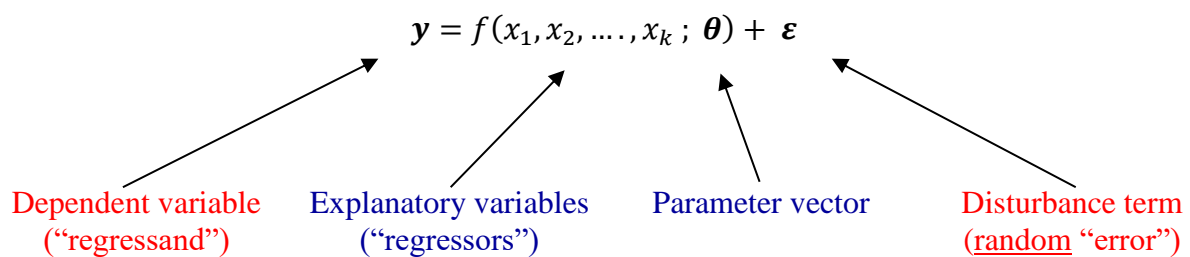


# 1. The OLS Estimator

OLS stands for “Ordinary Least Squares”. There are 6 assumptions ordinarily made, and the method of “fitting” a line through data is by least-squares. OLS is a common estimation methodology in Economics. While OLS is not appropriate for many contemporary economics problems, it is often the starting point for richer econometric models. OLS is also the pedagogical starting point for most econometric theory courses.

## 1.1 Population model and notation

Population “model” –



**Note:**

- The function, “ $f$ ”, may be linear or non-linear in the variables.
- The function, “ $f$ ”, may be linear or non-linear in the parameters.
- The function, “ $f$ ”, may be non-parametric, but we won’t consider this.
- We’ll focus on models that are parametric, and *usually* linear in the parameters.

**Questions:**

- Why is the error term needed?
- What is **random**, and what is **deterministic**?

What is **observable**, and what is **unobservable**?

## Examples:

1) Keynes' consumption function:

$$C = \beta_1 + \beta_2 Y + \varepsilon \quad (1)$$

2) Cobb-Douglas production function:

$$Y = AK^{\beta_2}L^{\beta_3}e^{\varepsilon} \quad (2)$$

By taking logs, the Cobb-Douglas production function can be rewritten as:

$$\log Y = \beta_1 + \beta_2 \log K + \beta_3 \log L + \varepsilon, \text{ where } \beta_1 = \log A$$

3) Wage equation:

$$\log Y = \beta_1 + \beta_2 \text{educ} + \beta_3 \text{age} + \dots + \varepsilon \quad (3)$$

## Sample Information

- Have a *sample* of “ $n$ ” observations:  $\{y_i ; x_{i1}, x_{i2}, \dots, x_{ik}\} ; i = 1, 2, \dots, n$
- We assume that these observed values are generated by the population model.

Let's take the case where the model is *linear in the parameters*:

$$y_i = \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + \varepsilon_i ; i = 1, \dots, n \quad (4)$$

Recall that the  $\beta$ 's and  $\varepsilon$  are unobservable. So,  $y_i$  is generated by 2 components:

1. Deterministic component:  $\sum_{j=1}^k \beta_j x_{ij}$ .
2. Stochastic component:  $\varepsilon_i$ .

So, the  $y_i$ 's must be “realized values” of a random variable.

Objectives:

- (i) Estimate unknown parameters
- (ii) Test hypotheses about parameters
- (iii) Predict values of  $y$  outside sample

## Interpreting the Parameters in a Model

Note that the  $\beta$ 's in equation (4) have an important economics interpretation:

$$\frac{\partial y_i}{\partial x_{1i}} = \beta_1; \text{ etc.}$$

The parameters are the *marginal effects* of the  $x$ 's on  $y$ , with other factors held constant (*ceteris paribus*). For example, from equation (1):

$$\partial C / \partial Y = \beta_2 = M.P.C.$$

We might wish to test the hypothesis that  $\beta_2 = 0.9$ , for example.

Depending on how the population model is specified, however, the  $\beta$ 's may *not* be interpreted as marginal effects. For example, after taking logs of the Cobb-Douglas production function in (2), we get the following population model:

$$\log Y = \beta_1 + \beta_2 \log K + \beta_3 \log L + \varepsilon,$$

and

$$\beta_2 = \frac{\partial Y / Y}{\partial K / K},$$

so that  $\beta_2$  is the elasticity of output with respect to capital. The point is that we need to be careful about how the parameters of the model are interpreted.

How could we test the hypothesis of constant returns to scale in the above Cobb-Douglas model?

So, we have a stochastic model that might be useful as a starting point to represent economics relationships. We need to be especially careful about the way in which we specify both parts of the model (the deterministic and stochastic parts).

## 1.2 Assumptions of OLS

All “models” are simplifications of reality. Presumably we want our model to be simple but “realistic” – able to explain actual data in a reliable and robust way.

To begin with we'll make a set of simplifying assumptions for our model. In fact, one of the main objectives of Econometrics is to re-consider these assumptions – are they realistic; can they be tested; what if they are wrong; can they be “relaxed”? The assumptions relate to: (1) functional form (parameters); (2) regressors; (3) disturbances.

### A.1: Linearity

The model is linear in the parameters:

$$y_i = \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + \varepsilon_i \quad ; \quad i = 1, \dots, n.$$

Linearity in the parameters allows the model to be written in matrix notation. Let,

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} ; \quad \mathbf{\beta} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} ; \quad X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1k} \\ x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix} ; \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix} .$$

$(n \times 1)$ 
 $(k \times 1)$ 
 $(n \times k)$ 
 $(n \times 1)$

Then, we can write the model, for the full sample, as:

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

If we take the  $i^{\text{th}}$  row (observation) of this model we have:

$$y_i = \mathbf{x}_i \boldsymbol{\beta} + \varepsilon_i \quad (\text{scalar})$$

### Notational points

- i. Vectors are in **bold**.
- ii. The dimensions of vectors/matrices are written (*rows*  $\times$  *columns*).
- iii. The first subscript denotes the row, the second subscript the column.
- iv. Some texts (including Greene, 2011), use the convention that vectors are columns. Hence, when an observation (row) is extracted from the  $X$  matrix, it is transformed into a column. Hence, the above equation would be expressed as  $y_i = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i$ .

### **A.2: Full Rank**

We assume that there are no exact linear dependencies among the columns of  $X$  (if there were, then one or more regressor is redundant). Note that  $X$  is  $(n \times k)$  and  $\text{Rank}(X) = k$ . So we are also implicitly assuming that  $n > k$ , since  $\text{Rank}(A) \leq \min\{\#\text{rows}, \#\text{cols}\}$ .

What does this assumption really mean? Suppose we had:

$$y_i = \beta_1 x_{i1} + \beta_2 (2x_{i1}) + \varepsilon_i$$

We can only identify, and estimate, the one function,  $(\beta_1 + 2\beta_2)$ . In this model,  $\text{Rank}(X) = k - 1 = 1$ . An example which is commonly found in undergraduate textbooks, of where A.2 is violated, is the dummy variable trap.

### **A.3: Errors Have a Zero Mean**

Assume that, *in the population*,  $E(\varepsilon_i) = 0 \quad ; \quad i = 1, 2, \dots, n$ . So,

$$E(\boldsymbol{\varepsilon}) = E \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix} = \mathbf{0} .$$

#### A.4: Spherical Errors

Assume that, in the population, the disturbances are generated by a process whose variance is constant ( $\sigma^2$ ), and that these disturbances are uncorrelated with each other:

$$var(\varepsilon_i) = \sigma^2 ; i = 1, 2, \dots, n \quad (\text{Homoskedasticity})$$

$$cov(\varepsilon_i, \varepsilon_j) = 0 ; \forall i \neq j \quad (\text{no Autocorrelation})$$

Putting these assumptions together it can be shown that the “covariance matrix” for the random vector,  $\boldsymbol{\varepsilon}$ , is:

$$V(\boldsymbol{\varepsilon}) = E [(\boldsymbol{\varepsilon} - E(\boldsymbol{\varepsilon}))(\boldsymbol{\varepsilon} - E(\boldsymbol{\varepsilon}))'] = E[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'] = \begin{bmatrix} E(\varepsilon_1\varepsilon_1) & \cdots & E(\varepsilon_1\varepsilon_n) \\ \vdots & \ddots & \vdots \\ E(\varepsilon_n\varepsilon_1) & \cdots & E(\varepsilon_n\varepsilon_n) \end{bmatrix}$$

but...

$$E(\varepsilon_i\varepsilon_i) = E(\varepsilon_i^2) = E[(\varepsilon_i - 0)^2] = var(\varepsilon_i) = \sigma^2$$

and

$$E(\varepsilon_i\varepsilon_j) = E[(\varepsilon_i - 0)(\varepsilon_j - 0)] = cov(\varepsilon_i, \varepsilon_j) = 0.$$

So:

$$V(\boldsymbol{\varepsilon}) = \begin{bmatrix} \sigma^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma^2 \end{bmatrix} = \sigma^2 I_n$$

a scalar matrix.

#### A.5: Generating Process for $\mathbf{X}$

The classical regression model assumes that the regressors are “fixed in repeated samples” (laboratory situation). We can assume this – very strong, though.

Alternatively, allow  $\mathbf{x}$ 's to be random, but restrict the form of their randomness – assume that the regressors are uncorrelated with the disturbances. The process that generates  $\mathbf{X}$  is unrelated to the process that generates  $\boldsymbol{\varepsilon}$  in the population.

## A.6: Normality of Errors

$$(\boldsymbol{\varepsilon}|X) \sim N[0, \sigma^2 I_n]$$

This assumption is not as strong as it seems:

- often reasonable due to the Central Limit Theorem (C.L.T.)
- often not needed
- when some distributional assumption is needed, often a more general one is ok

### Summary

The classical linear regression model is:

- $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
- $(\boldsymbol{\varepsilon}|X) \sim N[0, \sigma^2 I_n]$
- $\text{Rank}(X) = k$
- Data generating processes (D.G.P.s) of  $X$  and  $\boldsymbol{\varepsilon}$  are unrelated.

Implications for  $\mathbf{y}$  (if  $X$  is non-random; *or* conditional on  $X$ ):

$$E(\mathbf{y}) = X\boldsymbol{\beta} + E(\boldsymbol{\varepsilon}) = X\boldsymbol{\beta}$$

$$V(\mathbf{y}) = V(\boldsymbol{\varepsilon}) = \sigma^2 I_n$$

Because *linear* transformations of a Normal random variable are themselves Normal, we also have:  $\mathbf{y} \sim N[X\boldsymbol{\beta}, \sigma^2 I_n]$ .

### Some Questions

- How reasonable are the assumptions associated with the classical linear regression model?
- How do these assumptions affect the estimation of the model's parameters?
- How do these assumptions affect the way we test hypotheses about the model's parameters?
- Which of these assumptions are used to establish the various results we'll be concerned with?
- Which assumptions can be "relaxed" without affecting these results?

### 1.3 Deriving the OLS estimator

Our first task is to estimate the parameters of our model,

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad \boldsymbol{\varepsilon} \sim N[\mathbf{0}, \sigma^2 I_n] .$$

Note that there are  $(k + 1)$  parameters, including  $\sigma^2$ .

- Many possible procedures for estimating parameters.
- Choice should be based not only on computational convenience, but also on the “[sampling properties](#)” of the resulting estimator.
- To begin with, consider *one possible* estimation strategy – **Least Squares**.

For the  $i^{\text{th}}$  data-point, we have:

$$y_i = \mathbf{x}'_i \boldsymbol{\beta} + \varepsilon_i \quad ,$$

and the population regression is:

$$E(y_i | \mathbf{x}'_i) = \mathbf{x}'_i \boldsymbol{\beta} \quad .$$

We'll estimate  $E(y_i | \mathbf{x}'_i)$  by

$$\hat{y}_i = \mathbf{x}'_i \mathbf{b} .$$

*In the population*, the true (unobserved) disturbance is  $\varepsilon_i$  [ $= y_i - \mathbf{x}'_i \boldsymbol{\beta}$ ].

When we use  $\mathbf{b}$  to estimate  $\boldsymbol{\beta}$ , there will be some “estimation error”, and the value,  $e_i = y_i - \mathbf{x}'_i \mathbf{b}$  will be called the  $i^{\text{th}}$  “**residual**”.

So,

$$y_i = (\mathbf{x}'_i \boldsymbol{\beta} + \varepsilon_i) = (\mathbf{x}'_i \mathbf{b} + e_i) = (\hat{y}_i + e_i)$$

unobserved                      observed

[Population]                      [Sample]

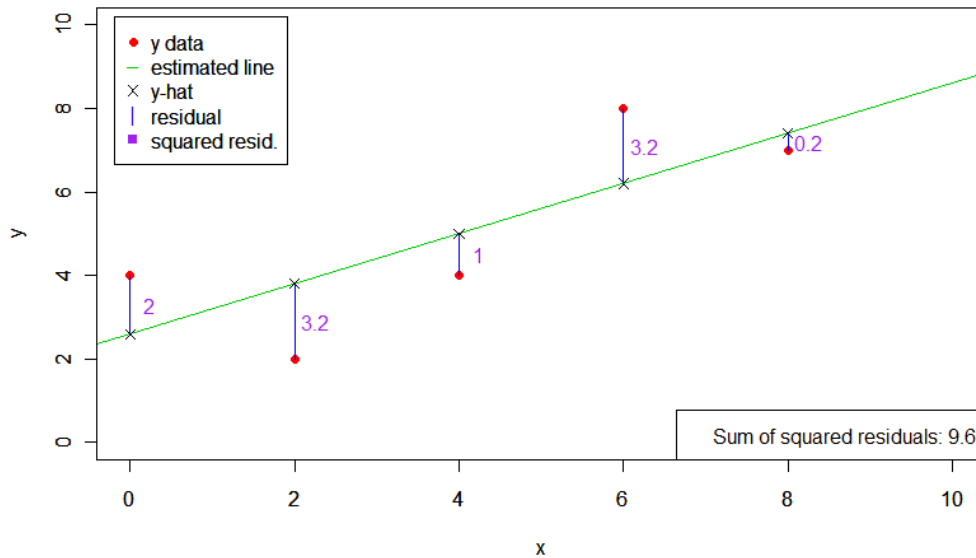
#### The Least Squares Criterion

“Choose  $\mathbf{b}$  so as to minimize the sum of the squared residuals.”

- Why *squared* residuals?
- Why not *absolute values* of residuals?

- Why not use a “minimum distance” criterion?

Fig 1.1. Minimizing the sum of squared residuals, for  $y = \{4, 2, 4, 8, 7\}$ ;  $x = \{0, 2, 4, 6, 8\}$ .



### Minimizing the Sum of Squared Residuals: An Optimization Problem

$$\begin{aligned} \text{Min.}_{(b)} \sum_{i=1}^n e_i^2 &\Leftrightarrow \text{Min.}_{(b)} (\mathbf{e}'\mathbf{e}) \\ &\Leftrightarrow \text{Min.}_{(b)} [(\mathbf{y} - \mathbf{X}\mathbf{b})'(\mathbf{y} - \mathbf{X}\mathbf{b})]. \end{aligned}$$

Now, let:

$$S = (\mathbf{y} - \mathbf{X}\mathbf{b})'(\mathbf{y} - \mathbf{X}\mathbf{b}) = \mathbf{y}'\mathbf{y} - \mathbf{b}'\mathbf{X}'\mathbf{y} - \mathbf{y}'\mathbf{X}\mathbf{b} + \mathbf{b}'\mathbf{X}'\mathbf{X}\mathbf{b}.$$

Note that,

$$\begin{aligned} \mathbf{b}'\mathbf{X}'\mathbf{y} &= \mathbf{y}'\mathbf{X}\mathbf{b}. \\ (1 \times k)(k \times n)(n \times 1) &\quad (1 \times 1) \end{aligned}$$

So,  $S = \mathbf{y}'\mathbf{y} - 2(\mathbf{y}'\mathbf{X})\mathbf{b} + \mathbf{b}'(\mathbf{X}'\mathbf{X})\mathbf{b}.$

Note:

(i)  $\partial(\mathbf{a}'\mathbf{x})/\partial\mathbf{x} = \mathbf{a}$

(ii)  $\partial(\mathbf{x}'\mathbf{A}\mathbf{x})/\partial\mathbf{x} = 2\mathbf{A}\mathbf{x}$  ; if A is *symmetric*



Applying these 2 results –

$$\partial S / \partial \mathbf{b} = \mathbf{0} - 2(\mathbf{y}'X)' + 2(X'X)\mathbf{b} = 2[X'X\mathbf{b} - X'\mathbf{y}] .$$

Set this to zero (for a turning point):

$$\begin{array}{lcl} X'X\mathbf{b} & = & X'\mathbf{y} , & (k \text{ equations in } k \text{ unknowns}) \\ (k \times n)(n \times k)(k \times 1) & & (k \times n)(n \times 1) & (\text{the “normal equations”}) \end{array}$$

so:

$$\mathbf{b} = (X'X)^{-1}X'\mathbf{y} \quad ; \quad \text{provided that } (X'X)^{-1} \text{ exists}$$

Notice that  $X'X$  is  $(k \times k)$ , and  $\text{rank}(X'X) = \text{rank}(X) = k$  (assumption).

This implies that  $(X'X)^{-1}$  exists.

We need the “full rank” assumption for the Least Squares estimator,  $\mathbf{b}$ , to *exist*.

*None of our other assumptions have been used so far.*

**Check** – have we *minimized*  $S$  ?

$$\left( \frac{\partial^2 S}{\partial \mathbf{b} \partial \mathbf{b}'} \right) = \partial / \partial \mathbf{b}' [2X'X\mathbf{b} - 2X'\mathbf{y}] = 2(X'X) \quad ; \quad \text{a } (k \times k) \text{ matrix.}$$

Note that  $X'X$  is at least positive *semi-definite* –

$$\eta'(X'X)\eta = (X\eta)'(X\eta) = (\mathbf{u}'\mathbf{u}) = \sum_{i=1}^n u_i^2 \geq 0 \quad ;$$

and so if  $X'X$  has full rank, it will be *positive-definite*, not negative-definite.

So, our assumption that  $X$  has full rank has two implications –

1. The Least Squares estimator,  $\mathbf{b}$ , *exists*.
2. Our optimization problem leads to the *minimization* of  $S$ , not its maximization!

### Aside – OLS formula in scalar form

For a population model with an intercept and a single regressor, you may have seen the following formulas used in undergraduate textbooks:

$$b_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{s_{X,Y}}{s_X^2} ,$$

$$b_0 = \bar{Y} - b_1 \bar{X} ,$$

where  $s_{X,Y}$  is the sample covariance between  $X_i$  and  $Y_i$ , and  $s_X^2$  is the sample variance of  $X_i$ .