

Regression Analysis I

Hüseyin Taştan¹

¹Department of Economics
Ph.D. Program - Advanced Econometrics
Yildiz Technical University

September 24, 2013

Linear Regression Model

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \dots + \beta_k x_{ik} + u_i, \quad i = 1, \dots, n$$

where y_i is i th observation on the dependent variable and x_{ij} , $j = 1, \dots, k$, is the i th observation on the j th explanatory variable, and u_i is the i th random error term. There are k unknown parameters and n equations which can also be written as:

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

Linear Regression Model

Using matrix notation:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & x_{12} & \dots & x_{1k} \\ 1 & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n2} & \dots & x_{nk} \end{bmatrix}, \boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix}, \mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}$$

and rewrite

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

Linear Regression Model

Equivalent notation: denote i th observation by

$$\mathbf{x}_i = [1 \quad x_{i2} \quad x_{i3} \quad \dots \quad x_{ik}]^\top$$

Using this notation CLRM becomes

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + u_i, \quad i = 1, \dots, n$$

Classical Regression Model

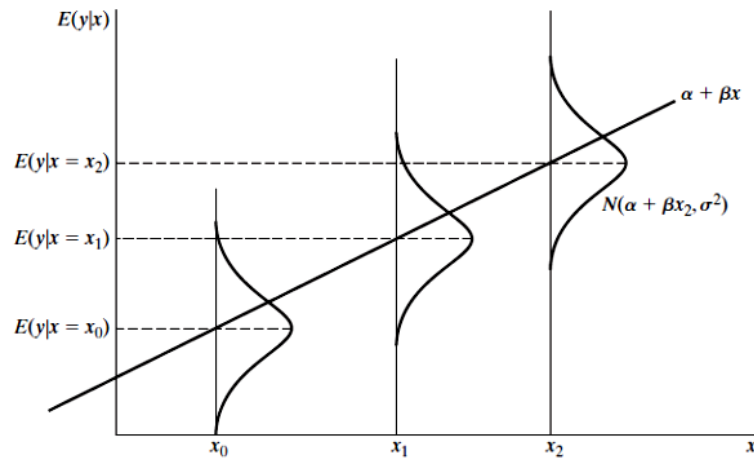


Figure 1 : Population Regression Function

Linear Regression Model

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

Assumptions of the CLRM (\mathbf{X} variables assumed to be stochastic, not fixed in repeated samples)

1. Linearity in parameters: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$
2. $\text{rank}(\mathbf{X}) = k$, (no perfect multicollinearity, columns of \mathbf{X} are linearly independent)
3. $E[\mathbf{u}|\mathbf{X}] = \mathbf{0}_{n \times 1}$, (zero conditional mean for the errors)
4. $\text{Var}[\mathbf{u}|\mathbf{X}] = E(\mathbf{u}\mathbf{u}^\top) = \sigma^2 \mathbf{I}_n$, (no heteroskedasticity, no autocorrelation)
5. $\mathbf{u}|\mathbf{X} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$, (errors follow multivariate normal distribution)

Adding the last assumption the model is called Classical Normal Linear Regression Model.

Ordinary Least Squares (OLS) Estimator

Sample Regression Function (SRF):

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

$\hat{\boldsymbol{\beta}}$ $k \times 1$: OLS estimator

$\hat{\mathbf{u}}$ $n \times 1$: residuals

OLS principle: minimize sum of squared residuals:

$$\hat{\boldsymbol{\beta}} = \arg \min_{\mathbf{b}} SSR(\mathbf{b})$$

where

$$SSR(\hat{\boldsymbol{\beta}}) = \sum_{i=1}^n \hat{u}_i^2 = \hat{\mathbf{u}}^\top \hat{\mathbf{u}}$$

OLS Estimator

$$SSR(\hat{\boldsymbol{\beta}}) = \sum_{i=1}^n \hat{u}_i^2 = \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}})^2$$

Using the first notation

$$\min_{\hat{\boldsymbol{\beta}}} SSR(\hat{\boldsymbol{\beta}}) = \hat{\mathbf{u}}^\top \hat{\mathbf{u}}$$

or the second notation:

$$\min_{\hat{\boldsymbol{\beta}}} SSR(\hat{\boldsymbol{\beta}}) = \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}})^2$$

From the first notation expanding SSR we get:

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

OLS Estimator

We need to evaluate the first derivative of SSR w.r.t. $\hat{\beta}$. Note that the second term is a linear combination and the last is a quadratic form. In general, for a $k \times 1$ vector z , and $k \times n$ matrix A , and $k \times k$ matrix B , we have:

$$\frac{\partial(z^\top A)}{\partial z} = A,$$

and

$$\frac{\partial(z^\top Bz)}{\partial z} = 2Bz$$

OLS Estimator

For example, let $z = \begin{bmatrix} z_1 & z_2 \end{bmatrix}$ and

$$A = \begin{bmatrix} 0 & 1 & -2 \\ 1 & 2 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$$

$$z^\top A = \begin{bmatrix} z_1 & z_2 \end{bmatrix} \begin{bmatrix} 0 & 1 & -2 \\ 1 & 2 & 0 \end{bmatrix} = \begin{bmatrix} z_2 & z_1 + 2z_2 & -2z_1 \end{bmatrix}$$

First derivatives wrt z_1 and z_2 :

$$\frac{\partial(z^\top A)}{\partial z_1} = \begin{bmatrix} 0 & 1 & -2 \end{bmatrix} \text{ and } \frac{\partial(z^\top A)}{\partial z_2} = \begin{bmatrix} 1 & 2 & 0 \end{bmatrix}$$

Collecting in a vector

$$\begin{aligned} \frac{\partial(z^\top A)}{\partial z} &= \begin{bmatrix} \frac{\partial(z^\top A)}{\partial z_1} \\ \frac{\partial(z^\top A)}{\partial z_2} \end{bmatrix} \\ &= \begin{bmatrix} 0 & 1 & -2 \\ 1 & 2 & 0 \end{bmatrix} = A \end{aligned}$$

OLS Estimator

Now take the derivative of $z^\top Bz$ wrt z vector.

$$\begin{aligned} z^\top Bz &= \begin{bmatrix} z_1 & z_2 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \\ &= 2z_1^2 + 2z_1z_2 + 2z_2^2 \end{aligned}$$

1st derivatives:

$$\begin{aligned} \frac{\partial(z^\top Bz)}{\partial z} &= \begin{bmatrix} \frac{\partial(z^\top Bz)}{\partial z_1} \\ \frac{\partial(z^\top Bz)}{\partial z_2} \end{bmatrix} \\ &= \begin{bmatrix} 4z_1 + 2z_2 \\ 2z_1 + 4z_2 \end{bmatrix} \\ &= 2 \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = 2Bz \end{aligned}$$

OLS Estimator

Back to SSR:

$$SSR(\hat{\beta}) = y^\top y - 2\hat{\beta}^\top X^\top y + \hat{\beta}^\top X^\top X \hat{\beta}$$

Hint: $\hat{\beta} = z$, $X^\top y = A$ ve $X^\top X = B$

OLS FOC:

$$\frac{\partial SSR(\hat{\beta})}{\partial \hat{\beta}} = -2X^\top y + 2X^\top X \hat{\beta} = 0_k$$

Normal equations:

$$X^\top X \hat{\beta} = X^\top y$$

Using the second CLRM assumption:

$$\text{rank}(X) = \text{rank}(X^\top X) = k$$

we can find the inverse of $X^\top X$. Multiplying both sides of normal equations by $(X^\top X)^{-1}$ OLS estimator is

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

OLS Estimator

Using the second notation:

$$\min_{\hat{\beta}} SSR(\hat{\beta}) = \sum_{i=1}^n (y_i - \mathbf{x}_i^{\top} \hat{\beta})^2$$

In this case FOC:

$$\frac{\partial SSR(\hat{\beta})}{\partial \hat{\beta}} = -2 \sum_{i=1}^n \mathbf{x}_i (y_i - \mathbf{x}_i^{\top} \hat{\beta}) = \mathbf{0}_k$$

and normal equations

$$\left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^{\top} \right) \hat{\beta} = \sum_{i=1}^n \mathbf{x}_i y_i$$

OLS estimator can be written as:

$$\begin{aligned} \hat{\beta} &= \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^{\top} \right)^{-1} \sum_{i=1}^n \mathbf{x}_i y_i \\ &\equiv (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{y} \end{aligned}$$

OLS Estimator

$k \times k$ matrix $\mathbf{x}_i \mathbf{x}_i^{\top}$ has the following elements:

$$\begin{aligned} \mathbf{x}_i \mathbf{x}_i^{\top} &= \begin{bmatrix} 1 \\ x_{i2} \\ \vdots \\ x_{ik} \end{bmatrix} \begin{bmatrix} 1 & x_{i2} & \dots & x_{ik} \end{bmatrix} \\ &= \begin{bmatrix} 1 & x_{i2} & x_{i3} & \dots & x_{ik} \\ x_{i2} & x_{i2}^2 & x_{i2}x_{i3} & \dots & x_{i2}x_{ik} \\ x_{i3} & x_{i3}x_{i2} & x_{i3}^2 & \dots & x_{i3}x_{ik} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{ik} & x_{ik}x_{i2} & x_{ik}x_{i3} & \dots & x_{ik}^2 \end{bmatrix} \end{aligned}$$

We have a sequence of n matrices defined analogously.

OLS Estimator

Summing these matrices

$$\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^{\top} = \begin{bmatrix} n & \sum x_{i2} & \sum x_{i3} & \dots & \sum x_{ik} \\ \sum x_{i2} & \sum x_{i2}^2 & \sum x_{i2}x_{i3} & \dots & \sum x_{i2}x_{ik} \\ \sum x_{i3} & \sum x_{i3}x_{i2} & \sum x_{i3}^2 & \dots & \sum x_{i3}x_{ik} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sum x_{ik} & \sum x_{ik}x_{i2} & \sum x_{ik}x_{i3} & \dots & \sum x_{ik}^2 \end{bmatrix} = \mathbf{X}^{\top} \mathbf{X}$$

This matrix is square, symmetric and positive definite. Similarly,

$$\begin{aligned} \sum_{i=1}^n \mathbf{x}_i y_i &= \sum_{i=1}^n \begin{bmatrix} 1 \\ x_{i2} \\ \vdots \\ x_{ik} \end{bmatrix} y_i = \begin{bmatrix} \sum y_i \\ \sum x_{i2} y_i \\ \vdots \\ \sum x_{ik} y_i \end{bmatrix} \\ &= \mathbf{X}^{\top} \mathbf{y} \end{aligned}$$

OLS Estimation

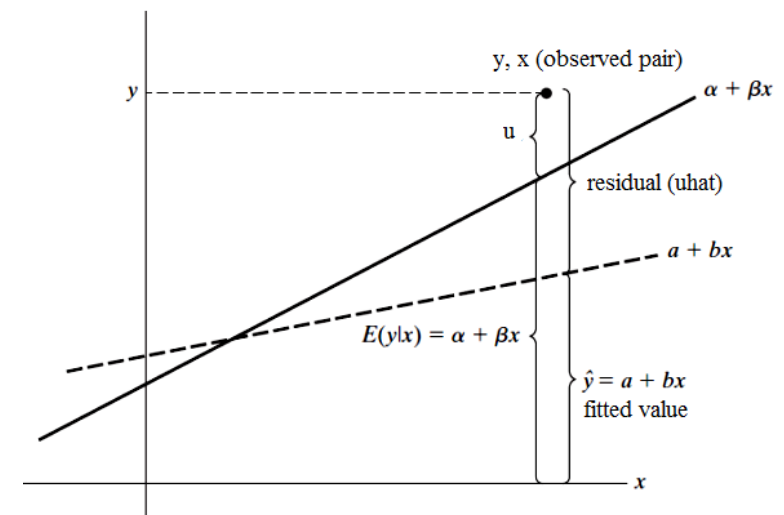


Figure 2 : Sample Regression Function (solid line is population regression function)

OLS Estimator

Example

Consider the following model where there are no explanatory variables:

$$y_i = \beta_1 + u_i, \quad i = 1, \dots, n$$

Now \mathbf{X} matrix only has a column of ones ($n \times 1$) which is denoted by \mathbf{z} :

$$\mathbf{z} = \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix}^\top = \mathbf{X}$$

Now the OLS estimator of β_1 :

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

OLS Estimator

Example (cont.d)

Note that

$$\mathbf{z}^\top \mathbf{z} = \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} = n$$

and

$$\mathbf{z}^\top \mathbf{y} = \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \sum_{i=1}^n y_i$$

OLS estimator is simply

$$\beta_1 = n^{-1} \sum_{i=1}^n y_i \equiv \bar{y}$$

the arithmetic mean.

OLS Estimation

Example

One dummy variable and an intercept model:

$$y_i = \delta_0 + \delta_1 D_i + u_i, \quad i = 1, \dots, n$$

Data consist of 5 observations for simplicity:

$$\mathbf{y} = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}, \quad D_i = \begin{cases} 1, & \text{if } y_i \leq 3; \\ 0, & \text{otherwise.} \end{cases} \Leftrightarrow \mathbf{D} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

Design matrix

$$\mathbf{X} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \\ 1 & 0 \\ 1 & 0 \end{bmatrix}$$

Parameter vector: $\hat{\beta} = [\hat{\delta}_0 \quad \hat{\delta}_1]^\top$ Find the OLS solution.

OLS Estimation

Example (cont.d)

$$\mathbf{X}^\top \mathbf{X} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \\ 1 & 0 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} 5 & 3 \\ 3 & 3 \end{bmatrix}$$

$$(\mathbf{X}^\top \mathbf{X})^{-1} = \begin{bmatrix} 1/2 & -1/2 \\ -1/2 & 5/6 \end{bmatrix}, \quad \mathbf{X}^\top \mathbf{y} = \begin{bmatrix} \sum y_i \\ \sum y_i D_i \end{bmatrix} = \begin{bmatrix} 15 \\ 6 \end{bmatrix}$$

$$\hat{\beta} = \begin{bmatrix} 1/2 & -1/2 \\ -1/2 & 5/6 \end{bmatrix} \begin{bmatrix} 15 \\ 6 \end{bmatrix} = \begin{bmatrix} 4.5 \\ -2.5 \end{bmatrix}$$

Sample regression function:

$$\hat{y}_i = 4.5 - 2.5 D_i$$

In this simple example, the intercept is simply the arithmetic mean of the base group (y is greater than 3) $((4 + 5)/2 = 4.5)$. The arithmetic mean of the other group is 2. The parameter estimate on D is simply the difference between these group means (-2.5) .

OLS Estimation

Example (cont.d)

Fitted values and residuals:

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \\ 1 & 0 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 4.5 \\ -2.5 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \\ 2 \\ 4.5 \\ 4.5 \end{bmatrix},$$

$$\hat{\mathbf{u}} = \mathbf{y} - \hat{\mathbf{y}} = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix} - \begin{bmatrix} 2 \\ 2 \\ 2 \\ 4.5 \\ 4.5 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ 1 \\ -0.5 \\ 0.5 \end{bmatrix}$$

OLS Estimation

Now the **same example** without the intercept:

$$y_i = \gamma_0 D_{i1} + \delta_0 D_{i2} + u_i, \quad i = 1, \dots, n$$

Now \mathbf{X} matrix and inner product matrix:

$$\mathbf{X} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{X}^\top \mathbf{X} = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix}$$

$$(\mathbf{X}^\top \mathbf{X})^{-1} = \begin{bmatrix} 1/3 & 0 \\ 0 & 1/2 \end{bmatrix}, \quad \mathbf{X}^\top \mathbf{y} = \begin{bmatrix} \sum y_i D_{i1} \\ \sum y_i D_{i2} \end{bmatrix} = \begin{bmatrix} 6 \\ 9 \end{bmatrix}$$

OLS estimates

$$\hat{\boldsymbol{\beta}} = \begin{bmatrix} \hat{\gamma}_0 \\ \hat{\delta}_0 \end{bmatrix} = \begin{bmatrix} 1/3 & 0 \\ 0 & 1/2 \end{bmatrix} \begin{bmatrix} 6 \\ 9 \end{bmatrix} = \begin{bmatrix} 2 \\ 4.5 \end{bmatrix}$$

Fitted regression equation:

$$\hat{y}_i = 2D_{i1} + 4.5D_{i2}$$

OLS Estimation

Example

Adding an intercept into the previous model we get:

$$y_i = \beta_0 + \gamma_0 D_{i1} + \delta_0 D_{i2} + u_i, \quad i = 1, \dots, n$$

$$\mathbf{X} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}$$

Obviously, the columns of \mathbf{X} are not linearly independent:
 $\text{rank}(\mathbf{X}) < 3$. Another way of seeing this:

$$\mathbf{X}^\top \mathbf{X} = \begin{bmatrix} 5 & 3 & 2 \\ 3 & 3 & 0 \\ 2 & 0 & 2 \end{bmatrix}, \quad |\mathbf{X}^\top \mathbf{X}| = 0$$

Again this matrix is **not full rank**.

Exercise

Single explanatory variable and a constant term:

$$y_i = \beta_1 + \beta_2 x_i + u_i, \quad i = 1, \dots, n$$

using matrix framework show that

$$\hat{\beta}_1 = \frac{\sum x_i^2 \sum y_i - \sum x_i \sum y_i x_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\hat{\beta}_2 = \frac{n \sum y_i x_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

Geometry of OLS

- ▶ Let us focus on numerical properties of OLS estimates.
- ▶ These properties are always valid regardless of the data generating process and related assumptions.
- ▶ In contrast, statistical properties, such as unbiasedness, efficiency, consistency, asymptotic normality etc, all depend on the validity of certain assumptions.
- ▶ To understand the geometry of OLS let us review basic concepts in Euclidean geometry.
- ▶ n -vector is defined as a column vector with n elements which can also be represented by a $n \times 1$ matrix.
- ▶ Euclidean space in n dimensions is denoted with \mathbb{E}^n .
- ▶ The set of n -vectors can also be denoted by \mathbb{R}^n (real line, \mathbb{R}).
- ▶ The difference is that wider set of operations are defined on \mathbb{E}^n .

Review of Euclidean Geometry

Figure 3 displays $\mathbf{x} = [x_1 \ x_2]^\top \in \mathbb{E}^2$ which is also known as the Cartesian coordinates.

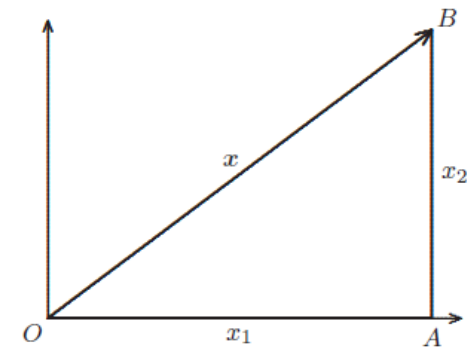


Figure 3 : A vector \mathbf{x} in 2-dim space

What is the length of \mathbf{x} ?

Review of Euclidean Geometry

- ▶ Let $\mathbf{x}, \mathbf{y} \in \mathbb{E}^n$, then the inner product is defined as

$$\mathbf{x}^\top \mathbf{y} \equiv \mathbf{y}^\top \mathbf{x}$$

which is obviously commutative.

- ▶ The length of any vector is defined as

$$\|\mathbf{x}\| = \sqrt{\mathbf{x}^\top \mathbf{x}} \equiv \left(\sum_{i=1}^n x_i^2 \right)^{1/2}$$

- ▶ Applying this to the case $n = 2$ one can see the length of the vector is just the hypotenuse formed by the coordinates of vector \mathbf{x} (see Figure 3).

Review of Euclidean Geometry

Addition: Let $\mathbf{x}, \mathbf{y} \in \mathbb{E}^2$, then

$$\mathbf{x} + \mathbf{y} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \end{bmatrix}$$

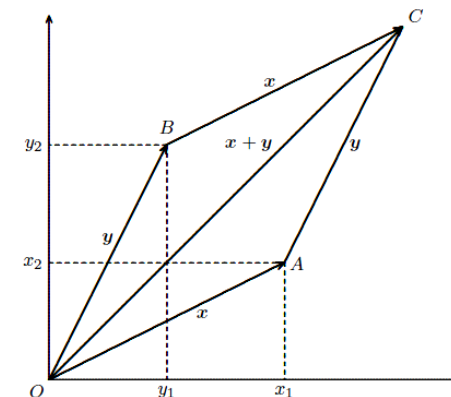


Figure 4 : Addition of two vectors

Review of Euclidean Geometry

Multiplication by a scalar is shown in Figure 5

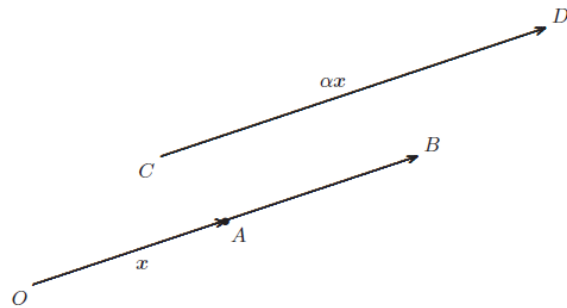


Figure 5 : Multiplication by a scalar

Let α be any scalar, then the length of x is

$$\|\alpha x\| = |\alpha| \sqrt{x^\top x} = |\alpha| \|x\|$$

Review of Euclidean Geometry

- ▶ Let us consider two vectors w and z both of length 1:

$$w = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \|z\| = 1$$

which is represented by a horizontal line in Figure 6.

- ▶ We know the length of z is 1 but what are the coordinates?
- ▶ Obviously from elementary trigonometry we have

$$z_1 = \cos \theta, \quad z_2 = \sin \theta$$

- ▶ From Pythagoras' Theorem

$$z_1^2 + z_2^2 \equiv \|z\|^2 = \cos^2 \theta + \sin^2 \theta = 1$$

Review of Euclidean Geometry

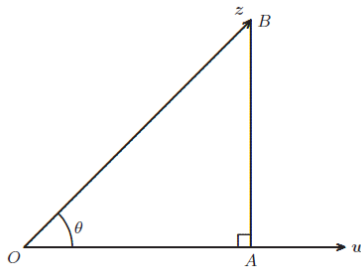


Figure 6 : The Angle between two vectors

Note in Figure 6: line OB represents z , the length of AB is $\sin \theta$, the length of OA is $\cos \theta$. What is the scalar product of w and z ?

$$w^\top z = \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} = \cos \theta$$

This result only holds for vectors of length 1.

Review of Euclidean Geometry

More generally, let

$$x = \alpha w, \quad y = \gamma z, \quad \alpha > 0, \quad \gamma > 0$$

The lengths are

$$\|x\| = \alpha, \quad \|y\| = \gamma$$

Thus,

$$\langle x, y \rangle = x^\top y = \alpha \gamma w^\top z$$

Because x is parallel to w , and y is parallel to z , the angle between x and y is the same as that between w and z , ie, θ . Therefore,

$$\langle x, y \rangle = \|x\| \|y\| \cos \theta$$

Cosine of the angle between two vectors measures how close two vectors are in terms of their directions.

Review of Euclidean Geometry

- Recall

$$-1 \leq \cos \theta \leq 1$$

and if θ is measured in radians we have

$$\cos 0 = 1, \quad \cos(\pi/2) = 0, \quad \cos \pi = -1$$

- Thus, $\cos \theta = 1$ for vectors that are parallel, 0 for vectors that are right angles (90°) to each other, and -1 for vectors that point in directly opposite directions.
- If the angle is $\pi/2$ (in radians) or 90° (in degrees) then the inner product is zero:

$$\mathbf{x}^\top \mathbf{y} = 0.$$

- This kind of vectors are said to be **orthogonal**.

Review of Euclidean Geometry

From

$$\langle \mathbf{x}, \mathbf{y} \rangle = \|\mathbf{x}\| \|\mathbf{y}\| \cos \theta$$

and

$$-1 \leq \cos \theta \leq 1$$

we have

Definition (Cauchy-Schwartz inequality)

$$|\mathbf{x}^\top \mathbf{y}| \leq \|\mathbf{x}\| \|\mathbf{y}\|$$

This says that the inner product can never be greater than the product of the lengths of vectors. If \mathbf{x} and \mathbf{y} parallel to each other the inequality becomes the equality.

Geometry of OLS

- In the regression model

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

the dependent variable and each column of \mathbf{X} can be thought of as vectors in \mathbb{E}^n .

- Further, the elements of n -vector can be thought of as the coordinates of a point in \mathbb{E}^n .
- Note that there are three vectors in the regression model: \mathbf{y} , $\mathbf{X}\boldsymbol{\beta}$ and \mathbf{u} .
- A **subspace** of \mathbb{E}^n can be defined in terms of a set of **basis** vectors.
- We are particularly interested in subspaces defined by the columns of \mathbf{X} as the basis.
- There are k basis vectors in \mathbf{X} : $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$. The subspace associated with these basis vectors is denoted $\mathcal{S}(\mathbf{X})$ or $\mathcal{S}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k)$.

Geometry of OLS

- The subspace $\mathcal{S}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k)$ consists of every vector that can be formed as a linear combination of $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$. This can be formally defined as follows

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

This is called **subspace spanned by** $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$ or the **column space** of \mathbf{X} .

- The **orthogonal complement** of $\mathcal{S}(\mathbf{X})$, denoted $\mathcal{S}^\perp(\mathbf{X})$, is the set of all vectors \mathbf{w} in \mathbb{E}^n that are orthogonal to everything in $\mathcal{S}(\mathbf{X})$. This implies that for every \mathbf{z} in $\mathcal{S}(\mathbf{X})$, $\mathbf{w}^\top \mathbf{z} = 0$:

$$\mathcal{S}^\perp(\mathbf{X}) \equiv \left\{ \mathbf{w} \in \mathbb{E}^n : \mathbf{w}^\top \mathbf{z} = 0 \right\}.$$

- If the dimension of $\mathcal{S}(\mathbf{X})$ is k then the dimension of $\mathcal{S}^\perp(\mathbf{X})$ is $n - k$. Figure 7 shows the case where $n = 2, k = 1$.

Geometry of OLS

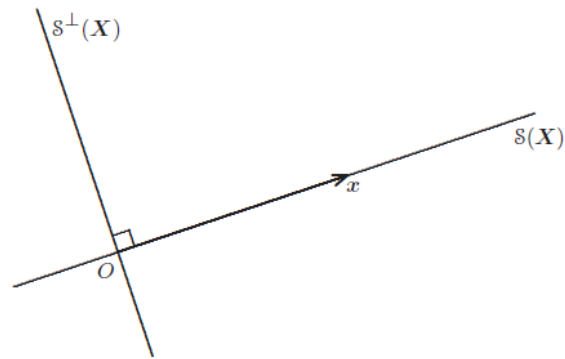


Figure 7 : Span of X and its orthogonal complement ($n = 2, k = 1$)

Geometry of OLS

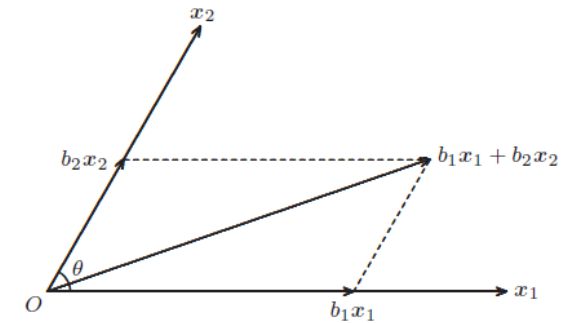


Figure 8 : 2-dimensional subspace

Figure 8 displays a 2-dimensional subspace.

Geometry of OLS

- ▶ Figure 9 represents linear regression model geometrically.
- ▶ The horizontal direction is chosen for the vector $X\beta$, and y and u are shown in the plane. Notice that u the error vector, is not orthogonal to $X\beta$.
- ▶ In order for the OLS estimator to be defined the square matrix $X^\top X$ must be invertible, ie nonsingular (full rank). This is another way of saying that the columns of X must be linearly independent. x_j is said to be linearly dependent if we can write it as a linear combination of other vectors in X :

$$x_j = \sum_{i \neq j} c_i x_i$$

Or,

$$Xb = 0_n$$

- ▶ If this is the case then premultiplying by X^\top

$$X^\top Xb = 0.$$

implying that $X^\top X$ cannot be inverted.

- ▶ If the k columns of X are not linearly independent then they will span a subspace of dimension less than k . This is called the rank of X .

Geometry of OLS

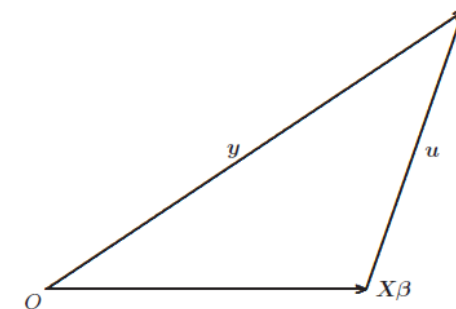


Figure 9 : Geometry of linear regression model

Geometry of OLS

- ▶ Note that every $\mathbf{X}\beta$ belongs to $\mathcal{S}(\mathbf{X})$, span of \mathbf{X} .
- ▶ The vector $\mathbf{X}\hat{\beta}$ constructed using the OLS estimator $\hat{\beta}$ belongs to this subspace. $\hat{\beta}$ was obtained by solving the system of equations represented by

$$\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\hat{\beta}) = 0$$

The i th element is

$$x_i^\top(\mathbf{y} - \mathbf{X}\hat{\beta}) = \langle x_i^\top, \mathbf{y} - \mathbf{X}\hat{\beta} \rangle$$

- ▶ The vector $\mathbf{y} - \mathbf{X}\hat{\beta}$ is orthogonal to all of the regressors.
- ▶ In other words the residual vector, $\hat{\mathbf{u}}$, is orthogonal to all the regressors. This is depicted in Figure 10.

Geometry of OLS

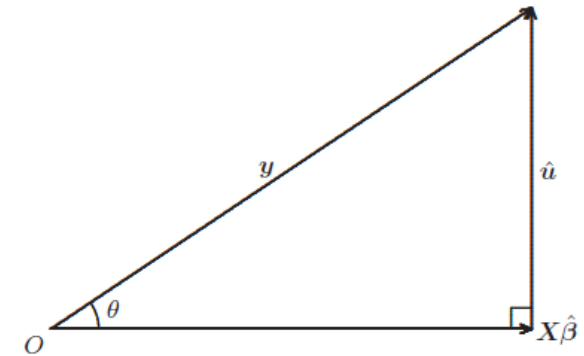


Figure 10 : Residuals and fitted values

Residual vector is orthogonal to all of the regressors:

$$\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\hat{\beta}) = 0$$

Geometry of OLS

- ▶ Adding a third dimension we obtain Figure 11.
- ▶ There are two regressors x_1 and x_2 which together span the horizontal plane labeled $\mathcal{S}(x_1, x_2)$. The shortest distance from \mathbf{y} to the horizontal plane is obtained by dropping a perpendicular. Minimizing SSR accomplishes this.
- ▶ Using Pythagoras' theorem we can write

$$\|\mathbf{y}\|^2 = \|\mathbf{X}\hat{\beta}\|^2 + \|\hat{\mathbf{u}}\|^2$$

or

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

i.e.

$$TSS = ESS + SSR$$

TSS = Total Sum of Squares; ESS = Explained Sum of Square; SSR = Residual Sum of Squares

Geometry of OLS

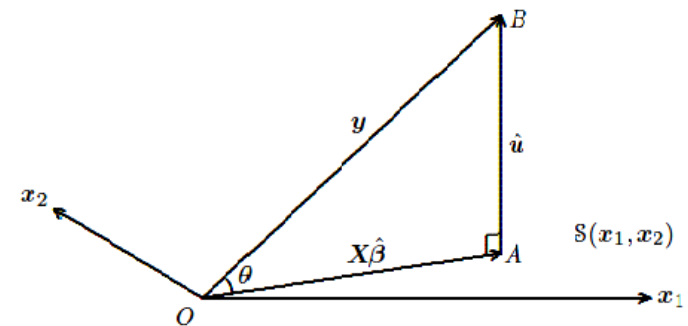


Figure 11 : \mathbf{y} projected on two regressors

Geometry of OLS

- ▶ When we estimate a linear regression model we implicitly map the regressand \mathbf{y} into a vector of fitted values $\mathbf{X}\hat{\beta}$ and a vector of residuals $\hat{\mathbf{u}} = \mathbf{y} - \mathbf{X}\hat{\beta}$.
- ▶ These mappings are examples of **orthogonal projections**.
- ▶ A **projection** is a mapping that takes each point \mathbb{E}^n into a point in a subspace of \mathbb{E}^n while leaving all points in that subspace unchanged.
- ▶ An **orthogonal projection** maps any point into the point of the subspace that is closest to it.

Geometry of OLS

- ▶ An orthogonal projection onto a given subspace can be performed by premultiplying the vector to be projected by a suitable **projection matrix**.
- ▶ In OLS

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

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

- ▶ Obviously,

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

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

\mathbf{P} projects onto $\mathcal{S}(\mathbf{X})$ and when applied to \mathbf{y} , yields the fitted values.

- ▶ Similarly, when \mathbf{y} is premultiplied by \mathbf{M} it yields residuals:

$$\mathbf{M}\mathbf{y} = \left(\mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \right) \mathbf{y} = \mathbf{y} - \mathbf{P}\mathbf{y} = \mathbf{y} - \mathbf{X}\hat{\beta} = \hat{\mathbf{u}}$$

Geometry of OLS

- ▶ Note that

$$\mathbf{P}\mathbf{X} = \mathbf{X}, \quad \mathbf{M}\mathbf{X} = \mathbf{0}, \quad \mathbf{P}\mathbf{P} = \mathbf{P}, \quad \mathbf{M}\mathbf{M} = \mathbf{M},$$

$$\mathbf{P} + \mathbf{M} = \mathbf{I}, \quad \mathbf{P}\mathbf{M} = \mathbf{0}$$

- ▶ Both \mathbf{P} and \mathbf{M} are symmetric and idempotent matrices. Orthogonal decomposition can be written as

$$\mathbf{y} = \mathbf{P}\mathbf{y} + \mathbf{M}\mathbf{y}$$

which can be represented by a right-angled triangle where \mathbf{y} is the hypotenuse and $\mathbf{P}\mathbf{y}$ and $\mathbf{M}\mathbf{y}$ are the other two sides.

- ▶ Using Pythagoras' theorem

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

i.e.

$$TSS = ESS + SSR$$

Coefficient of Determination, R^2

- ▶ The decomposition $TSS = ESS + SSR$ can be written as

$$\|\mathbf{y}\|^2 = \|\mathbf{P}\mathbf{y}\|^2 + \|\mathbf{M}\mathbf{y}\|^2.$$

- ▶ Using this we can define uncentered R^2

$$R_u^2 = \frac{ESS}{TSS} = \frac{\|\mathbf{P}\mathbf{y}\|^2}{\|\mathbf{y}\|^2} = 1 - \frac{\|\mathbf{M}\mathbf{y}\|^2}{\|\mathbf{y}\|^2} = 1 - \frac{SSR}{TSS} = \cos^2 \theta$$

where θ is the angle between \mathbf{y} and $\mathbf{P}\mathbf{y}$.

- ▶ However, R_u^2 is not invariant to the changes in units of measurement. Instead centered R^2 is more widely used:

$$R^2 = \frac{\|\mathbf{P}\mathbf{M}_i \mathbf{y}\|^2}{\|\mathbf{M}_i \mathbf{y}\|^2}$$

where

$$\mathbf{M}_i = \mathbf{I} - \mathbf{P}_i = \mathbf{I} - \mathbf{z}(\mathbf{z}^\top \mathbf{z})^{-1} \mathbf{z}^\top$$

where \mathbf{z} is n -vector of ones. Premultiplying \mathbf{y} by \mathbf{M}_i we obtain deviations from arithmetic mean. Since $-1 \leq \cos \theta \leq 1$ we have $0 \leq R^2 \leq 1$. Regression equation must contain a constant term.

Decomposition of y

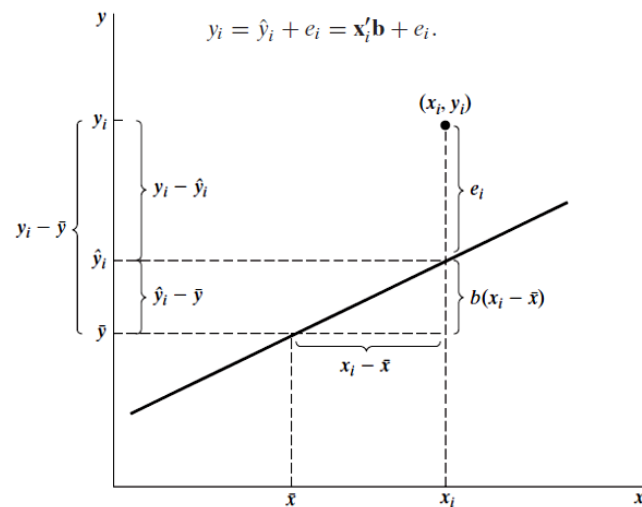


Figure 12 : Decomposition of y_i . Note that e_i is the i th residual (see Greene, p.40)

Goodness-of-Fit: R -Squared

- ▶ The coefficient of determination, R^2 , is simply an estimate of “how much variation in y is explained by x_1, x_2, \dots, x_k in the population”.
- ▶ A low R^2 value does not automatically imply that the classical assumptions fail.
- ▶ As the number of explanatory variables (k) increases, R^2 always increases (it never decreases). Thus, R^2 has a limited role in choosing between alternative models.
- ▶ The relative change in the R -squared when variables added to an equation may be very helpful (e.g. F-statistic for exclusion restrictions depends on the difference in R^2 s).

Adjusted R -Squared: \bar{R}^2

- ▶ Recall the definition of R^2 :

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

- ▶ Dividing the numerator and denominator by n :

$$R^2 = 1 - \frac{SSR/n}{TSS/n} = 1 - \frac{\sigma_u^2}{\sigma_y^2}$$

- ▶ Since TSS/n and SSR/n are biased estimators of respective population variances we will instead use:

$$\frac{TSS}{n-1}, \quad \frac{SSR}{n-k}$$

Adjusted R -Squared: \bar{R}^2

- ▶ Adjusted R -squared is defined as

$$\bar{R}^2 = 1 - \frac{SSR/(n-k)}{TSS/(n-1)} = 1 - (1-R^2) \frac{n-1}{n-k} = 1 - \frac{(n-1)\mathbf{y}^\top \mathbf{M}_X \mathbf{y}}{(n-k)\mathbf{y}^\top \mathbf{M}_I \mathbf{y}}$$

- ▶ Adjusted R^2 , or R -bar squared may increase or decrease when a new variable is added to the regression. Recall that, in contrast, R^2 never decreases.
- ▶ The reason is that when a new variable is added, while SSR decreases, the degrees of freedom ($n-k$) also decreases.
- ▶ Basically, it imposes a penalty for adding additional variables to a model. $SSR/(n-k)$ can go up or down.
- ▶ When a new x variable is added, R -bar square increases if, and only if, the t statistic on the new variable is greater than one in absolute value.
- ▶ Extension: when a group of x variables is added, \bar{R}^2 increases if, and only if, the F statistic for joint significance of the new variables is greater than 1.

Frisch-Waugh-Lovell (FWL) Theorem

- ▶ Rewrite the model by grouping regressors:

$$y = X_1\beta_1 + X_2\beta_2 + u, \quad X = \left[\underbrace{X_1}_{n \times k_1} \mid \underbrace{X_2}_{n \times k_2} \right]$$

- ▶ Define

$$M_1 = I - X_1(X_1^\top X_1)^{-1}X_1^\top$$

- ▶ $M_1 y$ gives the residuals from the regression of y on X_1
- ▶ $M_1 X_2$ gives the residuals from the regression of X_2 on X_1
- ▶ FWL theorem states that $\hat{\beta}_2$ can be obtained by regressing $M_1 y$ on $M_1 X_2$

$$M_1 y = M_1 X_2 \beta_2 + u_2$$

- ▶ The effect of X_1 is partialled out

Linear Transformations of Regressors

- ▶ What happens to OLS estimates, fitted values and residuals when take a linear transformation of explanatory variables?
- ▶ Let A be any nonsingular $k \times k$ matrix of constants and

$$\tilde{X} = XA$$

$$y = \tilde{X}\beta + u$$

- ▶ It can be shown that

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

- ▶ Thus the fitted values and residuals will remain the same, whereas OLS estimates will change.
- ▶ Residuals and fitted values are invariant under nonsingular transformations of the explanatory variables (eg consider changing units of measurement).

Influential Observations

- ▶ OLS parameter estimators are just a weighted average of the elements of the vector y .
- ▶ Defining the i th row $(X^\top X)^{-1} X^\top$ as c_i we see that

$$\hat{\beta}_i = c_i y$$

- ▶ As may be obvious from this relationship some of the observations may have much more influence than the other observations. This may be seen in Figure 11.

Influential Observations

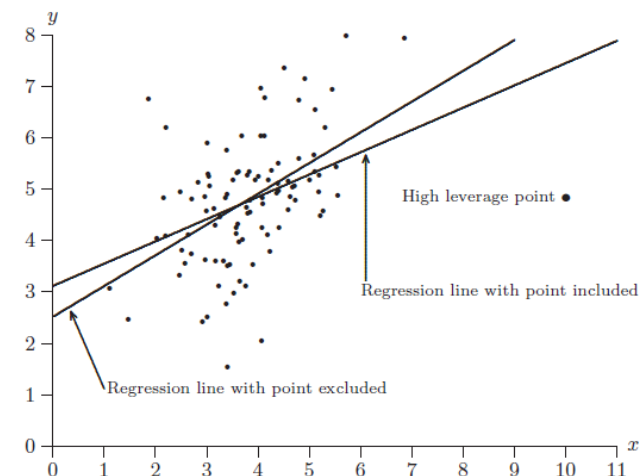


Figure 13 : An influential observation

Influential Observations

- ▶ Influential observations can be diagnosed by inspection the diagonal elements of the projection matrix \mathbf{P} , which is sometimes called the **hat matrix**. These observations are called high leverage or leverage points.
- ▶ The i th diagonal element is usually denoted h_i . Sum of these elements is just the trace of \mathbf{P} :

$$\begin{aligned}\sum_{i=1}^n h_i &= \text{tr}(\mathbf{P}) = \text{tr}(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) = \text{tr}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}) \\ &= \text{tr}(\mathbf{I}_k) = k\end{aligned}$$

where we used the property of trace:

$$\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{CAB}) = \text{tr}(\mathbf{BCA}).$$

- ▶ Average of h_i is just k/n . If the elements of h_i are close to their average value then no observations has very much leverage (balanced design). If some of h_i are much larger than k/n then they may be influential observations (unbalanced design).