

Econometrics I - Finite sample properties of the least squares estimator

Ryan T. Godwin

University of Manitoba

Finite sample properties of the least squares estimator

We'll derive some of the finite sample properties for the LS estimator. Finite sample properties \rightarrow statistical properties of \mathbf{b} for some finite value of n . Later, we'll consider properties of \mathbf{b} for when $n \rightarrow \infty$ (i.e. asymptotic properties). Recall that the population model is:

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

and that by assumptions A.3, A.4 and A.6, the error term is:

$$\boldsymbol{\epsilon} \sim N[\mathbf{0}, \sigma^2 I_n]$$

and also recall that our LS estimator is:

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

That is, \mathbf{b} is a function of the random sample data, so is itself random! This is a very important point.

ϵ is random \longrightarrow \mathbf{y} is random \longrightarrow \mathbf{b} is random

- ▶ \mathbf{b} is an estimator of β . 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 the sampling distribution:

- ▶ Repeatedly draw all possible samples of size n .
- ▶ Calculate values of \mathbf{b} each time.
- ▶ Construct a relative frequency distribution for the \mathbf{b} values and probability of occurrence.
- ▶ It is a hypothetical construct.

Question: Why is the sampling distribution a hypothetical construct?

The *sampling distribution* offers one basis for answering the question: “How good is \mathbf{b} as an estimator of β ?”

We will be assessing the quality of the estimator in terms of its performance in *repeated samples*. The finite sample properties that we derive in this chapter tell us nothing about the quality of the estimator for *one particular sample*.

We will explore some of the properties of the LS estimator, \mathbf{b} , and build up its sampling distribution. We'll introduce some general results, and then apply them to our specific problem. Before we do so, let's look at a *simulation experiment* using R.

For the experiment, I will estimate the simple model:

$$\mathbf{y} = \beta_1 + \beta_2 \mathbf{x} + \epsilon$$

where I will maintain all of the usual assumptions (A.1-A.6) and choose:

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = \begin{bmatrix} 2 \\ -4 \end{bmatrix} \quad ; \quad \sigma^2 = 1 \quad ; \quad n = 100$$

In addition I will use the unrealistic assumption that x is “fixed in repeated samples”, and will generate the x variable from a $N(0, 1)$ distribution (I will use the same distribution for ϵ as well).

In R, run the following code to set the parameters of the experiment:

```
1 beta1 <- 2
2 beta2 <- -4
3 n <- 100
4 x <- rnorm(n)
```

and now run the following code several times:

```
1 epsilon <- rnorm(n)
2 y <- beta1 + beta2 * x + epsilon
3 lm(y ~ x)
```

```
1 Coefficients:
2 (Intercept)          x
3      1.886      -4.060
```

Question: How can you use the above code to *simulate* the sampling distribution of b_2 ?

To simulate the sampling distribution, we need to run the above code many times, and collect the value of b_2 each time.

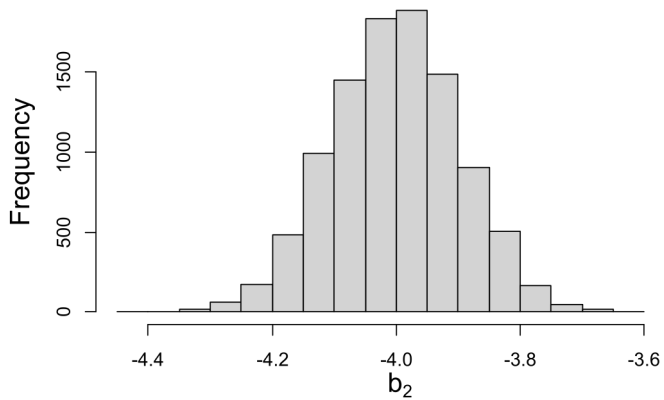
Then, we can plot b_2 in a histogram, take mean of b_2 , take the $\text{var}[b_2]$, etc. I set the *random seed* so that when you run the code on your computer, you get the same results as me. I use a for loop to repeat the experiment $\text{nrep} = 10000$ times.

```
1 set.seed(7010)
2 nrep <- 10000
3 n <- 100
4 x <- rnorm(n)
5 beta1 <- 2
6 beta2 <- -4
7 b2 <- numeric(nrep)
8 for(i in 1:nrep){
9   epsilon <- rnorm(n, mean=0, sd=1)
10  y <- beta1 + beta2 * x + epsilon
11  b2[i] <- lm(y ~ x)$coefficients[2]
12 }
13 mean(b2)
14 var(b2)
15 hist(b2)
```

```
1 > mean(b2)
2 [1] -4.000452
3 > var(b2)
4 [1] 0.01061148
```

The histogram of all 10,000 b_2 values is shown in Figure 1.

Figure: A simulated sampling distribution.



Questions:

1. What distribution appears to describe Figure 1?
2. Is the average value of the *estimator* close to the true *population* value?

Unbiased

Unbiased estimator. An estimator, $\hat{\theta}$, is an *unbiased* estimator of the parameter vector, θ , if:

$$E[\hat{\theta}] = \theta$$

That is, if $E[\hat{\theta}(\mathbf{y})] = \theta$, or $\int \hat{\theta}(\mathbf{y})p(\mathbf{y}|\theta)d\mathbf{y} = \theta$. The “Bias” of $\hat{\theta}$ is the quantity:

$$\text{Bias}(\theta, \mathbf{y}) = E[\hat{\theta}(\mathbf{y}) - \theta]$$

In words, an estimator is *unbiased* if it gives the right answer on average.

Unbiasedness of \bar{y} . Let $\{y_1, y_2, \dots, y_n\}$ be a random sample from population with a finite mean, μ , and a finite variance, σ^2 . Consider the *statistic*:

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

Then,

$$E[\bar{y}] = E\left[\frac{1}{n} \sum_{i=1}^n y_i\right] = \frac{1}{n} \sum_{i=1}^n E(y_i) = \frac{1}{n} \sum_{i=1}^n \mu = \left(\frac{1}{n} n \mu\right) = \mu$$

So, \bar{y} is an *unbiased estimator* of the parameter, μ . Here, there are lots of possible unbiased estimators of μ . So, we will need to consider additional characteristics of estimators to help us choose from among them.

Return to our LS problem:

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

Recall assumption A.5. We will use the strongest version of this assumption - that X is non-random. We could use the weaker versions of the assumption, and we'll get the same results, but the notation will be more cumbersome. Now, take the expected value of the random estimator \mathbf{b} :

$$\begin{aligned} E(\mathbf{b}) &= E \left[(X'X)^{-1} X'\mathbf{y} \right] = (X'X)^{-1} X'E(\mathbf{y}) \\ &= (X'X)^{-1} X'E[X\boldsymbol{\beta} + \boldsymbol{\epsilon}] = (X'X)^{-1} X'[X\boldsymbol{\beta} + E(\boldsymbol{\epsilon})] \\ &= (X'X)^{-1} X'[X\boldsymbol{\beta} + \mathbf{0}] = (X'X)^{-1} X'X\boldsymbol{\beta} \\ &= \boldsymbol{\beta} \end{aligned} \quad (1)$$

This proves that the **LS estimator of $\boldsymbol{\beta}$ is Unbiased.**

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

Sample average is linear. Let $\{y_1, y_2, \dots, y_n\}$ be 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:

$$\begin{array}{r} \mathbf{b} \\ (k \times 1) \end{array} = (X'X)^{-1} X'\mathbf{y} = \begin{array}{r} A\mathbf{y} \\ (k \times n)(n \times 1) \end{array}$$

Note that, under our (strictest form of the) assumptions, A is a *non-random* matrix. So,

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

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

Thus: **the LS estimator, \mathbf{b} , is a linear (and unbiased) estimator of β .**

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

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 $\text{var}(x_1), \dots, \text{var}(x_n)$.
- ▶ Off-diagonal elements are $\text{cov}(x_i, x_j)$; $i, j = 1, \dots, n$; $i \neq j$.
- ▶ We have already made use of the *covariance matrix* when we made assumption A.4.

Return to our LS problem. We have a $(k \times 1)$ random vector, \mathbf{b} , and we know that $E(\mathbf{b}) = \boldsymbol{\beta}$. The *covariance matrix* of \mathbf{b} , $V(\mathbf{b})$, is:

$$\begin{aligned} V(\mathbf{b}) &= E [(\mathbf{b} - E(\mathbf{b}))(\mathbf{b} - E(\mathbf{b}))'] \\ &= E [(\mathbf{b} - \boldsymbol{\beta})(\mathbf{b} - \boldsymbol{\beta})'] \end{aligned}$$

Now,

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

So,

$$(\mathbf{b} - \boldsymbol{\beta}) = (X'X)^{-1} X'\boldsymbol{\epsilon} \quad (2)$$

Using the result from 2 in $V(\mathbf{b})$ we have:

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

For assumption A.4, we showed earlier that because $E(\boldsymbol{\epsilon}) = \mathbf{0}$:

$$V(\boldsymbol{\epsilon}) = E(\boldsymbol{\epsilon}\boldsymbol{\epsilon}') = \sigma^2 I_n$$

So, we have:

$$\begin{aligned} V(\mathbf{b}) &= (X'X)^{-1} X' E[\boldsymbol{\epsilon}\boldsymbol{\epsilon}'] X (X'X)^{-1} \\ &= (X'X)^{-1} X' \sigma^2 I X (X'X)^{-1} \\ &= \sigma^2 (X'X)^{-1} (X'X) (X'X)^{-1} \\ &= \sigma^2 (X'X)^{-1}. \end{aligned}$$

So, the covariance matrix of \mathbf{b} is:

$$V(\mathbf{b}) = \sigma^2 (X'X)^{-1} \tag{3}$$

Question: What is the interpretation of the diagonal and off-diagonal elements of this matrix? What might the elements of this matrix be used for, in practice?

Finally, because the error term, ϵ is assumed to be Normally distributed,

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

Question: Why does the Normality of ϵ transfer to \mathbf{b} ?

We now have the full *sampling distribution* of the LS estimator, \mathbf{b} :

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

Note:

- ▶ This result depends on our *rigid* assumptions about the various components of the regression model.
- ▶ The Normal distribution here is a “multivariate Normal” distribution.
- ▶ As with estimation of the population mean μ , there are lots of other unbiased estimators of $\boldsymbol{\beta}$.

Question: How might we choose between the many possible linear and unbiased estimators for β ? Why is linearity desirable?

We need to consider other desirable properties that these unbiased estimators may have. One option to help discern the “best” estimator for β is to take into account the estimators’ *precisions*.

Suppose that we have two different *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 $\text{var}(\hat{\theta}_1) \leq \text{var}(\hat{\theta}_2)$

- ▶ The variance of an estimator is just the variance of its sampling distribution.
- ▶ “Efficiency” is a relative concept.

Question: What if there are 3 or more unbiased estimators being compared?

Mean Squared Error (MSE)

If one or more of the estimators being compared is biased we can take account of both variance, and any bias, at the same time by using “mean squared error” (MSE) of the estimators.

Suppose that $\hat{\theta}$ is an estimator of the (scalar) parameter, θ . Then the MSE of $\hat{\theta}$ is defined as:

$$\text{MSE}(\hat{\theta}) = E \left[(\hat{\theta} - \theta)^2 \right]$$

Note that:

$$\text{MSE}(\hat{\theta}) = \text{var}(\hat{\theta}) + [\text{Bias}(\hat{\theta})]^2$$

Efficiency and MSE for a vector of estimators

If we extend all of this to the case where we have a vector of parameters, $\boldsymbol{\theta}$, then we have the following definitions:

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

Suppose that we have two different (possibly) *biased* estimators, $\hat{\boldsymbol{\theta}}_1$ and $\hat{\boldsymbol{\theta}}_2$, of the parameter vector, $\boldsymbol{\theta}$. Then we say that $\hat{\boldsymbol{\theta}}_1$ is **at least as efficient** as $\hat{\boldsymbol{\theta}}_2$ if $\Delta = \text{MMSE}(\hat{\boldsymbol{\theta}}_2) - \text{MMSE}(\hat{\boldsymbol{\theta}}_1)$ is *positive semi-definite*.

Note: $\text{MMSE}(\hat{\boldsymbol{\theta}}) = E[(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})'] = V[\hat{\boldsymbol{\theta}}] + \text{Bias}(\hat{\boldsymbol{\theta}})\text{Bias}(\hat{\boldsymbol{\theta}})'$.

Gauss-Markhov theorem

Taking account of its *linearity*, *unbiasedness*, and its *precision*, in what sense is the LS estimator, \mathbf{b} , of $\boldsymbol{\beta}$, optimal?

Gauss-Markhov Theorem: In the “standard” linear regression model, $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$, the LS estimator, \mathbf{b} , of $\boldsymbol{\beta}$ is **Best Linear Unbiased (BLU)**. That is, it is *efficient* in the class of all linear and unbiased estimators of $\boldsymbol{\beta}$.

Question: Why is this an interesting result?

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

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

The *covariance matrix* for \mathbf{b}_0 is:

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

Take the difference between \mathbf{b}_0 and \mathbf{b} :

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

where

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

is the difference between how the other estimator uses the X data to “weight” the y data (C), and how the LS estimator uses the X data ($(X'X)^{-1} X'$). 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:

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

Solve for C in terms of D :

$$C = D + (X'X)^{-1} X',$$

and return to the covariance matrix of \mathbf{b}_0 :

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

or:

$$[V(\mathbf{b}_0) - V(\mathbf{b})] = \sigma^2 DD' \quad ; \quad \sigma^2 > 0 \quad (4)$$

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

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

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 (BLUE) of $\boldsymbol{\beta}$.

Questions:

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