

Statistical Properties of OLS

Linear regression model:

$$y_i = X_i\beta + u_i, \quad u_i \sim \text{IID}(0, \sigma^2), \quad i = 1, \dots, N. \quad (1)$$

We may assume that the data were actually generated by the **data-generating process**, or **DGP**,

$$y_i = X_i\beta_0 + u_i, \quad u_i \sim \text{NID}(0, \sigma_0^2), \quad i = 1, \dots, N. \quad (2)$$

Meanings of “ \sim ”, $\text{IID}(0, \sigma^2)$, $\text{NID}(0, \sigma_0^2)$, and the “0” subscripts.

A **model** is a set of DGPs, perhaps denoted by \mathbb{M} . It might consist of all DGPs like (2), where β takes some value in \mathbb{R}^k , σ^2 is a positive real number, and u_i follows a distribution with mean 0 and variance σ^2 .

The set of DGPs (2) defines the **classical normal linear model**.

Bias and Unbiasedness

Let $\hat{\theta}$ be an estimator of some parameter θ , with true value θ_0 . Then the **bias** of $\hat{\theta}$ is defined as $E(\hat{\theta}) - \theta_0$.

If the bias of $\hat{\theta}$ is zero for every admissible value of θ_0 , then $\hat{\theta}$ is said to be **unbiased**.

The model (1) can also be written, using matrix notation, as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad \mathbf{u} \sim \text{IID}(\mathbf{0}, \sigma^2 \mathbf{I}), \quad (3)$$

where \mathbf{y} and \mathbf{u} are N -vectors, \mathbf{X} is an $N \times k$ matrix, $\boldsymbol{\beta}$ is a k -vector, and $\text{IID}(\mathbf{0}, \sigma^2 \mathbf{I})$ is shorthand.

We can replace \mathbf{y} by $\mathbf{X}\boldsymbol{\beta}_0 + \mathbf{u}$ in $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$. This yields

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\beta}_0 + \mathbf{u}) = \boldsymbol{\beta}_0 + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u}. \quad (4)$$

The expectation of (4) is

$$E(\hat{\beta}) = \beta_0 + E((X^T X)^{-1} X^T u). \quad (5)$$

Thus $\hat{\beta}$ is unbiased if and only if the second term on the r.h.s. is $\mathbf{0}$.

Unfortunately, few estimators in econometrics are unbiased.

If the matrix X is **fixed** or **nonstochastic**, then

$$E((X^T X)^{-1} X^T u) = (X^T X)^{-1} X^T E(u). \quad (6)$$

For example, X might be $\mathbf{1}$, a vector of 1s.

In this case, $\hat{\beta}$ is unbiased whenever $E(u) = \mathbf{0}$.

Can we treat the matrix X as fixed? Yes, for experimental data, where the experimenter chose X before nature generated y .

But the assumption that X is fixed is often unreasonable. Some of the columns of X typically correspond to variables that are random.

A weaker assumption is that the regressors in \mathbf{X} are **exogenous**, so that \mathbf{X} is independent of \mathbf{u} . Therefore,

$$\mathbb{E}(\mathbf{u} \mid \mathbf{X}) = \mathbf{0}. \quad (7)$$

The expectation of the entire vector \mathbf{u} , that is, of every one of the u_i , is assumed to be zero conditional on the entire matrix \mathbf{X} .

Equation (7) is an **exogeneity assumption**. Given (7), it is easy to show that $\hat{\beta}$ is unbiased.

Because $\mathbb{E}(\mathbf{X}^\top \mid \mathbf{X}) = \mathbf{X}^\top$, and $\mathbb{E}(\mathbf{u} \mid \mathbf{X}) = \mathbf{0}$, it is clear that $\mathbb{E}(\mathbf{X}^\top \mathbf{u}) = \mathbf{0}$. Similarly,

$$\mathbb{E}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u} \mid \mathbf{X}) = \mathbf{0}. \quad (8)$$

We apply the Law of Iterated Expectations to obtain (8).

The exogeneity assumption (7) may be reasonable for **cross-section data**, but it makes no sense for **time-series data**.

For time-series data, we index observations by t .

Even if u_t is not related to current and past values of the regressors, it must be related to future values if current values of y_t affect them.

The **predeterminedness** condition

$$E(u_t | \mathbf{X}_t) = 0, \quad t = 1, \dots, T, \quad (9)$$

is much weaker than (7). We can say that the regressors are **predetermined** with respect to the disturbances.

The **first-order autoregressive**, or **AR(1)**, model is

$$y_t = \beta_1 + \beta_2 y_{t-1} + u_t, \quad u_t \sim \text{IID}(0, \sigma^2). \quad (10)$$

In matrix terms, this is

$$\mathbf{y} = \beta_1 \mathbf{1} + \beta_2 \mathbf{y}_1 + \mathbf{u}, \quad \mathbf{u} \sim \text{IID}(\mathbf{0}, \sigma^2 \mathbf{I}), \quad (11)$$

where \mathbf{y}_1 has typical element y_{t-1} .

For the AR(1) model, the predeterminedness condition (9) may well hold, but the exogeneity assumption (7) cannot possibly hold.

The AR(1) model is **dynamic** and involves a **lagged dependent variable**. OLS estimates of such models are never unbiased.

For the AR(1) model,

$$\hat{\beta}_2 = (\mathbf{y}_1^\top \mathbf{M}_\iota \mathbf{y}_1)^{-1} \mathbf{y}_1^\top \mathbf{M}_\iota \mathbf{y}, \quad \text{where} \quad \mathbf{M}_\iota = \mathbf{I} - \iota(\iota^\top \iota)^{-1} \iota^\top. \quad (12)$$

If we replace \mathbf{y} by $\beta_{10}\iota + \beta_{20}\mathbf{y}_1 + \mathbf{u}$, we find that

$$\hat{\beta}_2 = \beta_{20} + (\mathbf{y}_1^\top \mathbf{M}_\iota \mathbf{y}_1)^{-1} \mathbf{y}_1^\top \mathbf{M}_\iota \mathbf{u}. \quad (13)$$

The second term here does not have expectation 0.

Because \mathbf{y}_1 is stochastic, we cannot simply move the expectations operator and then take the unconditional expectation of \mathbf{u} .

Because $E(\mathbf{u} | \mathbf{y}_1) \neq \mathbf{0}$, we cannot take expectations conditional on \mathbf{y}_1 and then rely on the Law of Iterated Expectations.

Asymptotic Theory and Consistency

Exact results are often unavailable, so we rely on **asymptotic** results.

Asymptotic theory requires an **asymptotic construction**, which specifies how samples of arbitrarily large size are generated.

This is easy with cross-section data. Just draw more and more observations from a population of infinite size.

For pure time-series models, it is common to let the simulations run on forever. But does this make sense?

With an $M \times k$ regressor matrix, stack copies of X to create matrices of dimensions $2M \times k$, $3M \times k$, $4M \times k$, and so on.

In most cases, the parameters are held fixed as $N \rightarrow \infty$. However, it is sometimes necessary to let them change systematically with N .

Examples include time-series models with unit roots and cointegration and cross-section models with weak instruments.

Consider the model

$$y_t = \beta_1 + \beta_2 1/t + u_t, \quad u_t \sim \text{IID}(0, \sigma^2). \quad (14)$$

Here $\hat{\beta}_1$ and $\hat{\beta}_2$ are unbiased, and the values of the second regressor range from 1 to $1/N$.

If we use (14) directly as an asymptotic construction, the regressor gets smaller and smaller as $N \rightarrow \infty$.

If we stack copies of the original X , it remains bounded.

We discuss two forms of **stochastic convergence**: **Convergence in probability** and **convergence in distribution**.

The first allows us to find the **probability limit**, or **plim**, of a sequence of random variables.

Let $\{Y_N\}$ denote a sequence of random variables Y_N , $N = 1, \dots, \infty$. If the sequence converges in probability, then

$$\text{plim}_{N \rightarrow \infty} Y_N = Y_\infty. \quad (15)$$

For equation (15) to be true, what we need is that, for all $\epsilon > 0$,

$$\lim_{N \rightarrow \infty} \Pr(|Y_N - Y_\infty| > \epsilon) = 0. \quad (16)$$

The limit here is the ordinary limit of a sequence of real numbers.

For any chosen ϵ , however small, we can find N^* large enough so that, for all $N > N^*$, the probability that $|Y_N - Y_\infty|$ is greater than ϵ is arbitrarily small.

In (15), Y_∞ is often deterministic, but it may be stochastic.

For a sequence of random vectors, or matrices, denoted by $\{Y_N\}$, $\text{plim}_{N \rightarrow \infty} Y_N = Y_\infty$ means that

$$\lim_{N \rightarrow \infty} \Pr(\|Y_N - Y_\infty\| > \epsilon) = 0. \quad (17)$$

Here $\|\cdot\|$ denotes the Euclidean norm of a vector. It can also denote a matrix norm.

A second, weaker, form of convergence is **convergence in distribution**, or **convergence in law**, or **weak convergence**.

The random variables themselves no longer converge, but instead the sequence of CDFs of the random variables converges.

If the random variables are IID, then they all have the same CDF, and so the limiting CDF is simply the CDF of each element of the sequence.

In general, a sequence $\{Y_N\}$ of scalar random variables converges in distribution to the distribution characterized by the CDF F if

$$\lim_{N \rightarrow \infty} F_N(x) = F(x), \quad (18)$$

where F_N is the CDF of Y_N , for all real x at which F is continuous.

We can write

$$Y_N \xrightarrow{d} F. \quad (19)$$

If a sequence $\{Y_N\}$ converges in probability, it also converges in distribution.

But convergence in distribution of a sequence does not imply convergence in probability.

The two concepts coincide whenever the limiting distribution is **degenerate**. All the probability mass is concentrated on one single point, which is a **nonstochastic plim**.

Suppose Z_t is a random variable equal to 1 if a coin comes up heads, and equal to 0 if it comes up tails. After N tosses, the proportion of heads is just

$$Y_N \equiv \frac{1}{N} \sum_{t=1}^N Z_t. \quad (20)$$

If the coin really is unbiased, $E(Y_N) = 1/2$. Thus, if Y_N converges to anything, it must be to $1/2$.

Laws of Large Numbers

Suppose \bar{y}_N is the sample mean of y_i , $i = 1, \dots, N$, a sequence of random variables, each with expectation μ_y .

Provided the y_i are independent (or at least, not too dependent), a law of large numbers, or LLN, would state that

$$\text{plim}_{N \rightarrow \infty} \bar{y}_N = \text{plim}_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N y_i = \mu_y. \quad (21)$$

\bar{y}_N has a nonstochastic plim equal to the common expectation of each of the y_i .

If the y_i are IID, with variance σ^2 ,

$$E(\bar{y}_N) = \frac{1}{N} \sum_{i=1}^N E(y_i) = \frac{1}{N} \sum_{i=1}^N \mu_y = \mu_y. \quad (22)$$

It is easy to see that

$$\text{Var}(\bar{y}_N) = \frac{1}{N^2} \sum_{i=1}^N \sigma^2 = \frac{1}{N} \sigma^2. \quad (23)$$

Thus \bar{y}_N has mean μ_y and a variance which tends to zero as $N \rightarrow \infty$.

In the limit, \bar{y}_N becomes a nonstochastic quantity equal to μ_y .

As $N \rightarrow \infty$, we are collecting more information about $E(y_i)$, with each y_i providing a smaller fraction of that information.

Eventually, the random components of the individual y_i cancel out, and the sample mean \bar{y}_N converges to the population mean μ_y .

We need assumptions to prevent any of the y_i from having too much impact on \bar{y}_N . Thus $\text{Var}(y_i)$ must be bounded from above.

We also need to assume that there is not too much dependence among the y_i , to ensure that the random components $y_i - \mu_y$ cancel out.

The assumption that the y_i are $\text{IID}(0, \sigma^2)$ is sufficient.

Convenient Properties of Plims

Suppose that $\{Y_N\}$, $N = 1, \dots, \infty$, is a sequence of random variables with nonstochastic plim Y_∞ , and $\eta(Y_N)$ is a smooth function of Y_N . Then $\text{plim } \eta(Y_N) = \eta(Y_\infty)$.

If $\{Y_N\}$ and $\{Z_N\}$ converge in probability, then

$$\text{plim } Y_N Z_N = \text{plim } Y_N \text{plim } Z_N. \quad (24)$$

These features of plims are not shared by expectations. When $\eta(\cdot)$ is nonlinear,

$$E(\eta(Y)) \neq \eta(E(Y)). \quad (25)$$

Also, $E(YZ) \neq E(Y)E(Z)$ unless Y and Z are independent.

Many stochastic quantities do not have plims unless we divide them by N or, perhaps, by some power of N .

Consider the matrix $\mathbf{X}^\top \mathbf{X}$. Each element of this matrix is a sum of N products of, say, X_{ij} with X_{il} .

As $N \rightarrow \infty$, such a sum would tend to infinity as well. Therefore, the matrix $\mathbf{X}^\top \mathbf{X}$ does not generally have a plim. Instead, we assume that

$$\text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \mathbf{X} = \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}, \quad (26)$$

where $\mathbf{S}_{\mathbf{X}^\top \mathbf{X}}$ is a finite nonstochastic matrix.

Each element of the matrix $N^{-1} \mathbf{X}^\top \mathbf{X}$ is now an average of N numbers:

$$\left(\frac{1}{N} \mathbf{X}^\top \mathbf{X} \right)_{jl} = \frac{1}{N} \sum_{i=1}^N x_{ij} x_{il}. \quad (27)$$

We are implicitly assuming that an LLN holds for the sequences generated by the squares of the regressors and their cross-products.

The Same-Order Notation

Many quantities in econometrics change systematically as the sample size changes. They can do so at different rates.

The **same-order relation** provides a very convenient way to deal with such quantities.

Let $f(N)$ be a real-valued function of the integer $N > 0$, and r be a rational number.

Then $f(N)$ is of the same order as N^r if there exists a constant K , independent of N , and a positive integer N^* such that

$$\left| \frac{f(N)}{N^r} \right| < K \text{ for all } N > N^*. \quad (28)$$

When $f(N)$ is of the same order as N^r , we can write $f(N) = O(N^r)$.

Many quantities in econometrics are stochastic. For them, we need the **stochastic same-order relation**.

Let $\{a_N\}$ be a sequence of random variables indexed by N . Then a_N is of order N^r in probability if, for all $\epsilon > 0$, there exist a constant K and a positive integer N^* such that

$$\Pr \left(\left| \frac{a_N}{N^r} \right| > K \right) < \epsilon \text{ for all } N > N^*. \quad (29)$$

When a_N is of order N^r in probability, we can write $a_N = O_p(N^r)$.

The deterministic and stochastic same-order relations have the same properties. ETM often uses “ $O(\cdot)$ ” when it should be using “ $O_p(\cdot)$.”

We can manipulate the same-order relations as if they were simply powers of N . This makes them extremely useful.

Suppose $f(N)$ and $g(N)$ are $O(N^r)$ and $O(N^q)$, respectively. Then

$$f(N)g(N) = O(N^r)O(N^q) = O(N^{r+q}) \quad (30)$$

and

$$f(N) + g(N) = O(N^r) + O(N^q) = O(N^{\max(r,q)}). \quad (31)$$

We have assumed that $N^{-1}\mathbf{X}^\top\mathbf{X}$ has a probability limit of $\mathbf{S}_{\mathbf{X}^\top\mathbf{X}}$, which is a finite, positive definite, deterministic matrix.

Because each element of $N^{-1}\mathbf{X}^\top\mathbf{X}$ is $O_p(1)$ (or maybe $O(1)$ if the regressors are not stochastic), we can write $\mathbf{X}^\top\mathbf{X} = O_p(N)$.

Since $O(1) = O(N^0)$, the same-order relationship is still about N even when N does not explicitly appear.

In general, sums of random variables that *do not* have mean zero are $O_p(N)$. But sums of random variables that *do* have mean zero are $O_p(N^{1/2}) = O_p(\sqrt{N})$.

If we assume that $\mathbf{X}^\top \mathbf{y} = O_p(N)$, we see that

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = O_p(N^{-1}) O_p(N) = O_p(1). \quad (32)$$

Equation (32) says that $\hat{\boldsymbol{\beta}}$ does not systematically get larger or smaller as $N \rightarrow \infty$.

Econometricians mainly use same-order notation to make asymptotic theory easier. For example, suppose that $\hat{\theta} = (A + B) / \sqrt{C + D}$.

If we can show that $A = O_p(N^{1/2})$, $B = O_p(1)$, $C = O_p(N)$, and $D = O_p(N^{1/2})$, then it must be the case that

$$\hat{\theta} \stackrel{a}{=} \frac{A}{\sqrt{C}} = O_p(1). \quad (33)$$

The asymptotic approximation is likely to be good if $A - B \cong A$ and $D \ll C$. But it might be very poor if either or both of these conditions does not hold.

Consistency

Even when the OLS estimator is biased, it may turn out to be **consistent**.

Given a model \mathbb{M} , an estimator $\hat{\beta}$, and an asymptotic construction that allows $\hat{\beta}$ to be defined for arbitrary sample size N , $\hat{\beta}$ is consistent if

$$\text{plim}_{N \rightarrow \infty} \hat{\beta} = \beta_0, \quad (34)$$

where β_0 is the value of β associated with the actual DGP, which must belong to \mathbb{M} .

When the DGP is a special case of a regression model with $\beta = \beta_0$,

$$\hat{\beta} = \beta_0 + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u}. \quad (35)$$

To prove that $\hat{\beta}$ is consistent, we need to show that the plim of the second term is zero.

Neither $\mathbf{X}^\top \mathbf{X}$ nor $\mathbf{X}^\top \mathbf{u}$ has a plim. However, we can divide both of them by N . Then the plim of the second term in (35) is

$$\left(\text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \mathbf{X} \right)^{-1} \text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \mathbf{u} = (\mathbf{S}_{\mathbf{X}^\top \mathbf{X}})^{-1} \text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \mathbf{u} = \mathbf{0}. \quad (36)$$

Here we assume that $\mathbf{S}_{\mathbf{X}^\top \mathbf{X}}$ is nonsingular. Note that, even if $\mathbf{X}^\top \mathbf{X}$ is nonsingular for any finite $N \geq k$, $\mathbf{S}_{\mathbf{X}^\top \mathbf{X}}$ can be singular.

The second equality in (36) depends on the assumption that $E(\mathbf{X}_i^\top \mathbf{u}_i | \mathbf{X}_i) = \mathbf{0}$. Then the Law of Iterated Expectations tells us that $E(\mathbf{X}_i^\top \mathbf{u}_i) = \mathbf{0}$.

Assuming that we can apply a law of large numbers,

$$\text{plim}_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \mathbf{u} = \text{plim}_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \mathbf{X}_i^\top \mathbf{u}_i = \mathbf{0}. \quad (37)$$

Consistency is neither weaker nor stronger than unbiasedness.

Consider again the model

$$y_t = \beta_1 + \beta_2 1/t + u_t, \quad u_t \sim \text{IID}(0, \sigma^2). \quad (38)$$

$\hat{\beta}_1$ and $\hat{\beta}_2$ are evidently unbiased.

But $\hat{\beta}_2$ is not consistent if we use the, apparently natural, asymptotic construction in which $t \rightarrow \infty$.

As $N \rightarrow \infty$, each observation provides less and less information about β_2 , which causes $S_{X^\top X}$ to be singular.

Therefore, (36) does not hold, and the second term on the r.h.s. of (35) does not have a plim of zero.

Note that $\hat{\beta}_1$ is consistent even though $\hat{\beta}_2$ is not.

However, $\hat{\beta}_2$ is consistent under the alternative asymptotic construction in which we stack copies of the original X matrix.

There are two types of **inconsistent estimators**:

- An estimator, possibly unbiased, that does not tend to any nonstochastic plim.
- An estimator that tends to the wrong nonstochastic plim.

Suppose we estimate the population mean, μ , from a sample $y_i, 1 = 1, \dots, N$. The sample mean \bar{y} is unbiased and consistent under reasonable assumptions.

Now consider three (not very sensible) estimators:

$$\hat{\mu}_1 = \frac{1}{N+1} \sum_{i=1}^N y_i. \quad (39)$$

$\hat{\mu}_1$ is biased but consistent. It equals $N/(N+1)$ times \bar{y} . Thus its mean is $(N/(N+1))\mu$, which tends to μ as $N \rightarrow \infty$, and it is consistent whenever \bar{y} is.

$$\hat{\mu}_2 = \frac{1.01}{N} \sum_{i=1}^N y_i. \quad (40)$$

$\hat{\mu}_2$ is clearly biased and inconsistent. Its mean is 1.01μ , and it actually tends to a plim of 1.01μ as $N \rightarrow \infty$.

$$\hat{\mu}_3 = 0.01y_1 + \frac{0.99}{N-1} \sum_{i=2}^N y_i. \quad (41)$$

$\hat{\mu}_3$ is clearly unbiased, since it is a weighted average of two estimators, y_1 and the average of y_2 through y_N , each of which is unbiased.

But it is not consistent, because it does not converge to a nonstochastic plim. Instead, it converges to the random quantity $0.99\mu + 0.01y_1$.

Some types of inconsistency are unavoidable. For example, individual fixed effects generally cannot be estimated consistently.