}{\partial \beta_n \partial \beta_1} & \cdots & \frac{\partial^2 f}{\partial \beta_n \partial \beta_n} \end{bmatrix}.$$

Example:

$$f(\boldsymbol{\beta}) = \beta_1^2 + \cdots + \beta_n^2 = \sum_{i=1}^n \beta_i^2 = \boldsymbol{\beta}'\boldsymbol{\beta}$$

$$\frac{\partial f}{\partial \boldsymbol{\beta}} = \frac{\partial \boldsymbol{\beta}'\boldsymbol{\beta}}{\partial \boldsymbol{\beta}} = \begin{bmatrix} \frac{\partial f}{\partial \beta_1} \\ \vdots \\ \frac{\partial f}{\partial \beta_n} \end{bmatrix} = \begin{bmatrix} 2\beta_1 \\ \vdots \\ 2\beta_n \end{bmatrix} = 2 \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_n \end{bmatrix} = 2\boldsymbol{\beta}$$

$$\frac{\partial^2 f}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} = \frac{\partial \boldsymbol{\beta}'\boldsymbol{\beta}}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} = \begin{bmatrix} 2 & 0 & \cdots & 0 \\ 0 & 2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 2 \end{bmatrix} = 2I_n$$

**Rule 2:** Let

$$h(\beta) = \begin{bmatrix} h_1(\beta) \\ \vdots \\ h_m(\beta) \end{bmatrix}$$

be a  $(m \times 1)$  vector function that depends on the  $(n \times 1)$  vector  $\beta = (\beta_1, \dots, \beta_n)'$ .

Then the  $(m \times n)$  matrix of first partial derivatives is

$$\frac{\partial h}{\partial \beta'} := \begin{bmatrix} \frac{\partial h_1}{\partial \beta_1} & \dots & \frac{\partial h_1}{\partial \beta_n} \\ \vdots & & \vdots \\ \frac{\partial h_m}{\partial \beta_1} & \dots & \frac{\partial h_m}{\partial \beta_n} \end{bmatrix}.$$



Example:

$$h(\beta) = \mathbf{X}\beta = \begin{bmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & & \vdots \\ x_{m1} & \cdots & x_{mn} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_n \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n x_{1i}\beta_i \\ \vdots \\ \sum_{i=1}^n x_{mi}\beta_i \end{bmatrix}$$

$$\frac{\partial h}{\partial \beta'} = \frac{\partial \mathbf{X}\beta}{\partial \beta'} = \begin{bmatrix} \frac{\partial \sum_{i=1}^n x_{1i}\beta_i}{\partial \beta_1} & \cdots & \frac{\partial \sum_{i=1}^n x_{1i}\beta_i}{\partial \beta_n} \\ \vdots & & \vdots \\ \frac{\partial \sum_{i=1}^n x_{mi}\beta_i}{\partial \beta_1} & \cdots & \frac{\partial \sum_{i=1}^n x_{mi}\beta_i}{\partial \beta_n} \end{bmatrix} = \begin{bmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & & \vdots \\ x_{m1} & \cdots & x_{mn} \end{bmatrix} = \mathbf{X}$$

**Rule 3** (Chain rule for vector differentiation): Let  $\alpha$  and  $\beta$  be  $(m \times 1)$  and  $(n \times 1)$  vectors, respectively, and suppose  $h(\alpha)$  is  $(p \times 1)$  and  $g(\beta)$  is  $(m \times 1)$ . Further assume  $\alpha = g(\beta)$  and define  $f(\beta) = h(g(\beta))$ . Then

$$\frac{\partial f(\beta)}{\partial \beta'} = \frac{\partial h(g(\beta))}{\partial \beta'} = \frac{\partial h(\alpha)}{\partial \alpha'} \frac{\partial g(\beta)}{\partial \beta'}.$$

Example:

$$h(\alpha) = \alpha' \alpha, \quad \alpha = g(\beta) = y - X\beta, \quad f(\beta) = (y - X\beta)'(y - X\beta)$$

$$\frac{\partial f}{\partial \beta'} = \frac{\partial h(\alpha)}{\partial \alpha'} \frac{\partial g(\beta)}{\partial \beta'} = \frac{\partial \alpha' \alpha}{\partial \alpha'} \frac{\partial (y - X\beta)}{\partial \beta'} = -2\alpha' X = -2(y - X\beta)' X$$

### 3. Review of Basic OLS: Estimation

**This is a review of Bachelor level econometrics. Those not familiar with it should join the prep course offered in the first half of the semester.**

# The multiple regression model

## The population model

The population model we study is linear in its parameters,

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_K x_{Ki} + u_i$$

where

- $y, x_1, \dots, x_K$  are observable random scalars,
- $u$  is the unobservable random disturbance or error and
- $\beta_0, \beta_1, \beta_2, \dots, \beta_K$  are the parameters to be estimated.

In many cases—especially when we write the model in vector form—we want to absorb the intercept into the slope parameters. Then we write

$$y_i = \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_K x_{Ki} + u_i$$

and assume  $x_{1i} = 1$ .

# The multiple regression model

## Vector form

The regression model in vector form:

$$y_i = \mathbf{x}_i \boldsymbol{\beta} + u_i$$

where

$$\mathbf{x}_i \equiv (x_{1i}, \dots, x_{Ki}) \quad (1 \times K)$$

$$\boldsymbol{\beta} \equiv (\beta_1, \dots, \beta_K)' \quad (K \times 1).$$

Typically, in this formulation one regressor is a constant, e.g.,  $x_{1i} = 1$ .

# The multiple regression model

## Matrix form

The regression model in matrix form for a sample of  $N$  observations:

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{U} \quad (3)$$

with the observation matrices

$$\mathbf{Y} \equiv (y_1, \dots, y_N)' \quad (N \times 1)$$

$$\mathbf{X} \equiv (\mathbf{x}'_1, \dots, \mathbf{x}'_N)' \quad (N \times K)$$

and the disturbances

$$\mathbf{U} \equiv (u_1, \dots, u_N)' \quad (N \times 1).$$

# OLS estimation

## Derivation of the estimator

There is a sample of  $N$  observations  $\mathbf{Y}$  and  $\mathbf{X}$ , where  $\mathbf{X}$  has full column rank such that  $\mathbf{X}'\mathbf{X}$  is invertible. We want to find a vector  $\hat{\beta}$  that minimizes the squared distance between the elements of  $\mathbf{Y}$  and  $\mathbf{X}\hat{\beta}$ . Hence, we have to minimize the objective function

$$S(\hat{\beta}) = \sum_{i=1}^N (y_i - \mathbf{x}_i' \hat{\beta})^2 = \sum_{i=1}^N \hat{u}_i^2 = \hat{\mathbf{U}}' \hat{\mathbf{U}} = (\mathbf{Y} - \mathbf{X}\hat{\beta})' (\mathbf{Y} - \mathbf{X}\hat{\beta}).$$

First-order derivative:

$$\frac{\partial S(\hat{\beta})}{\partial \hat{\beta}'} = \frac{\partial S}{\partial \hat{\mathbf{U}}'} \frac{\partial \hat{\mathbf{U}}}{\partial \hat{\beta}'} = \frac{\partial \hat{\mathbf{U}}' \hat{\mathbf{U}}}{\partial \hat{\mathbf{U}}'} \frac{\partial (\mathbf{Y} - \mathbf{X}\hat{\beta})}{\partial \hat{\beta}'} = -2\hat{\mathbf{U}}' \mathbf{X} = -2(\mathbf{Y} - \mathbf{X}\hat{\beta})' \mathbf{X}$$

First order condition:

$$-2(\mathbf{Y} - \mathbf{X}\hat{\beta}_{LS})' \mathbf{X} \stackrel{!}{=} 0 \quad \Rightarrow \quad \mathbf{Y}' \mathbf{X} - \hat{\beta}_{LS}' \mathbf{X}' \mathbf{X} = 0 \quad \Rightarrow \quad \hat{\beta}_{LS}' = \mathbf{Y}' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1}$$

OLS estimator:

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

Does the OLS estimator find the minimum of the objective function  $S(\hat{\beta})$ ?

Check the Hessian:

$$\begin{aligned}\frac{\partial^2 S(\hat{\beta})}{\partial \hat{\beta} \partial \hat{\beta}'} &= \frac{\partial}{\partial \hat{\beta}'} \left( \frac{\partial S(\hat{\beta})}{\partial \hat{\beta}} \right) = \frac{\partial}{\partial \hat{\beta}'} \left( -2\mathbf{X}'(\mathbf{Y} - \mathbf{X}\hat{\beta}) \right) \\ &= \frac{\partial}{\partial \hat{\beta}'} \left( -2\mathbf{X}'\mathbf{Y} + 2\mathbf{X}'\mathbf{X}\hat{\beta} \right) = \frac{\partial(2\mathbf{X}'\mathbf{X}\hat{\beta})}{\partial \hat{\beta}'} = 2\mathbf{X}'\mathbf{X}.\end{aligned}$$

The Hessian is (for any value of  $\hat{\beta}_{LS}$ ) positive definite because the quadratic form  $\mathbf{X}'\mathbf{X}$  is positive definite. Hence, the OLS estimator found a (local) minimum of the objective function.



Note that the “matrix expression”

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

is equivalent to the “summation expression”

$$\hat{\beta}_{LS} = \left( N^{-1} \sum_{i=1}^N \mathbf{x}_i' \mathbf{x}_i \right)^{-1} \left( N^{-1} \sum_{i=1}^N \mathbf{x}_i' y_i \right).$$

When the regressor list includes an intercept (what it almost always does),

$$y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_K x_{Ki} + u_i = \beta_0 + \mathbf{x}_i \tilde{\boldsymbol{\beta}} + u_i,$$

we can derive an insightful expression for the OLS estimator if we separate the intercept  $\beta_0$  from the other coefficients  $\tilde{\boldsymbol{\beta}} = (\beta_1, \dots, \beta_K)'$  and define  $\mathbf{x}_i = (x_{1i}, \dots, x_{Ki})$ .

To this end, compute the average over all observations,

$$\bar{y} = \beta_0 + \bar{\mathbf{x}} \tilde{\boldsymbol{\beta}} + \bar{u},$$

and subtract it from the previous equation which yields the “de-meanned” equation

$$\underbrace{y_i - \bar{y}}_{\tilde{y}_i} = \underbrace{(\mathbf{x}_i - \bar{\mathbf{x}})}_{\tilde{\mathbf{x}}_i} \tilde{\boldsymbol{\beta}} + \underbrace{u_i - \bar{u}}_{\tilde{u}_i}.$$

Applying OLS to the de-meaned equation yields

$$\begin{aligned}\hat{\beta}_{LS} &= \left( N^{-1} \sum_{i=1}^N \tilde{\mathbf{x}}_i' \tilde{\mathbf{x}}_i \right)^{-1} \left( N^{-1} \sum_{i=1}^N \tilde{\mathbf{x}}_i' \tilde{y}_i \right) \\ &= \left( N^{-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})' (\mathbf{x}_i - \bar{\mathbf{x}}) \right)^{-1} \left( N^{-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})' (y_i - \bar{y}) \right) \\ &= \left( s_{\mathbf{x}}^2 \right)^{-1} s_{\mathbf{x}',y}\end{aligned}$$

where

- $s_{\mathbf{x}}^2$  is the sample variance matrix of  $\mathbf{x}$ , and
- $s_{\mathbf{x}',y}$  is the sample covariance matrix between  $\mathbf{x}'$  and  $y$ .

Define the OLS residuals as

$$\hat{\mathbf{U}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{LS} = \mathbf{Y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = (\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\mathbf{Y} = \mathbf{M}_X\mathbf{Y},$$

where

$$\mathbf{M}_X = \mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

is a symmetric, idempotent matrix of rank  $N - K$ . Hence,

$$\mathbf{M}_X = \mathbf{M}_X' \quad \text{and} \quad \mathbf{M}_X\mathbf{M}_X = \mathbf{M}_X.$$

For later use, the following reformulation is helpful

$$\hat{\mathbf{U}} = \mathbf{M}_X\mathbf{Y} = \mathbf{M}_X(\mathbf{X}\boldsymbol{\beta} + \mathbf{U}) = \mathbf{M}_X\mathbf{X}\boldsymbol{\beta} + \mathbf{M}_X\mathbf{U} = \mathbf{M}_X\mathbf{U}$$

because  $\mathbf{M}_X\mathbf{X} = (\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\mathbf{X} = \mathbf{X} - \mathbf{X} = \mathbf{0}$ .

A widely used measure of fit, the  $R^2$  is based on a variance decomposition which can be represented in terms of the de-meaned equation.

Defining the de-meaned observation matrices  $\tilde{\mathbf{Y}}$  and  $\tilde{\mathbf{X}}$ , the first order condition implies

$$\tilde{\mathbf{X}}'(\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\hat{\beta}_{LS}) = 0 \quad \Rightarrow \quad \tilde{\mathbf{X}}'\hat{\mathbf{U}} = 0.$$

Define  $\hat{\tilde{\mathbf{Y}}} = \tilde{\mathbf{X}}\hat{\beta}_{LS}$  and thus  $\tilde{\mathbf{Y}} = \hat{\tilde{\mathbf{Y}}} + \hat{\mathbf{U}}$ . Then the following decomposition holds:

$$\tilde{\mathbf{Y}}'\tilde{\mathbf{Y}} = (\hat{\tilde{\mathbf{Y}}} + \hat{\mathbf{U}})'(\hat{\tilde{\mathbf{Y}}} + \hat{\mathbf{U}}) = \hat{\tilde{\mathbf{Y}}}'\hat{\tilde{\mathbf{Y}}} + \hat{\mathbf{U}}'\hat{\mathbf{U}}$$

because the cross product is zero,  $\hat{\tilde{\mathbf{Y}}}'\hat{\mathbf{U}} = \hat{\beta}_{LS}'\tilde{\mathbf{X}}'\hat{\mathbf{U}} = 0$ .

Therefore, a measure of fit is

$$R^2 = \frac{\text{explained sum of squares}}{\text{total sum of squares}} = \frac{\hat{\tilde{\mathbf{Y}}}'\hat{\tilde{\mathbf{Y}}}}{\tilde{\mathbf{Y}}'\tilde{\mathbf{Y}}} = \frac{\tilde{\mathbf{Y}}'\tilde{\mathbf{Y}} - \hat{\mathbf{U}}'\hat{\mathbf{U}}}{\tilde{\mathbf{Y}}'\tilde{\mathbf{Y}}} = 1 - \frac{\hat{\mathbf{U}}'\hat{\mathbf{U}}}{\tilde{\mathbf{Y}}'\tilde{\mathbf{Y}}}.$$

In summation form, we obtain

$$R^2 = 1 - \frac{\hat{\mathbf{U}}' \hat{\mathbf{U}}}{\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}}} = 1 - \frac{\sum_{i=1}^N \hat{u}_i^2}{\sum_{i=1}^N \tilde{y}_i^2} = 1 - \frac{\sum_{i=1}^N \hat{u}_i^2}{\sum_{i=1}^N (y_i - \bar{y})^2}.$$

Using

$$\begin{aligned} \sum_{i=1}^N (y_i - \bar{y})^2 &= \sum_{i=1}^N (y_i^2 - 2\bar{y}y_i + \bar{y}^2) = \sum_{i=1}^N y_i^2 - 2\bar{y} \sum_{i=1}^N y_i + N\bar{y}^2 \\ &= \sum_{i=1}^N y_i^2 - 2N\bar{y}^2 + N\bar{y}^2 = \sum_{i=1}^N y_i^2 - N\bar{y}^2 = \mathbf{Y}'\mathbf{Y} - N\bar{y}^2 \end{aligned}$$

the  $R^2$  can also be written as

$$R^2 = 1 - \frac{\hat{\mathbf{U}}' \hat{\mathbf{U}}}{\mathbf{Y}'\mathbf{Y} - N\bar{y}^2}.$$

## 4. Review of Basic OLS: Inference

**This is a review of Bachelor level econometrics. Those not familiar with it should join the prep course offered in the first half of the semester.**

In this section (and only in this section), we make the following (strict) assumptions:

- The  $K$  regressors  $\mathbf{x}_i$  are not random but **fixed numbers**.
- The observation matrix  $\mathbf{X}$  has full column rank  $K$  so that  $\mathbf{X}'\mathbf{X}$  is invertible.
- The disturbances  $u_i$ ,  $i = 1, \dots, N$ , are independently and identically normally distributed with  $E[u_i] = 0$  and  $\text{Var}[u_i] = \sigma^2$ :

$$u_i \sim \text{Normal}(0, \sigma^2).$$



# Distribution of the OLS estimator

## Mean and variance of the disturbance vector $\mathbf{U}$

As all  $u_i$ ,  $i = 1, \dots, N$ , have mean zero, so has the vector  $\mathbf{U}$ :  $E[\mathbf{U}] = \mathbf{0}$ .

The i.i.d. assumption on  $u_i$  implies that any two  $u_i$  and  $u_j$ ,  $i \neq j$ , are uncorrelated:

$$\text{Cov}[u_i, u_j] = E[u_i u_j] = 0.$$

Hence, the variance of the vector of disturbances is

$$\begin{aligned}\text{Var}[\mathbf{U}] &= E[(\mathbf{U} - E[\mathbf{U}])(\mathbf{U} - E[\mathbf{U}])'] = E[\mathbf{U}\mathbf{U}'] \\ &= E\left[\begin{pmatrix} u_1 \\ \vdots \\ u_N \end{pmatrix} (u_1 \quad \dots \quad u_N)\right] = E\begin{pmatrix} u_1^2 & u_1 u_2 & \dots & u_1 u_N \\ u_2 u_1 & u_2^2 & \dots & u_2 u_N \\ \vdots & \vdots & \ddots & \vdots \\ u_N u_1 & u_N u_2 & \dots & u_N^2 \end{pmatrix} \\ &= \begin{pmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{pmatrix} = \sigma^2 \mathbf{I}_N\end{aligned}$$

# Distribution of the OLS estimator

## Distribution of the disturbance vector $\mathbf{U}$

Recall that all  $u_i$ ,  $i = 1, \dots, N$ , are individually distributed as

$$u_i \sim \text{Normal}(0, \sigma^2).$$

Given that  $E[\mathbf{U}] = \mathbf{0}$  and  $\text{Var}[\mathbf{U}] = \sigma^2 \mathbf{I}_N$ , the joint distribution of the vector of disturbances is

$$\mathbf{U} \sim \text{Normal}(\mathbf{0}, \sigma^2 \mathbf{I}_N).$$

# Distribution of the OLS estimator

## Mean of the OLS estimator

To find the mean of the OLS estimator

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

substitute the model

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{U}$$

which yields

$$\hat{\beta}_{LS} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\beta + \mathbf{U}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{U} = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{U}.$$

Now take expectations

$$E[\hat{\beta}_{LS}] = E[\beta] + E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{U}] = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' E[\mathbf{U}] = \beta$$

because the regressors are nonrandom and  $E[\mathbf{U}] = 0$ .

Result: the OLS estimator is unbiased.

The variance of the OLS estimator is

$$\begin{aligned}\text{Var}[\hat{\beta}_{LS}] &= E \left[ \hat{\beta}_{LS} - E(\hat{\beta}_{LS}) \right] \left[ \hat{\beta}_{LS} - E(\hat{\beta}_{LS}) \right]' \\&= E \left[ \hat{\beta}_{LS} - \beta \right] \left[ \hat{\beta}_{LS} - \beta \right]' \\&= E \left[ (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{U} \right] \left[ (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{U} \right]' \\&= E \left[ (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{U}\mathbf{U}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \right] \\&= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' E[\mathbf{U}\mathbf{U}']\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\&= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\sigma^2\mathbf{I}_N\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\&= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\&= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}.\end{aligned}$$

# Distribution of the OLS estimator

## Distribution of linear functions of normal random variables

Rule: Let  $\mathbf{W}$  be a  $N \times 1$  random vector with

$$\mathbf{W} \sim \text{Normal}(\mathbf{0}, \mathbf{V}),$$

and let  $\mathbf{A}$  be a  $K \times N$  matrix of full row rank  $K$  and  $\mathbf{B}$  a  $K$  vector. Then the  $K \times 1$  random vector

$$\mathbf{Z} = \mathbf{AW} + \mathbf{B}$$

is normally distributed with

$$E[\mathbf{Z}] = E[\mathbf{AW} + \mathbf{B}] = \mathbf{B}$$

$$\text{Var}[\mathbf{Z}] = E[(\mathbf{Z} - \mathbf{B})(\mathbf{Z} - \mathbf{B})'] = E[\mathbf{A}\mathbf{W}\mathbf{W}'\mathbf{A}'] = \mathbf{A} E[\mathbf{W}\mathbf{W}'] \mathbf{A}' = \mathbf{AVA}'.$$

Hence,

$$\mathbf{Z} \sim \text{Normal}(\mathbf{B}, \mathbf{AVA}').$$

# Distribution of the OLS estimator

## Case I: Known variance of the disturbances

To find the distribution of the OLS estimator, recall the expression

$$\hat{\beta}_{LS} = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{U}.$$

Since  $\beta$  and  $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$  are constant numbers,  $\hat{\beta}_{LS}$  is a linear function of the random vector  $\mathbf{U}$ .

Using the rule given on the previous slide,  $\hat{\beta}_{LS}$  is thus normally distributed with mean

$$E[\hat{\beta}_{LS}] = \beta$$

and variance

$$\text{Var}[\hat{\beta}_{LS}] = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}.$$

Result:

$$\hat{\beta}_{LS} \sim \text{Normal}(\beta, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}).$$

Unfortunately, this distribution is not very helpful in practice because the variance  $\sigma^2$  of the disturbances is unknown and must be estimated, inducing additional estimation uncertainty that must be accounted for.

# Distribution of the OLS estimator

## Estimation of the variance $\hat{\sigma}^2$

The variance  $\sigma^2$  is estimated as

$$\hat{\sigma}^2 = \frac{1}{N-K} \sum_{i=1}^N \hat{u}_i^2 = \frac{1}{N-K} \hat{\mathbf{U}}' \hat{\mathbf{U}}.$$

Due to the previous result  $\hat{\mathbf{U}} = \mathbf{M}_x \mathbf{U}$ , this is equivalent to

$$\hat{\sigma}^2 = \frac{1}{N-K} \hat{\mathbf{U}}' \hat{\mathbf{U}} = \frac{1}{N-K} \mathbf{U}' \mathbf{M}_x' \mathbf{M}_x \mathbf{U} = \frac{1}{N-K} \mathbf{U}' \mathbf{M}_x \mathbf{U}.$$

This is a quadratic form in normal random variables.

# Distribution of the OLS estimator

## Distribution of the variance $\hat{\sigma}^2$

Rule: For standard normally distributed random variables,  $\mathbf{W} \sim \text{Normal}(\mathbf{0}, \mathbf{I}_N)$ , the quadratic form

$$\mathbf{W}'\mathbf{C}\mathbf{W},$$

where  $\mathbf{C}$  is a symmetric, idempotent matrix of fixed numbers with rank  $r$ , is  $\chi^2$  distributed with  $r$  degrees of freedom.

To apply this rule, note that

$$\mathbf{U}/\sigma \sim \text{Normal}(\mathbf{0}, \mathbf{I}_N)$$

so that

$$\frac{\mathbf{U}'\mathbf{M}_X\mathbf{U}}{\sigma^2} = \frac{N-K}{\sigma^2} \times \frac{\mathbf{U}'\mathbf{M}_X\mathbf{U}}{N-K} = \frac{N-K}{\sigma^2} \times \hat{\sigma}^2 \sim \chi_{N-K}^2$$

because  $\mathbf{M}_X$  is a symmetric, idempotent matrix of fixed numbers with rank  $N-K$ .



# Distribution of the OLS estimator

## Student's $t$ distribution

Rule: If

- (i) the random variable  $Z$  has a standard normal distribution,
- (ii) the random variable  $W$  has a  $\chi_r^2$  distribution, and
- (iii)  $Z$  and  $W$  are independently distributed

then the random variable  $Z/\sqrt{W/r}$  is  $t$  distributed with  $r$  degrees of freedom,

$$\frac{Z}{\sqrt{W/r}} \sim t_r.$$

# Distribution of the OLS estimator

Case II: Distribution of the  $t$ -statistic when the variance is unknown

The  $t$ -statistic for a single parameter  $\beta_i$  is

$$t = \frac{\hat{\beta}_{LS,i} - \beta_i}{\widehat{\text{se}}(\hat{\beta}_{LS,i})},$$

where  $\widehat{\text{se}}(\hat{\beta}_{LS,i})$  is the  $i$ -th diagonal element of the root of the estimated variance of  $\hat{\beta}_{LS}$

$$\widehat{\text{Var}}(\hat{\beta}_{LS}) = \hat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1}.$$

Note that this can be rewritten as

$$\widehat{\text{Var}}(\hat{\beta}_{LS}) = \frac{\hat{\sigma}^2}{\sigma^2} \times \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} = \frac{\hat{\sigma}^2}{\sigma^2} \text{Var}(\hat{\beta}_{LS}).$$

Hence, the estimated variance of  $\hat{\beta}_{LS,i}$  is

$$\widehat{\text{Var}}(\hat{\beta}_{LS,i}) = \frac{\hat{\sigma}^2}{\sigma^2} \text{Var}(\hat{\beta}_{LS,i})$$

and the estimated standard error is thus

$$\widehat{\text{se}}(\hat{\beta}_{LS,i}) = \frac{\hat{\sigma}}{\sigma} \sqrt{\text{Var}(\hat{\beta}_{LS,i})}.$$

# Distribution of the OLS estimator

## Case II: Distribution of the $t$ -statistic when the variance is unknown

Therefore, the  $t$ -statistic for a single parameter  $\beta_i$  can be written as

$$\begin{aligned} t_i &= \frac{\hat{\beta}_{LS,i} - \beta_i}{\widehat{\text{se}}(\hat{\beta}_{LS,i})} = \frac{\hat{\beta}_{LS,i} - \beta_i}{\frac{\hat{\sigma}}{\sigma} \sqrt{\text{Var}(\hat{\beta}_{LS,i})}} = \frac{(\hat{\beta}_{LS,i} - \beta_i) / \sqrt{\text{Var}(\hat{\beta}_{LS,i})}}{\hat{\sigma} / \sigma} \\ &= \frac{(\hat{\beta}_{LS,i} - \beta_i) / \sqrt{\text{Var}(\hat{\beta}_{LS,i})}}{\sqrt{\frac{\mathbf{U}'\mathbf{M}_X\mathbf{U}}{\sigma^2} / (N - K)}} = \frac{Z}{\sqrt{W / (N - K)}}. \end{aligned}$$

Note that

- (i) the nominator  $Z = (\hat{\beta}_{LS,i} - \beta_i) / \sqrt{\text{Var}(\hat{\beta}_{LS,i})}$  has a standard normal distribution,
- (ii)  $W = \frac{\mathbf{U}'\mathbf{M}_X\mathbf{U}}{\sigma^2}$  has a  $\chi^2_{N-K}$  distribution, and
- (iii)  $Z$  and  $W$  are independently distributed (we do not show this here).

Hence,  $t_i$  is  $t$ -distributed with  $N - K$  degrees of freedom.

Since the  $t$  statistic does not depend on any unknowns, we can use it to derive confidence intervals and coefficient tests.

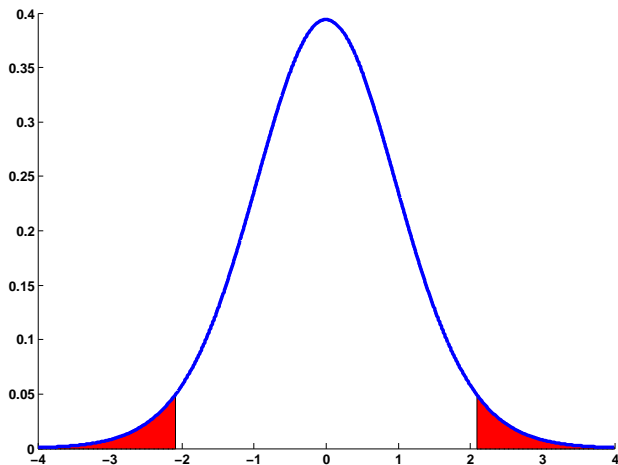
A  $(1 - \alpha)100\%$  two-sided confidence interval for  $\beta_i$  can be derived from

$$\Pr(t_{N-K, \alpha/2} \leq t \leq t_{N-K, 1-\alpha/2}) = 1 - \alpha,$$

where  $t_{N-K, \alpha/2}$  is the  $\alpha/2$ -quantile and  $t_{N-K, 1-\alpha/2}$  is the  $1 - \alpha/2$ -quantile of the  $t$  distribution with  $N - K$  degrees of freedom.

Note: the  $t$  distribution is symmetric, thus  $t_{N-K, \alpha/2} = -t_{N-K, 1-\alpha/2}$ :

$$\Pr(-t_{N-K, 1-\alpha/2} \leq t \leq t_{N-K, 1-\alpha/2}) = 1 - \alpha.$$



**Figure:  $t$  distribution with 20 degrees of freedom**

Use the symmetry and rearrange:

$$\Pr \left( -t_{N-K,1-\alpha/2} \leq (\hat{\beta}_i - \beta_i) / \widehat{\text{se}}(\hat{\beta}_{LS,i}) \leq t_{N-K,1-\alpha/2} \right) = 1 - \alpha$$
$$\Pr \left( \hat{\beta}_i - \widehat{\text{se}}(\hat{\beta}_{LS,i}) t_{N-K,1-\alpha/2} \leq \beta_i \leq \hat{\beta}_i + \widehat{\text{se}}(\hat{\beta}_{LS,i}) t_{N-K,1-\alpha/2} \right) = 1 - \alpha.$$

So a  $(1 - \alpha)100\%$  two-sided confidence interval for  $\beta_i$  is given by the set

$$\left[ \hat{\beta}_i - \widehat{\text{se}}(\hat{\beta}_{LS,i}) t_{N-K,1-\alpha/2}, \hat{\beta}_i + \widehat{\text{se}}(\hat{\beta}_{LS,i}) t_{N-K,1-\alpha/2} \right].$$

It contains the true value of  $\beta_i$  with a probability of  $(1 - \alpha)100\%$ .

Important: By “probability of  $(1 - \alpha)100\%$ ” we mean that if we could draw a large number of independent samples of  $\mathbf{Y}$  from the population (remember: the  $\mathbf{X}$  are fixed in this setting), in  $(1 - \alpha)100\%$  of the cases the confidence interval would contain the true value of  $\beta_i$ .

To test a two-sided hypothesis concerning a single parameter at the significance level of  $\alpha 100\%$ ,

$$H_0 : \beta_i = \beta_{i,0} \quad \text{vs.} \quad H_1 : \beta_i \neq \beta_{i,0}$$

we can use the test statistic

$$t = \frac{\hat{\beta}_{LS,i} - \beta_{i,0}}{\text{se}(\hat{\beta}_{LS,i})}.$$

Under the null hypothesis,  $t$  is  $t_{N-K}$  distributed. For a test decision, use again the symmetry of the  $t$  distribution and compare the absolute value of the  $t$  statistic with the  $(1 - \alpha/2)100\%$  quantile of the  $t_{N-K}$  distribution:

$$|t| > t_{N-K, 1-\alpha/2} \quad \Rightarrow \quad \text{reject } H_0.$$

# The $p$ -value

For a two-sided hypothesis the  $p$ -value is the probability of observing a value of  $\hat{\beta}_{LS,i}$  at least as different from  $\beta_i$  as the estimate actually computed from the data at hand ( $\hat{\beta}_{LS,i}^{act}$ ) if the null hypothesis ( $\beta_i = \beta_{i,0}$ ) is correct:

$$\begin{aligned} p\text{-value} &= \Pr_{H_0} \left[ |\hat{\beta}_i - \beta_{i,0}| > |\hat{\beta}_i^{act} - \beta_{i,0}| \right] \\ &= \Pr_{H_0} \left[ \left| \frac{\hat{\beta}_i - \beta_{i,0}}{\text{se}(\hat{\beta}_{LS,i})} \right| > \left| \frac{\hat{\beta}_i^{act} - \beta_{i,0}}{\text{se}(\hat{\beta}_{LS,i})} \right| \right] = \Pr_{H_0} [|t| > |t^{act}|] . \end{aligned}$$

Because under the null the  $t$ -statistic is  $t_{N-K}$  distributed, we have

$$p\text{-value} = \Pr(t < -|t^{act}|) + \Pr(t > |t^{act}|) = 2 \Pr(t > |t^{act}|).$$

The  $p$ -value is routinely computed by econometrics packages such as Stata.

For any test, the null hypothesis is rejected if the  $p$ -value is smaller than the significance level  $\alpha$ .



Consider a joint hypothesis that is linear in the coefficients and imposes  $q$  restrictions, where  $q \leq K$ . This hypothesis can be written in matrix notation as

$$\mathbf{R}\boldsymbol{\beta} = \mathbf{r},$$

where  $\mathbf{R}$  is a  $q \times K$  nonrandom matrix with full row rank and  $\mathbf{r}$  is a  $q \times 1$  nonrandom vector.

Example: After estimating the linear model

$$y_i = \beta_1 + \beta_2 x_{2i} + \beta_3 x_{3i} + u_i$$

(note:  $x_{1i} = 1$ ) we want to test the joint hypothesis that (a)  $\beta_1 = \beta_2$  and thus  $\beta_1 - \beta_2 = 0$ , and (b)  $\beta_2 + \beta_3 = 1$ . This gives rise to the restriction matrices

$$\begin{pmatrix} 1 & -1 & 0 \\ 0 & 1 & 1 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}.$$

To test the joint restrictions, we may use an  $F$  test. The test statistic is

$$\begin{aligned} F &= (\mathbf{R}\hat{\beta} - \mathbf{r})' \left[ \widehat{\mathbf{RVar}(\hat{\beta}_{LS})} \mathbf{R}' \right]^{-1} (\mathbf{R}\hat{\beta} - \mathbf{r})/q \\ &= \frac{(\mathbf{R}\hat{\beta} - \mathbf{r})' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1} (\mathbf{R}\hat{\beta} - \mathbf{r})/q}{\hat{\sigma}^2}. \end{aligned}$$

Under the null hypothesis it is  $F$  distributed with  $q$  and  $N - K$  degrees of freedom,

$$F \sim F_{q, N-K},$$

which is used to find the critical value.

Theorem: In the linear regression model with regressor matrix  $\mathbf{X}$ , the OLS estimator  $\hat{\beta}$  is the minimum variance linear unbiased estimator of  $\beta$  if the Gauss-Markov assumptions hold.

Comments:

- Hence, if we restrict our class of estimators to those which are a linear function of  $\mathbf{y}$  and unbiased, there is no better estimator than the OLS estimator in terms of estimation variance.
- The theorem can be extended to random regressors  $\mathbf{X}$ , see Greene (2012, p. 60-61).

## 5. Conditional Expectation

# The role of conditional expectations in econometrics

In most cases, we are interested in the average relationship between  $y$  and  $x$  holding a vector of control variables,  $\mathbf{c}$ , constant.

This may be called a structural conditional expectation:

$$E(y|x, \mathbf{c}).$$

Most OLS models are intended to estimate such a structural conditional expectation. Thus the importance of this concept.

Estimation is straightforward when you have a random sample.

But many complications can arise (which econometricians have to solve), in particular

- measurement error
- reverse causality
- omitted (unobserved) variables

Let us first review the concept of conditional expectation.

Let  $y$  be a random variable (*explained variable*), and let  $\mathbf{x} \equiv (x_1, x_2, \dots, x_K)$  be a  $1 \times \mathbf{K}$  random vector of *explanatory variables*. Then the conditional expectation of  $y$  given  $\mathbf{x}$  is

$$E(y|x_1, \dots, x_K) = E(y|\mathbf{x}) = \int_{-\infty}^{\infty} y f_{y|\mathbf{x}}(y|\mathbf{x}) dy$$

when the random variables are continuous and

$$E(y|x_1, \dots, x_K) = E(y|\mathbf{x}) = \sum_{y_j} y_j f_{y|\mathbf{x}}(y_j|\mathbf{x})$$

when the random variables are discrete.

We are interested in the way  $E(y|\mathbf{x})$  depends on  $\mathbf{x}$ .

To make this explicit, we write it as the function  $\mu(\mathbf{x})$ , i.e.,

$$E(y|x_1, \dots, x_K) = E(y|\mathbf{x}) = \mu(x_1, x_2, \dots, x_K) = \mu(\mathbf{x}).$$

In many cases (such as OLS), we assume  $\mu(\mathbf{x})$  is a function that depends on a finite number of (possibly unknown) parameters.

This yields a **parametric** model for  $E(y|\mathbf{x}) = \mu(\mathbf{x})$ .

For  $K = 2$  explanatory variables, here are some examples.

- linear in regressors and parameters:

$$E(y|x_1, x_2) = \mu(x_1, x_2) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

- nonlinear in regressors but linear in parameters:

$$E(y|x_1, x_2) = \mu(x_1, x_2) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_2^2$$

$$E(y|x_1, x_2) = \mu(x_1, x_2) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$$

- nonlinear in regressors and parameters

$$E(y|x_1, x_2) = \mu(x_1, x_2) = \exp[\beta_0 + \beta_1 \log(x_1) + \beta_2 x_2], \quad y \geq 0, x_1 > 0.$$



Note that

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

is a random variable because  $\mathbf{x}$  is a vector of random variables.

Sometimes it is necessary to concentrate on the case when  $\mathbf{x}$  takes on a particular (nonrandom) value  $\mathbf{x}_0$ . Then

$$E(y|\mathbf{x} = \mathbf{x}_0) = \mu(\mathbf{x}_0)$$

is a numeric value (a realization of a random variable).

In the following, we will often be a bit imprecise, mostly writing

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

even if a realization is meant. This should be clear from the context.

To study the effect of a causing variable on the dependent variable, one may calculate the partial effect of  $x_j$  on  $E(y|\mathbf{x})$  (or, to be somewhat imprecise, the partial effect of  $x_j$  on  $y$ ):

$$\text{Partial effect: } \frac{\partial E(y|\mathbf{x})}{\partial x_j} = \frac{\partial \mu(\mathbf{x})}{\partial x_j}.$$

The partial derivative is defined if  $\mu(\cdot)$  is differentiable and  $x_j$  is continuous.

The change in  $E(y|\mathbf{x})$  when  $x_j$  is increased by a small amount, holding  $x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_K$  constant, thus is approximately:

$$\Delta E(y|\mathbf{x}) \approx \frac{\partial \mu(\mathbf{x})}{\partial x_j} \cdot \Delta x_j, \quad \text{holding } x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_K \text{ fixed.} \quad (4)$$

Note: if  $x_j$  is discrete, partial effects are computed by comparing  $E(y|\mathbf{x})$  at different values of  $x_j$  holding the other variables fixed.

The partial effects in the purely linear model

$$E(y|x_1, x_2) = \mu(x_1, x_2) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

do not depend on the values of the regressors:

$$\frac{\partial E(y|\mathbf{x})}{\partial x_i} = \beta_i \quad \Rightarrow \quad E(dy|\mathbf{x}) = \beta_i dx_i.$$

For small (but not infinitesimally small) changes in  $x_i$  we use the approximation

$$E(\Delta y|\mathbf{x}) \approx \beta_i \Delta x_i.$$

The partial effects in the log-linear model

$$E(y|x_1, x_2) = \mu(x_1, x_2) = \beta_0 + \beta_1 \log(x_1) + \beta_2 \log(x_2)$$

depend on the value of the variable of interest:

$$\frac{\partial E(y|\mathbf{x})}{\partial x_i} = \beta_i / x_i \quad \Rightarrow \quad E(dy|\mathbf{x}) = \beta_i \frac{dx_i}{x_i}.$$

For small (but not infinitesimally small) changes in  $x_i$  we use the approximation

$$E(\Delta y|\mathbf{x}) \approx \beta_i \frac{\Delta x_i}{x_i}.$$

Note that  $100 \frac{\Delta x_i}{x_i}$  is a percent change and  $\beta_i$  is called a semi-elasticity.

Interpretation: an increase in  $x_i$  by 1 percent ( $\frac{\Delta x_i}{x_i} = 0.01$ ) leads to an expected change in  $y$  by  $0.01\beta_i$  units.

The partial effects in the log-log model

$$E(\log(y)|x_1, x_2) = \mu(x_1, x_2) = \beta_0 + \beta_1 \log(x_1) + \beta_2 \log(x_2)$$

depend on the values of the variable of interest and the lhs variable:

$$\frac{\partial E(\log(y)|\mathbf{x})}{\partial x_i} = \beta_i / x_i \quad \Rightarrow \quad E(d \log(y)|\mathbf{x}) = E(dy/y|\mathbf{x}) = \beta_i \frac{dx_i}{x_i}.$$

For small (but not infinitesimally small) changes in  $x_i$  we use the approximation

$$E\left(\frac{\Delta y}{y}|\mathbf{x}\right) \approx \beta_i \frac{\Delta x_i}{x_i}.$$

Note  $\beta_i$  is called an elasticity.

Interpretation: an increase in  $x_i$  by 1 percent ( $\frac{\Delta x_i}{x_i} = 0.01$ ) leads to an expected change in  $y$  by  $\beta_i$  percent.

The partial effects in the interaction model

$$E(y|x_1, x_2) = \mu(x_1, x_2) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$$

depend on the value of another variable, e.g.:

$$\frac{\partial E(y|\mathbf{x})}{\partial x_1} = \beta_1 + \beta_3 x_2.$$

Sometimes  $x_2$  is called a modifier of the partial effect of  $x_1$  on  $y$ .

# The error form of models of conditional expectations

When  $y$  is a random variable we would like to explain in terms of observable variables  $\mathbf{x}$ , it is useful to decompose  $y$  as

$$y = E(y|\mathbf{x}) + u \quad (5)$$

with

$$E(u|\mathbf{x}) = 0. \quad (6)$$

We can *always* write  $y$  as its conditional expectation,  $E(y|\mathbf{x})$ , plus an **error term** or **disturbance term** that has *conditional* mean zero. In fact, thereby we define  $u$ .

Two important properties of  $u$  follow from  $E(u|\mathbf{x}) = 0$ :

1.  $E(u) = 0$  (by the law of iterated expectations, see below)
2.  $u$  is uncorrelated with any function of  $\mathbf{x}$ . In particular, it is uncorrelated with each of  $x_1, \dots, x_K$ . Proof: Wooldridge, p. 31.

# The error form of models of conditional expectations

## Example

The error form allows us to write down population models in the usual way:

Consider again the conditional expectation

$$E(y|x_1, x_2) = \mu(x_1, x_2) = \beta_0 + \beta_1 x_1 + \beta_2 x_2.$$

The corresponding model in error form is

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u, \quad E(u|x_1, x_2) = 0.$$



# The law of iterated expectations

## Simple form

Suppose  $y$  is a random variable and  $\mathbf{x}$  is a random vector. The simple form of the law of iterated expectations (LIE) is

$$E[E(y|\mathbf{x})] = E(y).$$

The proof can be found in any textbook on statistics.

The LIE can be used to directly prove that  $E(u|\mathbf{x}) = 0$  implies  $E(u) = 0$ :

$$E(u) = E[E(u|\mathbf{x})] = E[0] = 0.$$

# The law of iterated expectations

## General formulation

Suppose  $y$  is a random variable,  $\mathbf{w}$  is a random vector and  $\mathbf{x} = \mathbf{f}(\mathbf{w})$  is a function of  $\mathbf{w}$  and thus also a random vector. For example, the vector  $\mathbf{x}$  could be a subset of  $\mathbf{w}$ .

A general formulation of the LIE is

$$E(y|\mathbf{x}) = E[E(y|\mathbf{w})|\mathbf{x}]. \quad (7)$$

In addition, it holds that

$$E(y|\mathbf{x}) = E[E(y|\mathbf{x})|\mathbf{w}]. \quad (8)$$

Both statements are variants of the general rule: “The smaller information set always dominates.”

# The law of iterated expectations

## Examples

Why is the LIE useful?

Let us apply the LIE to the case where  $\mathbf{w} = \{\mathbf{x}, \mathbf{z}\}$ . Then

$$E[E(y|\mathbf{x}, \mathbf{z})|\mathbf{x}] = E(y|\mathbf{x}).$$

This can be very helpful when our theoretical model includes variables  $\mathbf{z}$  that are unobservable.

For example, consider the model

$$E(y|x_1, x_2, z) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 z \quad (9)$$

but where  $z$  is unobserved.

How can we hope to estimate  $\beta_1$  and  $\beta_2$ ?

We need to find the conditional expectation solely with respect to  $x_1$  and  $x_2$ .

# The law of iterated expectations

## Examples

By the LIE, and the linearity of the CE operator,

$$\begin{aligned} E(y|x_1, x_2) &= E(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 z | x_1, x_2) \\ &= \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 E(z|x_1, x_2). \end{aligned} \quad (10)$$

So far, nothing is won. But now let us make a (hopefully economically justified) assumption about  $E(z|x_1, x_2)$ . For example, it might be linear in  $x_1$  and  $x_2$ :

$$E(z|x_1, x_2) = \delta_0 + \delta_1 x_1 + \delta_2 x_2. \quad (11)$$

Then we can plug this into the previous equation and rearrange:

$$\begin{aligned} E(y|x_1, x_2) &= \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 (\delta_0 + \delta_1 x_1 + \delta_2 x_2) \\ &= (\beta_0 + \beta_3 \delta_0) + (\beta_1 + \beta_3 \delta_1) x_1 + (\beta_2 + \beta_3 \delta_2) x_2 \end{aligned}$$

This last expression is linear in  $(x_1, x_2)$  and estimable based on data for  $(x_1, x_2)$ .

Without further restrictions, however, we cannot estimate  $\beta_1$  and  $\beta_2$ . This is a variant of the *identification problem* in econometrics which we cover later.

Rule:

If  $E(u) = 0$  and  $u$  and  $\mathbf{x}$  are independent, then  $E(u|\mathbf{x}) = 0$ .

The converse is not true:

If  $E(u|\mathbf{x}) = 0$ , then  $E(u) = 0$  and  $u$  and  $\mathbf{x}$  are uncorrelated but possibly dependent (for example, their variances might be related).

## 6. Linear Projections

Recall that we use OLS to find the best linear relationship between data  $\mathbf{Y}$  and  $\mathbf{X}$ , where best means that the squared distance

$$S(\hat{\beta}) = (\mathbf{Y} - \mathbf{X}\hat{\beta})'(\mathbf{Y} - \mathbf{X}\hat{\beta})$$

is minimized. Important: this is a linear relationship within a *sample*.

Finding a linear relationship between the random variable  $y$  and the  $(1 \times K)$ -dimensional random vector  $\mathbf{x}$  is also interesting in *population*. This is achieved by minimizing the expected squared distance

$$\min_{b_0, \mathbf{b} \in \mathbb{R}^K} E \left[ (y - b_0 - \mathbf{x}\mathbf{b})^2 \right]. \quad (12)$$

This is called a *linear projection* (or *minimum mean square linear predictor* or *least squares linear predictor*).

Hence, OLS is the sample analogue of linear prediction which is about the population.

# Finding the parameters of the linear projection\*

To minimize

$$E[(y - b_0 - \mathbf{x}\mathbf{b})^2]$$

take the first derivatives,

$$\frac{\partial E[(y - b_0 - \mathbf{x}\mathbf{b})^2]}{\partial b_0} = -2 E[y - b_0 - \mathbf{x}\mathbf{b}]$$

$$\frac{\partial E[(y - b_0 - \mathbf{x}\mathbf{b})^2]}{\partial \mathbf{b}} = -2 E[(y - b_0 - \mathbf{x}\mathbf{b}) \mathbf{x}'],$$

and equate them to zero,

$$E[y - \beta_0 - \mathbf{x}\boldsymbol{\beta}] \stackrel{!}{=} 0$$

$$E[(y - \beta_0 - \mathbf{x}\boldsymbol{\beta}) \mathbf{x}'] \stackrel{!}{=} 0,$$

denoting the solutions to  $b_0$  and  $\mathbf{b}$  by  $\beta_0$  and  $\boldsymbol{\beta}$ , respectively.



For  $\beta_0$  we obtain

$$E[y] - \beta_0 - E[x]\beta = 0 \quad \Rightarrow \quad \beta_0 = E[y] - E[x]\beta$$

and for  $\beta$ , using the solution for  $\beta_0$ ,

$$0 = E[(y - \beta_0 - x\beta) x']$$

$$0 = E[x'y - x'\beta_0 - x'x\beta]$$

$$0 = E[x'y] - E[x']\beta_0 - E[x'x]\beta$$

$$0 = E[x'y] - E[x'](E[y] - E[x]\beta) - E[x'x]\beta$$

$$0 = E[x'y] - E[x']E[y] - E[x'x]\beta + E[x']E[x]\beta$$

$$0 = (E[x'y] - E[x']E[y]) - (E[x'x] - E[x']E[x])\beta$$

$$0 = \text{Cov}(x', y) - \text{Var}(x)\beta$$

$$\beta = [\text{Var}(x)]^{-1} \text{Cov}(x', y)$$

Assume in population there is some possibly nonlinear relationship between  $y$  and  $\mathbf{x}$ .

Generally, we are interested in the conditional expectation

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

but it may be difficult to estimate or is just unknown. Therefore, we may focus on the linear projection

$$L(y|1, \mathbf{x}) = L(y|1, x_1, \dots, x_K) = \beta_0 + \mathbf{x}\beta$$

as a kind of linear approximation in population.

Note that the linear projection is identical to the conditional expectation whenever the conditional expectation is linear in  $\mathbf{x}$ .

Recall: The linear projection minimizes the expected squared deviations from the linear relationship between  $y$  and  $\mathbf{x}$  while OLS minimizes the sample average of the squared deviations from the linear relationship between the observations  $y_i$  and  $\mathbf{x}_i$ ,  $i = 1, \dots, N$ . Hence, the OLS objective function is the sample analogue of the linear projection objective function.

Therefore, it is not surprising that the linear projection parameters

$$\beta_0 = E[y] - E[\mathbf{x}]\beta \quad \text{and} \quad \beta = [\text{Var}(\mathbf{x})]^{-1} \text{Cov}(\mathbf{x}', y)$$

and the OLS parameters

$$\hat{\beta}_{0,OLS} = \bar{y} - \bar{\mathbf{x}}\hat{\beta}_{OLS} \quad \text{and} \quad \hat{\beta}_{OLS} = [\mathbf{s}_{\mathbf{x}}^2]^{-1} \mathbf{s}_{\mathbf{x}',y}.$$

are in the same relationship:

**The OLS parameters are the sample analogues of the linear projection parameters.**

The relationship between OLS and linear projection becomes even more obvious when we define

$$\mathbf{x} = [1, x_1, \dots, x_K].$$

Then the so-defined parameter vector  $\beta$  of the linear projection

$$L(y|\mathbf{x}) = L(y|1, x_1, \dots, x_K) = \mathbf{x}\beta$$

turns out to be

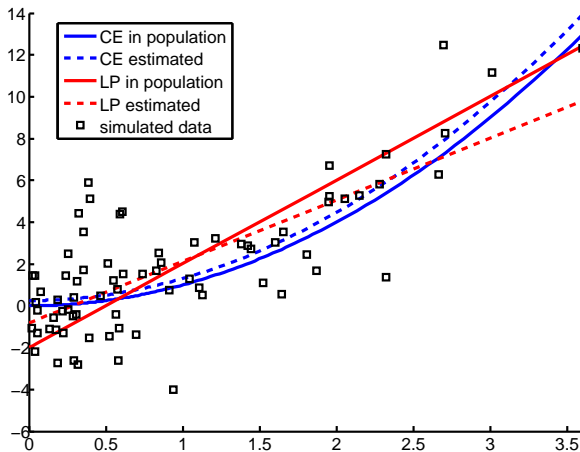
$$\beta = [E(\mathbf{x}'\mathbf{x})]^{-1} E(\mathbf{x}'y).$$

The obvious sample counterpart is the OLS estimator

$$\hat{\beta}_{LS} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} = \left( N^{-1} \sum_{i=1}^N \mathbf{x}_i' \mathbf{x}_i \right)^{-1} \left( N^{-1} \sum_{i=1}^N \mathbf{x}_i' y_i \right).$$

# Example

Model:  $y = x^2 + u$ ,  $x \sim \text{exp}(1)$ ,  $u \sim N(0, 4) \Rightarrow E(y|x) = x^2$ ,  $L(y|1, x) = 4x - 2$



# The error form for linear projections

## Properties

As for the conditional expectation, we can write the linear projection in error form

$$y = L(y|1, \mathbf{x}) + u = \beta_0 + \mathbf{x}\boldsymbol{\beta} + u.$$

By construction of the linear projection, the error  $u$  has the following properties:

1. It has mean zero:  $E(u) = 0$ .
2. It is uncorrelated with  $\mathbf{x}$ :  $\text{Cov}(\mathbf{x}, u) = 0$ .

# The error form for linear projections

Proof\*

Property #1:

Taking expectations of the linear projections in error form yields

$$E(y) = \beta_0 + E(\mathbf{x})\beta + E(u).$$

Comparing this with the first-order condition for  $\beta_0$ ,

$$\beta_0 = E[y] - E[\mathbf{x}]\beta \quad \Rightarrow \quad E(y) = \beta_0 + E(\mathbf{x})\beta$$

shows that

$$E(u) = 0.$$

# The error form for linear projections

Proof\*

Property #2:

Starting from the error form  $y = \beta_0 + \mathbf{x}\beta + u$  and subtracting its expectation  $E(y) = \beta_0 + E(\mathbf{x})\beta$  yields

$$y - E(y) = [\mathbf{x} - E(\mathbf{x})]\beta + u.$$

Multiply from the left by  $[\mathbf{x} - E(\mathbf{x})]'$  and take expectations:

$$\underbrace{E\{[\mathbf{x} - E(\mathbf{x})]'[y - E(y)]\}}_{\text{Cov}(\mathbf{x}', y)} = \underbrace{E\{[\mathbf{x} - E(\mathbf{x})]'[\mathbf{x} - E(\mathbf{x})]\}}_{\text{Var}(\mathbf{x})}\beta + \underbrace{E\{[\mathbf{x} - E(\mathbf{x})]'u\}}_{\text{Cov}(\mathbf{x}', u)}$$

Using the definition of  $\beta$  yields:

$$\text{Cov}(\mathbf{x}', y) = \text{Var}(\mathbf{x}) \underbrace{[\text{Var}(\mathbf{x})]^{-1} \text{Cov}(\mathbf{x}', y)}_{\beta} + \text{Cov}(\mathbf{x}', u)$$

and thus

$$\text{Cov}(\mathbf{x}', y) = \text{Cov}(\mathbf{x}', y) + \text{Cov}(\mathbf{x}', u) \quad \Rightarrow \quad \text{Cov}(\mathbf{x}', u) = 0.$$