

The sampling distribution of the OLS estimator,
and the statistical properties of the OLS
estimator in finite and large samples

Econometrics (35B206), Lecture 2

Tunga Kantarcı, TiSEM, Tilburg University, Spring 2019

Sampling distribution

In this lecture $\hat{\beta}_{OLS} \equiv \hat{\beta}$.

Sampling distribution

Consider the model

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

Suppose we have n observations for \mathbf{y} and \mathbf{X} . Using this **one sample** we estimate $\boldsymbol{\beta}$. Let $\hat{\boldsymbol{\beta}}$ represent this estimate. The estimation results in **one $\hat{\boldsymbol{\beta}}$** .

Sampling distribution

```
. regress wage educ
```

Source	SS	df	MS	Number of obs	=	997
Model	7842.35455	1	7842.35455	F(1, 995)	=	251.46
Residual	31031.0745	995	31.1870095	Prob > F	=	0.0000
				R-squared	=	0.2017
				Adj R-squared	=	0.2009
Total	38873.429	996	39.0295472	Root MSE	=	5.5845

wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
educ	1.135645	.0716154	15.86	0.000	.9951106	1.27618
_cons	-4.860424	.9679821	-5.02	0.000	-6.759944	-2.960903

Sampling distribution

If there is only one $\hat{\beta}$, how can $\hat{\beta}$ has a standard error? How can it have a distribution? How can we talk about the statistical properties of $\hat{\beta}$?

Sampling distribution

The distribution of $\hat{\beta}$ results from a **conceptual experiment**.

Sampling distribution

The experiment is about taking samples from the population repeatedly.

Consider the LRM

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

Take a sample of n observations for \mathbf{y} and \mathbf{X} from the population.

Using the \mathbf{y} and \mathbf{X} , obtain $\hat{\boldsymbol{\beta}}$.

Sampling distribution

Take a new sample of n observations for \mathbf{y} and \mathbf{X} from the population.

Using the new \mathbf{y} and \mathbf{X} , obtain a new $\hat{\beta}$.

Repeatedly take all possible samples from the population.

Obtain many $\hat{\beta}$.

$\hat{\beta}$ now has a distribution. This is the **sampling distribution** of $\hat{\beta}$.

The statistical properties of $\hat{\beta}$ is about this sampling distribution of $\hat{\beta}$.

Sampling distribution

We will study the statistical properties of $\hat{\beta}$ in theory.

But we will also study the statistical properties in a more applied manner. E.g., we will plot the sampling distribution of $\hat{\beta}$, and see how it behaves if we violate an assumption of the SLM.

To conduct this study, we should take samples for \mathbf{y} and \mathbf{X} from the population repeatedly, obtain many $\hat{\beta}$, and plot the sampling distribution of $\hat{\beta}$.

But taking repeated samples from the population is expensive.

What can we do?

Simulating the sampling distribution

We can simulate the sampling distribution of $\hat{\beta}$. The simulation experiment is about obtaining \mathbf{y} in the LRM

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

Take a sample of n observations for \mathbf{y} and \mathbf{X} from the population.

Using the \mathbf{y} and \mathbf{X} , obtain $\hat{\beta}$.

Simulating the sampling distribution

Draw n random numbers for ϵ from a probability distribution.

Given the original observations for \mathbf{X} , the initial estimate of $\hat{\beta}$, and the random numbers for ϵ , generate n observations for \mathbf{y} .

Using the generated \mathbf{y} and original \mathbf{X} , obtain a new $\hat{\beta}$.

Repeat the procedure to generate new sets of \mathbf{y} .

Obtain many $\hat{\beta}$.

Obtain a simulated sampling distribution for $\hat{\beta}$.

Simulating the sampling distribution

In the experiment, we keep the n observations of \mathbf{X} constant as we repeatedly generate new \mathbf{y} . This simplifies the experiment because then we can attribute the response of the sampling distribution to interesting counterfactual scenarios rather than to the sampling variance of \mathbf{X} . This is the same as what we do in statistical derivations. We condition on \mathbf{X} meaning that we keep \mathbf{X} constant. This simplifies the derivations very much. Treating \mathbf{X} constant in repeated sampling is not realistic unless \mathbf{X} is collected in an experimental setting where the experimenter had chosen the value for \mathbf{X} before \mathbf{y} was determined. We justify this treatment with the random sampling assumption.

Simulating the sampling distribution

In the experiment, we assume that $E[\epsilon \mid \mathbf{X}] = 0$ holds.

Simulating the sampling distribution

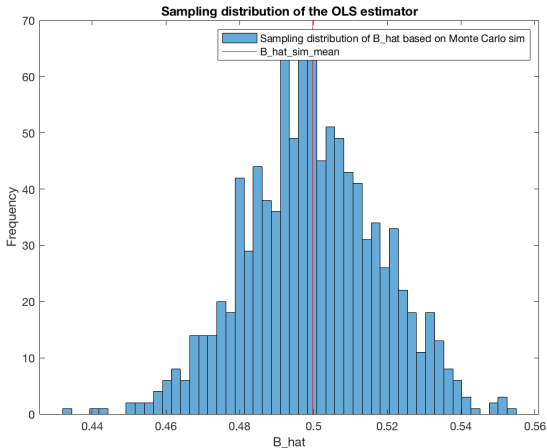
If there is no initial sample data for \mathbf{X} , we can keep the experiment purely hypothetical. Take random draws for \mathbf{X} from a distribution. Do the same for ε . Assume a value for β . Generate \mathbf{y} . Obtain a first $\hat{\beta}$. Generate new ε . Using the original \mathbf{X} and the assumed β , generate new \mathbf{y} . Obtain a new $\hat{\beta}$. Repeat the procedure to obtain many $\hat{\beta}$. This is what we do next.

Simulating the sampling distribution

```
# Simulate the sampling distribution of the OLS estimator
N_sim = 1000
N_obs = 9000
B_true = 0.5
x = unifrnd(-1,1,N_obs,1)
B_hat_sim = NaN(1,N_sim)
for i = 1:N_sim
    e = normrnd(0,1,N_obs,1)
    y = x*B_true+e
    B_hat_sim(1,i) = inv(x'*x)*x'*y
end
```

Simulating the sampling distribution

```
histogram(B_hat_sim)
```



Statistical properties

How do we want the sampling distribution of $\hat{\beta}$ to behave? How do we want the mean and the variance of this distribution to behave?

Statistical properties

We make a distinction between a small and large sample. We want the sampling distribution of $\hat{\beta}$ to behave in certain ways in a finite sample. We want the sampling distribution of $\hat{\beta}$ to behave in certain ways in a large sample.

Statistical properties

Behaviour in a finite sample means that the behaviour does not depend on n . We fix n , and study the behaviour. If we change n , the behaviour is not affected. Behaviour in a large sample means that the behaviour depends on n . We increase n , and study how the sampling distribution behaves. Why we differentiate between a small and large sample will become clear later in these slides.

Statistical properties in finite samples, unbiasedness

Suppose that we calculate a $\hat{\beta}$ with the sample at hand. $\hat{\beta}$ is an estimate of the true β . We want to believe that this particular $\hat{\beta}$ is close to β in some criterion.

Statistical properties in finite samples, unbiasedness

Consider the mean of the sampling distribution of $\hat{\beta}$, conditional on \mathbf{X}

$$E \left[\hat{\beta} \mid \mathbf{X} \right].$$

Note that sampling distribution is created by repeatedly taking all possible samples from the **population**. The mean in the population is the **expected value**!

Statistical properties in finite samples, unbiasedness

The criterion we want $\hat{\beta}$ to satisfy is

$$E \left[\hat{\beta} \mid \mathbf{X} \right] = \beta.$$

It says that the mean of the sampling distribution of $\hat{\beta}$ is equal to β . It says that on average $\hat{\beta}$ will correctly estimate β . If this is true, we say that $\hat{\beta}$ is an **unbiased estimator** of β .

Statistical properties in finite samples, unbiasedness

In the simulation above, `mean(B_hat_sim)` gives 0.4998. `B_true` was 0.5.

Statistical properties in finite samples, unbiasedness

In practice, what does unbiasedness imply? Suppose that you draw an unlucky sample from the population, and obtain a bad $\hat{\beta}$. Or think of our simulation experiment. Suppose that you take n draws for ε from its assumed distribution which turn out to be extreme. $\hat{\beta}$ will be far from its population mean $E[\hat{\beta} | \mathbf{X}]$ which is equal to β . Hence, in practice, to satisfy unbiasedness as much as possible, the sample we draw should be typical.

Statistical properties in finite samples, unbiasedness, proof

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

Taking the expectation conditional on \mathbf{X} ,

$$\begin{aligned}\mathbb{E}[\hat{\beta} \mid \mathbf{X}] &= \mathbb{E}[\beta \mid \mathbf{X}] + \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon \mid \mathbf{X}] \\ &= \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbb{E}[\varepsilon \mid \mathbf{X}] \\ &= \beta\end{aligned}$$

if $\mathbb{E}[\varepsilon \mid \mathbf{X}] = 0$.

Hence, $\mathbb{E}[\varepsilon \mid \mathbf{X}] = 0$ is another requirement for unbiasedness.

Statistical properties in finite samples, unbiasedness

By the LIE, it is also true that

$$\begin{aligned} E \left[\hat{\beta} \right] &= E_{\mathbf{x}} \left[E \left[\hat{\beta} \mid \mathbf{x} \right] \right] \\ &= E_{\mathbf{x}} \left[\beta \right] \\ &= \beta. \end{aligned}$$

Statistical properties in finite samples, unbiasedness

Is the OLS estimator the only unbiased estimator? Consider a competitor estimator

$$\hat{\beta}_0 = \mathbf{C}\mathbf{y}.$$

\mathbf{C} is some $K \times n$ matrix that depends on \mathbf{X} . Taking the expectation conditional on \mathbf{X} ,

$$\begin{aligned} E[\hat{\beta}_0 \mid \mathbf{X}] &= E[\mathbf{C}\mathbf{y} \mid \mathbf{X}] \\ &= E[\mathbf{C}(\mathbf{X}\beta + \varepsilon) \mid \mathbf{X}] \\ &= \mathbf{C}\mathbf{X}\beta + \mathbf{C}E[\varepsilon \mid \mathbf{X}] \\ &= \mathbf{C}\mathbf{X}\beta \\ &= \beta \end{aligned}$$

if $E[\varepsilon \mid \mathbf{X}] = 0$, and if $\mathbf{C}\mathbf{X} = \mathbf{I}$.

The OLS estimator is not the only unbiased estimator!

Statistical properties in finite samples, unbiasedness

The OLS estimator is not the only unbiased estimator. But recall that we wanted to believe in the OLS estimator in some criterion, and we have considered unbiasedness as a criterion. But now we have more than one estimator that is unbiased. Why should we still believe in the OLS estimator $\hat{\beta}$?

Statistical properties in finite samples, unbiasedness

Another objection to the unbiasedness criterion is summarised nicely by the story of three econometricians who go duck hunting. The first shoots about a foot in front of the duck, the second about a foot behind. The third yells: "We got him!"

Statistical properties in finite samples, efficiency

We need to judge $\hat{\beta}$ on an additional criterion than unbiasedness to preserve our belief in $\hat{\beta}$. This new criterion is

$$\text{Var} \left[\hat{\beta}_0 \mid \mathbf{X} \right] \geq \text{Var} \left[\hat{\beta} \mid \mathbf{X} \right].$$

It says that the variance of the sampling distribution of the OLS estimator $\hat{\beta}$ is the smallest when compared to the variance of the sampling distribution of any other competing unbiased estimator $\hat{\beta}_0$.

If this is true, we say that $\hat{\beta}$ is the best unbiased estimator.

Statistical properties in finite samples, efficiency, proof

Derive the variance-covariance matrix of $\hat{\beta}$, conditional on \mathbf{X} .

$$\hat{\beta} = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\epsilon,$$

and hence

$$\hat{\beta} - \beta = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\epsilon.$$

Using the general variance formula, conditional on \mathbf{X} ,

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

if $\text{E} [\epsilon\epsilon' \mid \mathbf{X}] = \sigma^2 \mathbf{I}$.

Statistical properties in finite samples, efficiency, proof

Derive the variance-covariance matrix of $\hat{\beta}_0$, conditional on \mathbf{X} .

$$\hat{\beta}_0 = \mathbf{C}\mathbf{y} = \mathbf{C}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \mathbf{C}\mathbf{X}\boldsymbol{\beta} + \mathbf{C}\boldsymbol{\varepsilon} = \boldsymbol{\beta} + \mathbf{C}\boldsymbol{\varepsilon}$$

where $\mathbf{C}\mathbf{X} = \mathbf{I}$ by unbiasedness of $\hat{\beta}_0$. Hence,

$$\hat{\beta}_0 - \boldsymbol{\beta} = \mathbf{C}\boldsymbol{\varepsilon}.$$

Using the general variance formula, conditional on \mathbf{X} ,

$$\begin{aligned}\text{Var} [\hat{\beta}_0 \mid \mathbf{X}] &= \text{E} \left[\left(\hat{\beta}_0 - \boldsymbol{\beta} \right) \left(\hat{\beta}_0 - \boldsymbol{\beta} \right)' \mid \mathbf{X} \right] \\ &= \text{E} [\mathbf{C}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'\mathbf{C}' \mid \mathbf{X}] \\ &= \mathbf{C}\text{E} [\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}' \mid \mathbf{X}] \mathbf{C}' \\ &= \mathbf{C}\sigma^2\mathbf{I}\mathbf{C}' \\ &= \sigma^2\mathbf{C}\mathbf{C}'\end{aligned}$$

since \mathbf{C} depends on \mathbf{X} , and if $\text{E} [\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}' \mid \mathbf{X}] = \sigma^2\mathbf{I}$.

Statistical properties in finite samples, efficiency, proof

$$\begin{aligned}\text{Var} [\hat{\beta}_0 | \mathbf{X}] - \text{Var} [\hat{\beta} | \mathbf{X}] &= \sigma^2 \mathbf{C} \mathbf{C}' - \sigma^2 (\mathbf{X}' \mathbf{X})^{-1} \\ &= \sigma^2 [\mathbf{C} \mathbf{C}' - (\mathbf{X}' \mathbf{X})^{-1}] \\ &= \sigma^2 [\mathbf{C} \mathbf{C}' - \mathbf{I} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{I}'] \\ &= \sigma^2 [\mathbf{C} \mathbf{C}' - \mathbf{C} \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{C}'] \\ &= \sigma^2 \mathbf{C} [\mathbf{C}' - \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{C}'] \\ &= \sigma^2 \mathbf{C} [\mathbf{I} - \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'] \mathbf{C}' \\ &= \sigma^2 \mathbf{C} [\mathbf{I} - \mathbf{P}] \mathbf{C}' \\ &= \sigma^2 \mathbf{C} \mathbf{M} \mathbf{C}' \\ &= \sigma^2 \mathbf{C} \mathbf{M} \mathbf{M}' \mathbf{C}' \\ &= \sigma^2 \mathbf{C} \mathbf{M} (\mathbf{C} \mathbf{M})' \\ &\geq 0\end{aligned}$$

since $\mathbf{C} \mathbf{X} = \mathbf{I}$, \mathbf{M} is symmetric and idempotent, and $\mathbf{C} \mathbf{M} (\mathbf{C} \mathbf{M})'$ is positive semidefinite (or nonnegative definite) (Greene, A-114).

Statistical properties in finite samples, efficiency

Let \mathbf{q} be a $K \times 1$ vector of constants. Use \mathbf{q} to obtain a linear combination of the true coefficients of a LRM such that

$$\mathbf{q}'\boldsymbol{\beta} = c_1\beta_1 + c_2\beta_2 + \dots + c_K\beta_K.$$

It is easy to verify that $\mathbf{q}'\hat{\beta}_0$ and $\mathbf{q}'\hat{\beta}$ are unbiased estimators of $\mathbf{q}'\boldsymbol{\beta}$. The difference between the variances of these estimators is

$$\begin{aligned}\text{Var} \left[\mathbf{q}'\hat{\beta}_0 \mid \mathbf{X} \right] - \text{Var} \left[\mathbf{q}'\hat{\beta} \mid \mathbf{X} \right] &= \mathbf{q}'\text{Var} \left[\hat{\beta}_0 \mid \mathbf{X} \right] \mathbf{q} - \mathbf{q}'\text{Var} \left[\hat{\beta} \mid \mathbf{X} \right] \mathbf{q} \\ &= \mathbf{q}' \left[\text{Var} \left[\hat{\beta}_0 \mid \mathbf{X} \right] - \text{Var} \left[\hat{\beta} \mid \mathbf{X} \right] \right] \mathbf{q} \\ &= \mathbf{q}'\sigma^2 \mathbf{CM}(\mathbf{CM})' \mathbf{q} \\ &= \sigma^2 \mathbf{q}' \mathbf{CM}(\mathbf{CM})' \mathbf{q} \\ &\geq 0\end{aligned}$$

since $\text{Var}[\mathbf{aX}] = \mathbf{a}\text{Var}[\mathbf{X}]\mathbf{a}'$ for the constant and random vectors \mathbf{a} and \mathbf{X} , and $\mathbf{q}'\mathbf{CM}(\mathbf{CM})'\mathbf{q}$ is positive semidefinite.

Statistical properties in finite samples, efficiency

This shows that the OLS estimator $\hat{\beta}$ yields the smallest variance when it is used to estimate any linear combination of β .

This result allows us to carry out a joint hypothesis test on β using a suitable test statistic which will need to employ the variance of $\mathbf{q}'\hat{\beta}$. E.g., suppose we want to test if

$$\mathbf{q}'\beta = [0 \quad 1 \quad -1] \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = \beta_2 - \beta_3$$

is 0. In a suitable test statistic, we now know that we can use

$$\text{Var} \left[\mathbf{q}'\hat{\beta} \mid \mathbf{X} \right]$$

because we know that

$$\text{Var} \left[\mathbf{q}'\hat{\beta}_0 \mid \mathbf{X} \right] \geq \text{Var} \left[\mathbf{q}'\hat{\beta} \mid \mathbf{X} \right].$$

Is the OLS estimator a linear estimator?

If an estimator is a linear function of the dependent variable, it is a linear estimator. Is $\hat{\beta}$ a linear estimator?

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

$\hat{\beta}$ is a linear function of the values of \mathbf{y} . The values of \mathbf{y} are linearly combined using weights that are a non-linear function of the values of \mathbf{X} . Hence, $\hat{\beta}$ is a **linear estimator** with respect to how it uses the values of the dependent variable only, irrespective of how it uses the values of the regressors.

Is the OLS estimator a linear estimator?

Consider the bivariate LRM

$$\mathbf{y} = \mathbf{x}_0\beta_0 + \mathbf{x}_1\beta_1 + \boldsymbol{\varepsilon}.$$

\mathbf{x}_0 is a column of ones. In this model

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \left(\frac{\sum_{i=1}^n (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) y_i - \left(\frac{\sum_{i=1}^n (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) \bar{y}.$$

$\hat{\beta}_1$ is a linear function of the values of \mathbf{y} .

Statistical properties in finite samples

Gauss-Markov Theorem. In the LRM with regressor matrix \mathbf{X} , the OLS estimator, $\hat{\beta}$, is the minimum variance, linear, unbiased estimator of β . For any vector of constants \mathbf{q} , the minimum variance linear unbiased estimator of $\mathbf{q}'\beta$ in the regression model is $\mathbf{q}'\hat{\beta}$.

Notes on the variance of the OLS estimator

Consider the linear model

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

\mathbf{X} contains the column of ones \mathbf{x}_0 and other k variables.

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

and it can be shown that

$$\text{Var} [\hat{\beta}_j \mid \mathbf{X}] = \frac{\sigma^2}{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 (1 - R_j^2)}, \quad j = 0, 1, \dots, k.$$

where R_j^2 is the coefficient of determination from regressing \mathbf{x}_j on all the other regressors. The expression shows that the variance of the OLS estimator is (i) increasing with the variance of ε , (ii) decreasing with the sample size, (iii) decreasing with the sample variance of \mathbf{X} , and (iv) increasing with R_j^2 .

Notes on the variance of the OLS estimator

By the LIE, it is also true that

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

Notes on the variance of the OLS estimator

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

The problem is that σ^2 is unobserved. An unbiased estimator of σ^2 is

$$\hat{\sigma}^2 = \frac{\hat{\varepsilon}'\hat{\varepsilon}}{n - K}.$$

s is the **standard error of the regression**.

The **estimator** of the variance-covariance matrix of $\hat{\beta}$ is then given by

$$\text{Est. Var} [\hat{\beta} \mid \mathbf{X}] = \hat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1} = \frac{\hat{\varepsilon}'\hat{\varepsilon}}{n - K} (\mathbf{X}'\mathbf{X})^{-1}.$$

Notes on the variance of the OLS estimator

The **estimator** of the standard error of $\hat{\beta}$ is

$$\text{Est. S.E.} [\hat{\beta} \mid \mathbf{X}] = \sqrt{\hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1}}.$$

What is a standard error? Recall the notion of the sampling distribution of $\hat{\beta}$. The standard deviation of the sampling distribution of $\hat{\beta}$ is the standard error of $\hat{\beta}$. This is important! Hence it is an exercise in the lab.

Statistical properties in finite samples, normality

We have shown that

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

Assume for the **first time** that ε is multivariate normal (A6). That is,

$$\varepsilon \mid \mathbf{X} \sim N[\mathbf{0}, \sigma^2 \mathbf{I}].$$

Is $\hat{\beta}$ multivariate normal?

Statistical properties in finite samples, normality

We condition on \mathbf{X} and hence treat it as given. The matrix

$$(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}',$$

is $K \times n$. Recast it as a $K \times n$ matrix

$$\begin{bmatrix} \mathbf{w}_1 & \mathbf{w}_2 & \dots & \mathbf{w}_n \end{bmatrix}.$$

ε is $n \times 1$. $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon$ becomes

$$\mathbf{w}_1\varepsilon_1 + \mathbf{w}_2\varepsilon_2 + \dots + \mathbf{w}_n\varepsilon_n.$$

Hence, $\hat{\beta}$ is a linear combination of the elements of ε . A linear combination of normal random variables is normal. Hence, $\hat{\beta}$ is multivariate normal.

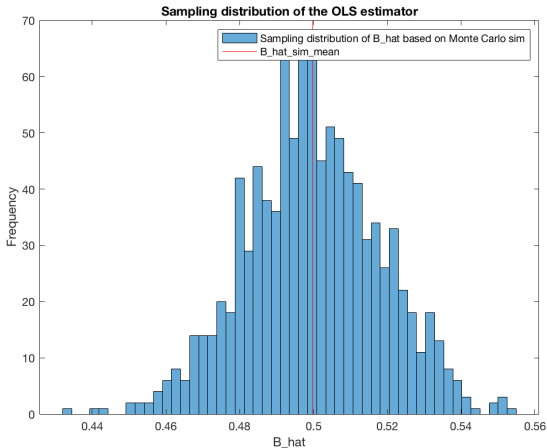
Statistical properties in finite samples, normality

Using the mean and variance-covariance matrix of $\hat{\beta}$ derived above,

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

Statistical properties in finite samples, normality

```
histogram(B_hat_sim)
```



Recall that in our simulation we assumed that the error is normal.

Statistical properties in finite samples, normality

In **finite** sample analysis, the normal distribution of $\hat{\beta}$ is a consequence of the assumption that ε is normal, and that \mathbf{X} is constant. In **large** sample analysis, we will obtain an approximate normal distribution for $\hat{\beta}$, **without assuming** that ε is normal, and \mathbf{X} is constant.

Statistical properties in large samples

We used unbiasedness and efficiency as criteria to judge if $\hat{\beta}$ is a good estimator. These criteria do not depend on n . We would consider new criteria to judge $\hat{\beta}$ based on n . But why?

Statistical properties in large samples

First reason. $\hat{\beta}$ should come closer to the population parameter β if we increase n and come closer to the population N . Who wants an estimator that does not satisfy this?

Statistical properties in large samples

Second reason. Estimators developed to estimate true parameters in complicated models are usually biased. We could check if a biased estimator in a small sample becomes an unbiased estimator in a large sample. E.g., the OLS estimator is biased if the lagged dependent variable is an explanatory variable in the model. However, it is consistent.

Statistical properties in large samples

Third reason. Often the derivation of a property of an estimator is not tractable in a small sample but in a large sample. This is because the expected value of a non-linear function of a statistic is **not** the non-linear function of the expected value of that statistic. But the plim of a non-linear function of a statistic is the non-linear function of the plim of that statistic.

Statistical properties in large samples

Recall the sampling distribution of $\hat{\beta}$ obtained in the simulation experiment. Imagine creating a sequence of sampling distributions of $\hat{\beta}$ with successively larger n . If the distributions in this sequence become more and more similar in form to some specific distribution as n becomes extremely large, this specific distribution is called the **asymptotic distribution** of $\hat{\beta}$.

Statistical properties in large samples, consistency

If the asymptotic distribution of $\hat{\beta}$ becomes concentrated on the particular value β as n approaches infinity, β is said to be the probability limit of $\hat{\beta}$. We then write

$$\hat{\beta} \xrightarrow{p} \beta,$$

or

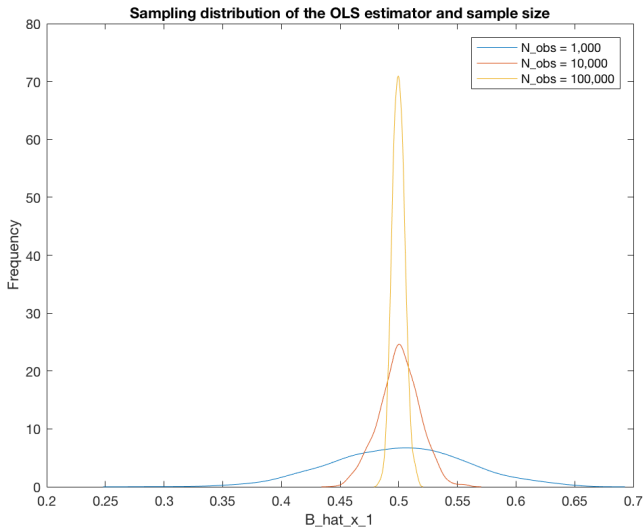
$$\text{plim } \hat{\beta} = \beta,$$

or

$$\hat{\beta} - \beta = o_p(1).$$

We then say that $\hat{\beta}$ is **consistent**. This is our first large sample criterion.

Statistical properties in large samples, consistency



Statistical properties in large samples, consistency

$$\begin{aligned}\hat{\beta} &= \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon \\ &= \beta + \left(\frac{1}{n}\mathbf{X}'\mathbf{X}\right)^{-1}\frac{1}{n}\mathbf{X}'\varepsilon.\end{aligned}$$

Statistical properties in large samples, consistency

$$\begin{aligned}
 \mathbf{X}'\mathbf{X} &= \begin{bmatrix} 1 & 1 & \dots & 1 & \dots & 1 \\ x_{21} & x_{22} & \dots & x_{2i} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{j1} & x_{j2} & \dots & x_{ji} & \dots & x_{jn} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{k1} & x_{k2} & \dots & x_{ki} & \dots & x_{kn} \end{bmatrix} \begin{bmatrix} 1 & x_{21} & \dots & x_{j1} & \dots & x_{k1} \\ 1 & x_{22} & \dots & x_{j2} & \dots & x_{k2} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_{2i} & \dots & x_{ji} & \dots & x_{ki} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_{2n} & \dots & x_{jn} & \dots & x_{kn} \end{bmatrix} \\
 &= \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_i & \dots & \mathbf{x}_n \end{bmatrix} \begin{bmatrix} \mathbf{x}'_1 \\ \mathbf{x}'_2 \\ \vdots \\ \mathbf{x}'_i \\ \vdots \\ \mathbf{x}'_n \end{bmatrix} \\
 &= \sum_{i=1}^n \mathbf{x}_i \mathbf{x}'_i.
 \end{aligned}$$

Statistical properties in large samples, consistency

$$\begin{aligned} \mathbf{X}'\boldsymbol{\varepsilon} &= \begin{bmatrix} 1 & 1 & \dots & 1 & \dots & 1 \\ x_{21} & x_{22} & \dots & x_{2i} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{j1} & x_{j2} & \dots & x_{ji} & \dots & x_{jn} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{k1} & x_{k2} & \dots & x_{ki} & \dots & x_{kn} \end{bmatrix} \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_i \\ \vdots \\ \varepsilon_n \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_i & \dots & \mathbf{x}_n \end{bmatrix} \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_i \\ \vdots \\ \varepsilon_n \end{bmatrix} \\ &= \sum_{i=1}^n \mathbf{x}_i \varepsilon_i. \end{aligned}$$

Statistical properties in large samples, consistency

Hence,

$$\hat{\beta} = \beta + \left(\frac{1}{n} \mathbf{X}' \mathbf{X} \right)^{-1} \frac{1}{n} \mathbf{X}' \varepsilon$$

becomes

$$\hat{\beta} = \beta + \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i.$$

Take the plim of both sides of the equation. When taking the plim, we do not condition on \mathbf{X} . To derive asymptotic results, we do not need the technical simplification brought by fixing \mathbf{X} in repeated samples.

Statistical properties in large samples, consistency

Using the sum rule of plim (Greene, Theorem D.14),

$$\text{plim } \hat{\beta} = \text{plim } \beta + \text{plim } \left[\left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i \right].$$

Using the product rule of plim (Greene, Theorem D.14),

$$\text{plim } \hat{\beta} = \beta + \text{plim } \left[\left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \right] \text{plim } \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i.$$

Statistical properties in large samples, consistency

Assuming that \mathbf{x}_i is i.i.d. (A5), and using the WLLN (Greene, Theorem D.5),

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \xrightarrow{p} E [\mathbf{x}_i \mathbf{x}_i'] .$$

Statistical properties in large samples, consistency

Using the ratio rule of plim (Greene, Theorem D.14), and assuming that $(\mathbf{X}'\mathbf{X})^{-1}$ exists (A2),

$$\left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \xrightarrow{p} (\mathbb{E} [\mathbf{x}_i \mathbf{x}_i'])^{-1}.$$

Statistical properties in large samples, consistency

Assuming that ε_i is i.i.d. (A5), assuming that $E[\mathbf{x}_i \varepsilon_i] = \mathbf{0}$ (A3), and using the WLLN,

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i \xrightarrow{p} E[\mathbf{x}_i \varepsilon_i] = \mathbf{0}.$$

Statistical properties in large samples, consistency

$$\text{plim } \hat{\beta} = \beta + \underbrace{\text{plim} \left[\left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \right]}_{(\mathbb{E}[\mathbf{x}_i \mathbf{x}_i'])^{-1}} \underbrace{\text{plim} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i}_0$$

Hence,

$$\hat{\beta} \xrightarrow{p} \beta.$$

Statistical properties in large samples, consistency

The probability limit can be seen as the large-sample equivalent of the expected value. Hence, $\hat{\beta} \xrightarrow{p} \beta$ can be seen as the large-sample equivalent of unbiasedness.

Statistical properties in large samples, consistency

$\hat{\beta}$ is also consistent under weaker versions of A5 and A3. The former is an exercise in the tutorial. The latter is as follows. To show consistency, we assume

$$E[\mathbf{x}_i \varepsilon_i] = \mathbf{0}$$

which is weak exogeneity, and not

$$E[\mathbf{X} \varepsilon_i] = \mathbf{0}$$

which is strict exogeneity. But remember that in finite sample analysis, for unbiasedness, we have assumed strict exogeneity. This shows that as n increases, we enjoy a weaker model assumption.

Statistical properties in large samples: asy. efficiency

The variance of the asymptotic distribution of $\hat{\beta}$ is called the **asymptotic variance** of $\hat{\beta}$. Among the consistent estimators, if the asymptotic variance of $\hat{\beta}$ is smaller than the asymptotic variance of any other estimator, $\hat{\beta}$ is said to be **asymptotically efficient**. This is our second large sample criterion.

Statistical properties in large samples: asy. efficiency

Suppose that the asymptotic variance of a competitor estimator $\hat{\beta}_0$ is $\mathbf{\Omega}$. Then, under assumptions A1, A2, A3, A4, and A5, it can be shown that

$$\begin{aligned}\text{Asy. Var}(\hat{\beta}_0) - \text{Asy. Var}(\hat{\beta}) &= \mathbf{\Omega} - \frac{\sigma^2}{n} (\text{E} [\mathbf{x}_i \mathbf{x}_i'])^{-1} \\ &\geq 0.\end{aligned}$$

This is the large-sample equivalent of the efficiency criteria we have considered in finite samples. We do not consider the proof.

Statistical properties in large samples: asy. normality

Consider our earlier result

$$\hat{\beta} = \beta + \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i.$$

Rearrange the terms, and multiply both sides of the equation with \sqrt{n} to obtain

$$\sqrt{n} (\hat{\beta} - \beta) = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \sqrt{n} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i.$$

Statistical properties in large samples: asy. normality

We already know that

$$\left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \xrightarrow{p} (\mathbb{E} [\mathbf{x}_i \mathbf{x}_i'])^{-1}.$$

Convergence in probability implies convergence in distribution.

Hence,

$$\left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \xrightarrow{d} (\mathbb{E} [\mathbf{x}_i \mathbf{x}_i'])^{-1}.$$

Note that $\mathbb{E} [\mathbf{x}_i \mathbf{x}_i']$ is a constant matrix and has no variance.

Statistical properties in large samples: asy. normality

Assuming that $E[\mathbf{x}_i \varepsilon_i] = 0$ (A3); assuming that \mathbf{x}_i and ε_i are both i.i.d. (A5), and applying the CLT (Greene, Theorem D.19A),

$$\sqrt{n} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i \xrightarrow{d} N[\mathbf{0}, \sigma^2 E[\mathbf{x}_i \mathbf{x}_i']] .$$

We did **not** assume that ε_i is normal (A6). We are enjoying the CLT!

Statistical properties in large samples: asy. normality

$$\sqrt{n} \left(\hat{\beta} - \beta \right) = \underbrace{\left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1}}_{\xrightarrow{d} (\mathbb{E}[\mathbf{x}_i \mathbf{x}_i'])^{-1}} \underbrace{\sqrt{n} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i}_{\xrightarrow{d} N[\mathbf{0}, \sigma^2 \mathbb{E}[\mathbf{x}_i \mathbf{x}_i']]} .$$

Using the product rule of limiting distributions (Greene, Theorem D.16),

$$\sqrt{n} \left(\hat{\beta} - \beta \right) \xrightarrow{d} (\mathbb{E} [\mathbf{x}_i \mathbf{x}_i'])^{-1} N [\mathbf{0}, \sigma^2 \mathbb{E} [\mathbf{x}_i \mathbf{x}_i']] .$$

Statistical properties in large samples: asy. normality

$$\begin{aligned}\sqrt{n}(\hat{\beta} - \beta) &\xrightarrow{d} (E[\mathbf{x}_i \mathbf{x}_i'])^{-1} N[\mathbf{0}, \sigma^2 E[\mathbf{x}_i \mathbf{x}_i']] \\ &\xrightarrow{d} N\left[\mathbf{0}, \sigma^2 \left((E[\mathbf{x}_i \mathbf{x}_i'])^{-1} \right) E[\mathbf{x}_i \mathbf{x}_i'] \left((E[\mathbf{x}_i \mathbf{x}_i'])^{-1} \right)'\right] \\ &\xrightarrow{d} N\left[\mathbf{0}, \sigma^2 (E[\mathbf{x}_i \mathbf{x}_i'])^{-1} E[\mathbf{x}_i \mathbf{x}_i'] (E[\mathbf{x}_i \mathbf{x}_i'])^{-1}\right] \\ &\xrightarrow{d} N\left[\mathbf{0}, \sigma^2 (E[\mathbf{x}_i \mathbf{x}_i'])^{-1}\right],\end{aligned}$$

using the property that the transpose and inverse operations commute from the second to the third line.

Statistical properties in large samples: asy. normality

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N\left[\mathbf{0}, \sigma^2 (E[\mathbf{x}_i \mathbf{x}_i'])^{-1}\right]$$

is a **limiting distribution**.

This means that when n is **infinite**, $\sqrt{n}(\hat{\beta} - \beta)$ has an **exact** distribution, which is the limiting distribution.

This implies that when n is **finite**, $\sqrt{n}(\hat{\beta} - \beta)$ has an **approximate** distribution, which is close to the limiting distribution. At least this is what we want to believe in, and therefore we assume this. This approximate distribution is also called the asymptotic distribution. The point is that we know that if we let n increase, the approximation is better. We know that if we let n go to infinity, we will get the exact distribution (Greene, Section D.3).

Statistical properties in large samples: asy. normality

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N\left[\mathbf{0}, \sigma^2 (\mathbb{E}[\mathbf{x}_i \mathbf{x}_i'])^{-1}\right]$$

Assuming that this limiting distribution holds approximately for **finite** n ,

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{a} N\left[\mathbf{0}, \sigma^2 (\mathbb{E}[\mathbf{x}_i \mathbf{x}_i'])^{-1}\right].$$

$$(\hat{\beta} - \beta) \xrightarrow{a} N\left[\mathbf{0}, \sigma^2 \frac{1}{n} (\mathbb{E}[\mathbf{x}_i \mathbf{x}_i'])^{-1}\right].$$

$$\hat{\beta} \xrightarrow{a} N\left[\beta, \sigma^2 \frac{1}{n} (\mathbb{E}[\mathbf{x}_i \mathbf{x}_i'])^{-1}\right],$$

or

$$\hat{\beta} \stackrel{a}{\approx} N\left[\beta, \sigma^2 \frac{1}{n} (\mathbb{E}[\mathbf{x}_i \mathbf{x}_i'])^{-1}\right].$$

Statistical properties in large samples: asy. normality

$$\hat{\beta} \stackrel{a}{\sim} N \left[\beta, \sigma^2 \frac{1}{n} (E [\mathbf{x}_i \mathbf{x}_i'])^{-1} \right].$$

σ^2 and $1/n (E [\mathbf{x}_i \mathbf{x}_i'])^{-1}$ are population terms of the limiting distribution. They are unobserved. In practice, they are estimated with $\hat{\sigma}^2$ and $(\mathbf{X}'\mathbf{X})^{-1}$ given sample data, respectively.