

3. OLS Part III

In this section we derive some finite-sample properties of the OLS estimator.

3.1 The Sampling Distribution of the OLS Estimator

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

$$\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = f(\mathbf{y})$$

$\boldsymbol{\varepsilon}$ is random \Rightarrow \mathbf{y} is random \Rightarrow \mathbf{b} is random

- \mathbf{b} is an *estimator* of $\boldsymbol{\beta}$. It is a function of the *random* sample data.
- \mathbf{b} is a “statistic”.
- \mathbf{b} has a probability distribution – called its *Sampling Distribution*.
- Interpretation of *sampling distribution* –

Repeatedly draw all possible samples of size n .

Calculate values of \mathbf{b} each time.

Construct relative frequency distribution for the \mathbf{b} values and probability of occurrence.

It is a *hypothetical* construct. Why?

- Sampling distribution offers *one* basis for answering the question:

“How good is \mathbf{b} as an estimator of $\boldsymbol{\beta}$?”

Note:

Quality of estimator is being assessed in terms of performance in *repeated samples*. Tells us nothing about quality of estimator for *one particular sample*.

- Let's explore some of the properties of the LS estimator, \mathbf{b} , and build up its sampling distribution.
- Introduce some general results, and apply them to our problem.

Definition: An estimator, $\hat{\boldsymbol{\theta}}$ is an *unbiased* estimator of the parameter vector, $\boldsymbol{\theta}$, if $E[\hat{\boldsymbol{\theta}}] = \boldsymbol{\theta}$.

That is, $E[\hat{\boldsymbol{\theta}}(\mathbf{y})] = \boldsymbol{\theta}$.

That is, $\int \hat{\boldsymbol{\theta}}(\mathbf{y})p(\mathbf{y} | \boldsymbol{\theta})d\mathbf{y} = \boldsymbol{\theta}$.

The quantity, $\mathbf{B}(\boldsymbol{\theta}, \mathbf{y}) = E[\hat{\boldsymbol{\theta}}(\mathbf{y}) - \boldsymbol{\theta}]$, is called the “Bias” of $\hat{\boldsymbol{\theta}}$.

Example: $\{y_1, y_2, \dots, y_n\}$ is a random sample from population with a finite mean, μ , and a finite variance, σ^2 .

Consider the *statistic* $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$.

$$\begin{aligned} \text{Then, } E[\bar{y}] &= E\left[\frac{1}{n} \sum_{i=1}^n y_i\right] = \frac{1}{n} \sum_{i=1}^n E(y_i) \\ &= \frac{1}{n} \sum_{i=1}^n \mu = \left(\frac{1}{n} n\mu\right) = \mu. \end{aligned}$$

So, \bar{y} is an *unbiased estimator* of the parameter, μ .

- Here, there are lots of possible unbiased estimators of μ .
- So, need to consider additional characteristics of estimators to help choose.

Return to our LS problem –

$$\mathbf{b} = (X'X)^{-1}X'\mathbf{y}$$

- Recall – either assume that X is *non-random*, or condition on X .
- We'll assume X is non-random – get same result if we condition on X .

Then: $E(\mathbf{b}) = E[(X'X)^{-1}X'\mathbf{y}] = (X'X)^{-1}X'E(\mathbf{y})$

So,

$$\begin{aligned} E(\mathbf{b}) &= (X'X)^{-1}X'E[X\boldsymbol{\beta} + \boldsymbol{\varepsilon}] = (X'X)^{-1}X'[X\boldsymbol{\beta} + E(\boldsymbol{\varepsilon})] \\ &= (X'X)^{-1}X'[X\boldsymbol{\beta} + \mathbf{0}] = (X'X)^{-1}X'X\boldsymbol{\beta} \\ &= \boldsymbol{\beta}. \end{aligned}$$

The LS estimator of $\boldsymbol{\beta}$ is Unbiased

Definition: Any estimator that is a *linear function* of the random sample data is called a *Linear Estimator*.

Example: $\{y_1, y_2, \dots, y_n\}$ is a random sample from population with a finite mean, μ , and a finite variance, σ^2 .

Consider the *statistic* $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} [y_1 + y_2 + \dots + y_n]$.

This statistic is a *linear estimator* of μ .

(Note that the “weights” are non-random.)

Return to our LS problem –

$$\mathbf{b} = (X'X)^{-1}X'\mathbf{y} = A\mathbf{y}$$

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

Note that, under our assumptions, A is a *non-random* matrix.

So,

$$\begin{pmatrix} b_1 \\ \vdots \\ b_k \end{pmatrix} = \begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \cdots & \vdots \\ a_{k1} & \cdots & a_{kn} \end{bmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}.$$

For example, $b_1 = [a_{11}y_1 + a_{12}y_2 + \dots + a_{1n}y_n]$; *etc.*

The LS estimator, \mathbf{b} , is a linear (& unbiased) estimator of $\boldsymbol{\beta}$

Now let's consider the dispersion (variability) of \mathbf{b} , as an estimator of $\boldsymbol{\beta}$.

Definition: Suppose we have an $(n \times 1)$ random vector, \mathbf{x} . Then the *Covariance Matrix* of \mathbf{x} is defined as the $(n \times n)$ matrix:

$$V(\mathbf{x}) = E[(\mathbf{x} - E(\mathbf{x}))(\mathbf{x} - E(\mathbf{x}))'].$$

- Diagonal elements of $V(\mathbf{x})$ are $var.(x_1), \dots, var.(x_n)$.
- Off-diagonal elements are $covar.(x_i, x_j); i, j = 1, \dots, n; i \neq j$.

Return to our LS problem –

We have a $(k \times 1)$ random vector, \mathbf{b} , and we know that $E(\mathbf{b}) = \boldsymbol{\beta}$.

$$V(\mathbf{b}) = E[(\mathbf{b} - E(\mathbf{b}))(\mathbf{b} - E(\mathbf{b}))']$$

Now,

$$\begin{aligned} \mathbf{b} &= (X'X)^{-1}X'\mathbf{y} = (X'X)^{-1}X'(X\boldsymbol{\beta} + \boldsymbol{\varepsilon}) \\ &= (X'X)^{-1}(X'X)\boldsymbol{\beta} + (X'X)^{-1}X'\boldsymbol{\varepsilon} \\ &= I\boldsymbol{\beta} + (X'X)^{-1}X'\boldsymbol{\varepsilon}. \end{aligned}$$

So,

$$(\mathbf{b} - \boldsymbol{\beta}) = (X'X)^{-1}X'\boldsymbol{\varepsilon}. \quad [*]$$

Using the result, $[*]$, in $V(\mathbf{b})$, we have:

$$\begin{aligned} V(\mathbf{b}) &= E\{[(X'X)^{-1}X'\boldsymbol{\varepsilon}][(X'X)^{-1}X'\boldsymbol{\varepsilon}']\} \\ &= (X'X)^{-1}X'E[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}']X(X'X)^{-1}. \end{aligned}$$

We showed, earlier, that because $E(\boldsymbol{\varepsilon}) = \mathbf{0}$, $V(\boldsymbol{\varepsilon}) = E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}') = \sigma^2 I_n$.

(What other assumptions did we use to get this result?)

So, we have:

$$V(\mathbf{b}) = (X'X)^{-1}X'E[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}']X(X'X)^{-1} = (X'X)^{-1}X'\sigma^2IX(X'X)^{-1} = \sigma^2(X'X)^{-1}(X'X)(X'X)^{-1} \\ = \sigma^2(X'X)^{-1}.$$

$$V(\mathbf{b}) = \sigma^2(X'X)^{-1} \\ (k \times k)$$

Interpret diagonal and off-diagonal elements of this matrix.

Finally, because the error term, $\boldsymbol{\varepsilon}$ is assumed to be Normally distributed,

1. $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$: this implies that \mathbf{y} is also Normally distributed. (Why?)
2. $\mathbf{b} = (X'X)^{-1}X'\mathbf{y} = A\mathbf{y}$: this implies that \mathbf{b} is also Normally distributed.

So, we now have the full **Sampling Distribution** of the LS estimator, \mathbf{b} :

$$\mathbf{b} \sim N[\boldsymbol{\beta}, \sigma^2(X'X)^{-1}]$$

Note:

- This result depends on our various, *rigid*, assumptions about the various components of the regression model.
- The Normal distribution here is a “*multivariate* Normal” distribution. (*See handout on “Spherical Distributions”.*)
- As with estimation of population mean, $\boldsymbol{\mu}$, in previous example, there are lots of other *unbiased* estimators of $\boldsymbol{\beta}$ in the model $= X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$.
- How might we choose between these possibilities? Is *linearity* desirable?
- We need to consider other *desirable* properties that these unbiased estimators may have.
- *One option* is to take account of estimators' *precisions*.

3.2 The Efficiency of OLS

Definition: Suppose we have two *unbiased* estimators, $\hat{\theta}_1$ and $\hat{\theta}_2$, of the (scalar) parameter, θ . Then we say that $\hat{\theta}_1$ is **at least as efficient** as $\hat{\theta}_2$ if $\text{var.}(\hat{\theta}_1) \leq \text{var.}(\hat{\theta}_2)$.

Note:

1. The variance of an estimator is just the variance of its sampling distribution.
2. "Efficiency" is a *relative* concept.
3. What if there are 3 or more unbiased estimators being compared?
 - What if one or more of the estimators being compared is *biased*?
 - In this case we can take account of both variance, and any bias, at the same time by using "*mean squared error*" (MSE) of the estimators.

Definition: Suppose we have two *unbiased* estimators, $\hat{\theta}_1$ and $\hat{\theta}_2$, of the parameter vector, θ . Then we say that $\hat{\theta}_1$ is **at least as efficient** as $\hat{\theta}_2$ if $\Delta = V(\hat{\theta}_2) - V(\hat{\theta}_1)$ is *at least positive semi-definite*.

Taking account of its *linearity*, *unbiasedness*, and its *precision*, in what sense is the LS estimator, b , of β *optimal*?

Theorem (Gauss-Markhov):

In the "standard" linear regression model, $y = X\beta + \varepsilon$, the LS estimator, b , of β is **Best Linear Unbiased** (BLU). That is, it is **Efficient** in the class of all linear and unbiased estimators of β .

1. Is this an *interesting* result?
2. What *assumptions* about the "standard" model are we going to exploit?

Proof

Let \mathbf{b}_0 be any other *linear* estimator of $\boldsymbol{\beta}$:

$$\mathbf{b}_0 = C\mathbf{y} \quad ; \quad \text{for some non-random } C .$$

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

Now, $V(\mathbf{b}_0) = CV(\mathbf{y})C' = C(\sigma^2 I_n)C' = \sigma^2 CC'$

$$(k \times k)$$

Define: $D = C - (X'X)^{-1}X'$

so that $D\mathbf{y} = C\mathbf{y} - (X'X)^{-1}X'\mathbf{y} = \mathbf{b}_0 - \mathbf{b}$.

Now restrict \mathbf{b}_0 to be *unbiased*, so that $E(\mathbf{b}_0) = E(C\mathbf{y}) = CX\boldsymbol{\beta} = \boldsymbol{\beta}$.

This requires that $CX = I$, which in turn implies that

$$DX = [C - (X'X)^{-1}X']X = CX - I = \mathbf{0} \quad (\text{and } D'X' = 0)$$

(What assumptions have we used so far?)

Now, focus on covariance matrix of \mathbf{b}_0 :

$$\begin{aligned} V(\mathbf{b}_0) &= \sigma^2 [D + (X'X)^{-1}X'] [D + (X'X)^{-1}X']' \\ &= \sigma^2 [DD' + (X'X)^{-1}X'X(X'X)^{-1}] \quad ; \quad DX = 0 \\ &= \sigma^2 DD' + \sigma^2 (X'X)^{-1} \\ &= \sigma^2 DD' + V(\mathbf{b}), \end{aligned}$$

or, $[V(\mathbf{b}_0) - V(\mathbf{b})] = \sigma^2 DD' \quad ; \quad \sigma^2 > 0$

Now we just have to "sign" this (matrix) difference:

$$\boldsymbol{\eta}'(DD')\boldsymbol{\eta} = (D'\boldsymbol{\eta})'(D'\boldsymbol{\eta}) = \mathbf{v}'\mathbf{v} = \sum_{i=1}^n v_i^2 \geq 0 .$$

So, $\Delta = [V(\mathbf{b}_0) - V(\mathbf{b})]$ is a p.s.d. matrix, implying that \mathbf{b}_0 is *relatively less efficient* than \mathbf{b} .

Result:

The LS estimator is the Best Linear Unbiased estimator of β .

- What assumptions did we use, and where?
- Were there any standard assumptions that we *didn't* use?
- What does this suggest?