

Chapter 3: the statistical property of OLS

Compared to the numeric property of OLS (e.g., it minimizes the MSE), the statistical property of OLS (e.g., unbiasedness, efficiency, consistency etc.) can be hold only when the assumption made is true.

We denote the Data Generating Process (**DGP**) as the mechanism that dominates the reality (e.g., to generate the data). In econometrics, we develop models to describe the DGP or at least to partially describe it. Therefore, we may assume that our regression model is correctly specified, i.e., it embeds the DGP (in practice, however, the real DGP for the economy can be very sophisticated and usually we do not pretend to know every detail of it – that is why ‘all the models are wrong, though some are useful’). For example, suppose that we assume that the **DGP** is:

$$y_t = X_t\beta_0 + u_t, u_t \sim NID(0, \sigma_0^2) \quad (3.02)$$

where X_t represent a column vector for the t^{th} observation. we can develop the following **model**:

$$y_t = X_t\beta + u_t, u_t \sim IID(0, \sigma^2) \quad (3.01)$$

(3.02) is embedded in (3.01) as (3.02) is simply (3.01) by taking specific values for the parameters (e.g., β, σ_0) and restricting the distribution which the error term follows to be the Normal distribution.

We assume that our model is correctly specified – only then we can talk about the statistical properties of our models. The challenge is that in practice the DGP may not be generated by (3.02), and the property of the OLS estimator for the model (3.01) becomes questionable.

3.2 Are OLS estimators unbiased?

Following the example which we show in (3.01) and (3.02), we can rewrite our model in (3.01) in a matrix notation as:

$$y = X\beta + u, u \sim IID(0, \sigma^2 I)$$

Or

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} \beta + \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix}, u_t \sim IID(0, \sigma^2)$$

The OLS estimator can be written as:

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

Suppose our model is correctly specified and the true GDP is indeed generated by (3.02). We can insert (3.02) into the OLS estimator, and we have:

$$\begin{aligned} \hat{\beta} &= (X'X)^{-1}X'(X\beta_0 + u) \\ \hat{\beta} &= \beta_0 + (X'X)^{-1}X'u \end{aligned}$$

If we take expectations on both sides, we have:

$$E(\hat{\beta}) = \beta_0 + E((X'X)^{-1}X'u) \quad (3.06)$$

Thus, the OLS estimator is unbiased when the second component of (3.06) equals to zero.

Sometimes it is reasonable to consider X as fixed or non-stochastic. For example, in an experiment (though this normally happens in the statistical, not econometrical, contexts), we can control the values of X . Thus, we have:

$$E(\hat{\beta}) = \beta_0 + (X'X)^{-1}X'E(u) = \beta_0$$

In this case, the OLS estimator is unbiased. However, having fixed values of X is very rarely the case in practice. A weaker assumption is **exogeneity**, which suggests that any randomness in the DGP which generate X is **(mean) independent** of the error term u , or alternatively speaking, the expected values of all the error terms are mean independent to the all the independent variables, e.g.,

$$E(u|X) = 0$$

We can write the equation in a vector format:

$$E\left(\begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix} \middle| \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}\right) = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} = 0$$

where u_1 and X_1 are the error term and the column vector for the explanatory variables for the 1st observation. The equation suggests that the mean of the entire vector u , e.g., $\begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix}$, or alternatively speaking, the mean of every single element of the \mathbf{u} (e.g., u_1, u_2, \dots, u_n) is zero conditional on the entire matrix of \mathbf{X} , e.g., $\begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}$. Since by definition $E(\mathbf{u}|\mathbf{X})$ is a vector of variables (e.g., $E(u_1|\mathbf{X}), E(u_2|\mathbf{X}), \dots, E(u_n|\mathbf{X})$) which depend on the value of \mathbf{X} (e.g., all the elements including X_1, X_2, \dots, X_n), the equation $E(\mathbf{u}|\mathbf{X}) = \mathbf{0}$ suggests that the mean of every single element of the \mathbf{u} vector would be zero no matter what values any element of the matrix \mathbf{X} (e.g., and thus, each of the X_1, X_2, \dots, X_n) could possibly take. Or alternatively speaking, the vector variable \mathbf{u} (e.g., $\begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix}$) is mean independent of the vector variable \mathbf{X} (e.g., $\begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}$).

If the assumption of exogeneity holds, by using the Law of Iterated Expectations, it can be proven that:

$$E((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}) = \mathbf{0}$$

Thus, we have:

$$E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}_0 + E((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}) = \boldsymbol{\beta}_0$$

Thus, if we assume exogeneity, we can conclude that the OLS estimator is unbiased.

Usually, it is reasonable to assume exogeneity for cross-sectional data. For example, each observation represents an individual firm, person, or city. We may usually consider the observations to be collected from random samples from the population. Thus, we can assume that the error term for any of the observations/individuals tends to have a mean value (of zero) irrelevant to the values of the independent variable for **all** the observations/individuals (e.g., X_1, X_2, \dots, X_n). e.g.,

$$E\left(\begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix} \middle| \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}\right) = E(\mathbf{u}|\mathbf{X}) = E(\mathbf{u}) = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} = \mathbf{0}$$

For example, we may model a person's salary based on his/her working experience. The error term contains other influencing factors to the person's salary but not included in the model, such as the person's intelligence. It may be reasonable to assume that the expected impact by one person's intelligence on the person's salary is mean independent of everyone's working experience (and we assume the impact by the person's intelligence conditional on everyone's

working experience equals to zero, as we assume that the expected impact by one person's intelligence on the person's salary is zero). Cross-sectional data are heavily involved in applied microeconomics and other social sciences.

The assumption of exogeneity does not hold for time series data when there are lagged dependent variables used as independent variables. See next.

The OLS estimator can be biased

The OLS estimator can be biased when there are lagged dependent variables in the model (e.g., we call this type of model as '**autoregressive**' models), e.g.,

$$y_t = \beta y_{t-1} + u_t, u_t \sim IID(0, \sigma^2)$$

This equation can be written in a vector/matrix form:

$$\mathbf{y} = \mathbf{y}_1 \boldsymbol{\beta} + \mathbf{u}, \mathbf{u} \sim IID(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (3.11)$$

(3.11) can be written equivalently as (again, suppose we have a sample of size n):

$$\begin{bmatrix} y_2 \\ \vdots \\ y_n \\ y_{n+1} \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_{n-1} \\ y_n \end{bmatrix} \boldsymbol{\beta} + \begin{bmatrix} u_2 \\ \vdots \\ u_n \\ u_{n+1} \end{bmatrix}$$

Thus, we can assume predeterminedness for (3.11): that is, $E(u_t | x_t) = E(u_t) = 0$, in this case, we have the lagged dependent variable as our independent variable (e.g., y_{t-1}), thus we have the predeterminedness assumption as: $E(u_t | y_{t-1}) = E(u_t) = 0$, where $t = 1, \dots, n+1$, for a sample of size n . This assumption could be reasonable because we always realize the value of y_{t-1} before u_t , and we can assume that u_t is mean independent of y_{t-1} , or **alternatively speaking we assume that y_{t-1} has no impact on the expected value of u_t .**

However, in this case, we cannot assume exogeneity. The assumption of exogeneity is $E(\mathbf{u} | \mathbf{X}) = \mathbf{0}$. In this case, we have the lagged dependent variable as our independent variable, thus the assumption of exogeneity is $E(\mathbf{u} | \mathbf{y}_1) = E(\mathbf{u})$, which requires that \mathbf{y}_1 has no impact on the expected value of any element of the vector \mathbf{u} . This can be written as:

$$E(\mathbf{u} | \mathbf{y}_1) = E \left(\begin{bmatrix} u_2 \\ \vdots \\ u_n \\ u_{n+1} \end{bmatrix} \middle| \begin{bmatrix} y_1 \\ \vdots \\ y_{n-1} \\ y_n \end{bmatrix} \right) = E \left(\begin{bmatrix} u_2 \\ \vdots \\ u_n \\ u_{n+1} \end{bmatrix} \right)$$

We already assumed predeterminedness, e.g., $E(u_t|y_{t-1}) = E(u_t) = 0$, thus we have $E(u_t y_{t-1}) = 0$. However, we have:

$$\begin{aligned} E(u_t y_t) &= E(u_t(\beta y_{t-1} + u_t)) = \beta E(u_t y_{t-1} + u_t u_t) \\ &= \beta E(u_t y_{t-1}) + E(u_t u_t) \\ &= E(u_t u_t) = E(u_t^2) \end{aligned}$$

Thus, unless the error term is always zero, $E(u_t y_t) \neq 0$. However, y_t is the regressor for the error term u_{t+1} . Thus, the exogeneity of $E(\mathbf{u}|\mathbf{y}_1) = E(\mathbf{u})$ cannot hold. Intuitively, this is because y_t depends on y_{t-1} , and thus u_{t-1} (according to the model itself by definition, $y_t = \beta_1 + \beta_2 y_{t-1} + u_t$). Thus, we cannot assume that the error terms are mean independent of the regressor.

Thus, if we have the lagged dependent variables as explanatory variables, the second component corresponding to (3.06) do not equal to zero, and the OLS estimator is biased. It can be proved that, not only the parameter for the lagged dependent variable (as the independent variable) would be biased but all the parameters in the autoregressive models from an OLS estimator will be biased.

3.3 Are OLS estimator consistent?

Consistency means that when the sample size goes to infinity, the estimate tends to be the quantity being estimated. That is, when the sample size is large enough, we will be confident that the estimate will be close to the true value. The OLS estimator will be often consistent even when being biased.

How should we think of the sample size of infinity? For cross-sectional data we can pretend that we draw the sample from a population of infinity size (e.g., unlimited number of firms or persons). We simply draw more and more observations from the population. For time series data we can pretend that the time goes back to unlimited history and goes forward to the unlimited future so that we have unlimited observations, and we simply try to observe for longer and longer periods of time.

Probability limit

The probability limit, or **plim**, of the function of a vector y^n is defined as:

$$\text{plim}_{n \rightarrow \infty} \alpha(y^n) = \alpha_0$$

y^n is usually written as y , omitting the vector size n .

Thus when $n \rightarrow \infty$, the vector $\alpha(y^n)$ approaches the vector of α_0 . α_0 could be a vector of random or non-random variables.

Under the definition of **plim**, $\alpha(y^n)$ and α_0 have the following relationship, for $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} \Pr(\| \alpha(y^n) - \alpha_0 \| < \varepsilon) = 1$$

Thus, for a specific tolerance level ε , no matter how small, when $n \rightarrow \infty$, the probability that the norm of the discrepancy (e.g., Euclidean distance) between the two vectors $\alpha(y^n)$ and α_0 is less than ε goes to unity.

Some vectors originally do not have **plim**, but they may have **plim** when they are divided by some quantity (e.g., n or the power of n), e.g., $\mathbf{X}'\mathbf{X}$ does not have a plim, but we can divide $\mathbf{X}'\mathbf{X}$ by n and we have:

$$\text{plim}_{n \rightarrow \infty} \frac{\mathbf{X}'\mathbf{X}}{n} = \mathbf{S}_{\mathbf{X},\mathbf{X}} \quad (3.17)$$

Each element in $\mathbf{X}'\mathbf{X}$ is a scalar product between two columns in \mathbf{X} .

Thus, we have $\mathbf{S}_{\mathbf{X},\mathbf{X}}$ as a finite non-random matrix with a full rank of k . This is because each element in $\frac{\mathbf{X}'\mathbf{X}}{n}$ is an average of n numbers:

$$\left(\frac{\mathbf{X}'\mathbf{X}}{n} \right)_{ij} = \frac{1}{n} \sum_{t=1}^n x_{ti} x_{tj}$$

we have made implicit assumptions that multiple variables in \mathbf{X} do not have too strong dependence (no perfect multicollinearity).

The OLS estimator is consistent

When the model is correctly specified, we have:

$$\hat{\beta} = \beta_0 + (X'X)^{-1}X'u$$

The OLS estimator is consistent when $(X'X)^{-1}X'u$ has a plim of zero. e.g., we know that:

$$\text{plim}_{n \rightarrow \infty} (X'X)^{-1}X'u = \frac{\text{plim}_{n \rightarrow \infty} X'u}{\text{plim}_{n \rightarrow \infty} X'X} = \frac{\text{plim}_{n \rightarrow \infty} \frac{X'u}{n}}{\text{plim}_{n \rightarrow \infty} \frac{X'X}{n}} = \frac{\text{plim}_{n \rightarrow \infty} \frac{X'u}{n}}{S_{X'X}} = \frac{\text{plim}_{n \rightarrow \infty} \frac{\sum_{t=1}^n X'_t u_t}{n}}{S_{X'X}}$$

As we assume predeterminedness, e.g.,

$$E(u_t | X_t) = 0$$

According to LLN, we know that $E(X_t u_t) = 0$, thus the previous equation becomes:

$$\text{plim}_{n \rightarrow \infty} (X'X)^{-1}X'u = \frac{\text{plim}_{n \rightarrow \infty} \frac{\sum_{t=1}^n X'_t u_t}{n}}{S_{X'X}} = \frac{0}{S_{X'X}} = 0$$

Thus, the OLS estimator is consistent, although it may not be unbiased (e.g., when we have lagged dependent variable as the explanatory variable). In this kind, when we have large sample, we can trust the OLS estimate to be very close to the true value of the parameters.

Unbiasedness and consistency are two different things. Sometimes the estimator is biased but consistent, and sometimes it is unbiased but inconsistent (e.g., either the estimator does not tend to any non-random probability limit or tend to a wrong probability limit). See examples in Davison and Mackinnon (2003).

The covariance matrix of the OLS estimator

Until now, we know that the OLS estimator of the parameters to be unbiased and consistent (with the assumption of exogeneity) or biased and consistent (only with predeterminedness), but that is not enough. We would like to know the distribution of the OLS estimate (e.g., $\hat{\beta}$). We are usually interested in the central second moment of the distribution of $\hat{\beta}$. e.g.,

The OLS covariance matrix

Previously we have introduced the variance-covariance of the error term, e.g.,

$$Var(\mathbf{u}) = E(\mathbf{u}\mathbf{u}') = \sigma^2 \mathbf{I} \quad (3.26)$$

Note that, this only requires the covariance of the error terms to be zero, this does not require the error terms are independent (e.g., **IID** is a stronger assumption).

Thus, we would like to calculate $Var(\hat{\boldsymbol{\beta}})$, as we know that:

$$\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}_0 + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}$$

Thus,

$$\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0 = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}$$

As we assume the OLS estimator is unbiased, based on (3.22), we have the variance matrix of $\hat{\boldsymbol{\beta}}$ as follows:

$$\begin{aligned} Var(\hat{\boldsymbol{\beta}}) &= E\left((\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)'\right) = E(\hat{\boldsymbol{\beta}}\hat{\boldsymbol{\beta}}') = E((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u})' \\ &= E((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}\mathbf{u}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}) \\ &= E(\mathbf{u}\mathbf{u}')E((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}) \\ &= \sigma_0^2 E((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}) \\ &= \sigma_0^2 E((\mathbf{X}'\mathbf{X})^{-1}) \end{aligned}$$

(3.28p)

where $Var(\hat{\boldsymbol{\beta}})$ is random variable depending on \mathbf{X} .

We usually use $(\mathbf{X}'\mathbf{X})^{-1}$ as an estimate of $E((\mathbf{X}'\mathbf{X})^{-1})$, and write (3.28p) as:

$$Var(\hat{\boldsymbol{\beta}}) = \sigma_0^2 E((\mathbf{X}'\mathbf{X})^{-1}) = \sigma_0^2 (\mathbf{X}'\mathbf{X})^{-1} \quad (3.28)$$

However, note that $E((\mathbf{X}'\mathbf{X})^{-1})$ is the expected value of a complicated nonlinear function of r.v, and thus it does not equal to the nonlinear function of the expected value. That means, strictly speaking, $(\mathbf{X}'\mathbf{X})^{-1}$ is a biased estimate of $E((\mathbf{X}'\mathbf{X})^{-1})$.

Econometricians thus take alternative approaches to justify the equation in (3.28). First, they may argue that $(\mathbf{X}'\mathbf{X})^{-1}$ is a consistent estimate of $\sigma_0^2 E((\mathbf{X}'\mathbf{X})^{-1})$. Thus, $Var(\hat{\beta}) = \sigma_0^2 E((\mathbf{X}'\mathbf{X})^{-1}) = \sigma_0^2 (\mathbf{X}'\mathbf{X})^{-1}$ asymptotically.

Or they may assume that \mathbf{X} is fixed (e.g., stronger than independence between the regressor and the error terms). This also justifies (3.28), though this is a very strong assumption. Thus, strictly speaking, when we have small samples, estimation of the variance (3.28) is biased downward because we ignore the variability coming from the change in the explanatory variable observations over repeated samples – but normally we take it as a compromise.

The precision of the LS estimates

For a scalar parameter, the accuracy of an estimator is often taken to be proportional to the inverse of its variance, which is called the **precision** of the parameter. The **precision matrix** is defined as the inverse of the covariance matrix of the estimator. (3.28) can be rewritten as:

$$Var(\hat{\beta}) = \frac{\sigma_0^2}{n} \left(\frac{1}{n} \mathbf{X}'\mathbf{X} \right)^{-1}$$

Thus, we find that $Var(\hat{\beta})$ is proportional to three things.

First, $Var(\hat{\beta})$ is proportional to the error variance σ_0^2 . The more random variation in the error term, the more random variation in the parameter estimate. Also, as (3.17) suggests that, when the sample size n is reasonably large, we have:

$$\text{plim}_{n \rightarrow \infty} \frac{\mathbf{X}'\mathbf{X}}{n} = \mathbf{S}_{X,X}$$

Thus $\left(\frac{1}{n} \mathbf{X}'\mathbf{X} \right)^{-1}$ tends to $\mathbf{S}_{X,X}^{-1}$ which is a non-random vector variable and does not vary with the sample size n . Thus, we find that $Var(\hat{\beta})$ also tends to be proportional to $\frac{1}{n}$. Thus, if we double the sample size, we expect the variance of the estimate to be halved.

$Var(\hat{\beta})$ also depends on the matrix \mathbf{X} . When the explanatory variables contain a lot of information, $\mathbf{S}_{X,X}$ tends to be large (e.g., this can be intuitively and analytically shown if we look at an example where we only have one regressor, i.e., $\mathbf{S}_{X,X}$ is directly related to the variance of the regressor). Some elements in the $Var(\hat{\beta})$ matrix tend to be large because some explanatory variables contain overlapped information – which we call **multicollinearity** (see more details in Davison and Mackinnon, 2003).

We can use $Var(\hat{\beta})$ to make inference about the OLS estimator $\hat{\beta}$.

Linear function of parameter estimates

Once we know that covariance matrix of the OLS estimator, e.g., $Var(\hat{\beta})$, we can then know the variance of a linear combination of the OLS estimator, e.g.,

$$\begin{aligned}
 Var(\omega' \hat{\beta}) &= E \left(\omega' (\hat{\beta} - \beta_0) (\hat{\beta} - \beta_0)' \omega | X \right) \\
 &= \omega' E \left((\hat{\beta} - \beta_0) (\hat{\beta} - \beta_0)' | X \right) \omega \\
 &= \omega' \sigma_0^2 (X'X)^{-1} \omega \\
 &= \sigma_0^2 \omega' (X'X)^{-1} \omega
 \end{aligned} \tag{3.33}$$

We can use this finding to derive the variance of the forecast error. We denote that y_s and X_s represent the row vector for the dependent variable and the matrix for the explanatory variables for $t = s$ where s is beyond the in-sample data for $t = 1$ to n . We assume that the model is correctly specified, suppose \hat{y}_s is the forecast, therefore we have the forecast error variance as follows:

$$\begin{aligned}
 E(y_s - \hat{y}_s)^2 &= E(y_s - X_s \hat{\beta})^2 \\
 &= E(X_s \beta_0 + u_s - X_s \hat{\beta})^2 \\
 &= E(u_s + X_s \beta_0 - X_s \hat{\beta})^2 \\
 &= E(u_s^2 + (X_s \beta_0 - X_s \hat{\beta})^2 + 2u_s(X_s \beta_0 - X_s \hat{\beta}))^2 \\
 &= E(u_s^2) + E(u_s + X_s \beta_0 - X_s \hat{\beta})^2 + 2E(u_s(X_s \beta_0 - X_s \hat{\beta}))
 \end{aligned}$$

If we can assume exogeneity (e.g., no serial correlation), we have:

$$u_s X_s = 0$$

Thus, the previous equation becomes:

$$E(y_s - X_s \hat{\beta})^2 = E(u_s^2) + E(X_s \beta_0 - X_s \hat{\beta})^2$$

If we assume $\hat{\beta}$ as unbiased, we have:

$$E(y_s - X_s \hat{\beta})^2 = \sigma_0^2 + Var(X_s \hat{\beta})$$

This suggests that the forecasting error variance, when we assume the model is correctly specified and exogeneity to be hold, is the error variance σ_0^2 plus a penalty of $Var(X_s \hat{\beta})$ due to

the fact we do not use the true parameter β_0 but the estimate, e.g., $\hat{\beta}$ to generate the forecast.

Note that, since \hat{y}_s is a linear combination of X_s , and we assumed that X_s is uncorrelated with u_s . We know that \hat{y}_s is uncorrelated with u_s .

Efficiency of the OLS estimator

One estimator is said to be more efficient than another when its estimates are more accurate. That is, it uses the information available more efficiently. e.g.,

1. For scalar parameters, one estimator is more efficient than another if its precision is larger than that of the latter.
2. For a vector parameter, e.g., $\hat{\beta}$ and $\tilde{\beta}$ are both unbiased estimators of the vector parameter β . Thus, $\hat{\beta}$ is more efficient than $\tilde{\beta}$ if and only if $Var(\tilde{\beta}) - Var(\hat{\beta})$ is a nonzero positive semidefinite matrix.

The fact that $\hat{\beta}$ is more efficient than $\tilde{\beta}$ also suggests that every element in β or any combination of those parameters, is estimated at least as efficiently by using $\hat{\beta}$ as using $\tilde{\beta}$.

Proof: suppose γ is a linear combination of β , e.g., $\gamma = \omega' \beta$. Here in a special case, γ could be a particular parameter in β if the elements in ω are all zero except for that parameter in β . As we see in (3.33) that,

$$Var(\omega' \hat{\beta}) = \omega' (\sigma_0^2 X' X)^{-1} \omega = \omega' Var(\hat{\beta}) \omega$$

Accordingly,

$$Var(\omega' \tilde{\beta}) = \omega' Var(\tilde{\beta}) \omega$$

Thus, if we take the difference, we have:

$$Var(\omega' \tilde{\beta}) - Var(\omega' \hat{\beta}) = \omega' (Var(\tilde{\beta}) - Var(\hat{\beta})) \omega \tag{3.36}$$

Therefore, if $\hat{\beta}$ is more efficient than $\tilde{\beta}$, the RHS of (3.36) will be non-negative, which suggests that the LHS of (3.36) must be non-negative. This suggests that $Var(\omega' \tilde{\beta}) - Var(\omega' \hat{\beta}) \geq 0$.

Thus, any combination of parameters in β is estimated at least as efficiently by using $\hat{\beta}$ as using $\tilde{\beta}$.

The OLS estimator is more efficient than other linear unbiased estimators

We say an estimator to be a 'linear' estimator if we can write the estimator as a linear function of the vector of observations y for the dependent variable. e.g., We know that the OLS estimator is:

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

We can denote that $\tilde{\beta}$ to be a linear estimator but is NOT the OLS estimator, e.g.,

$$\tilde{\beta} = Ay = (X'X)^{-1}X'y + Cy$$

where A and C are $k \times n$ matrices that depend on X . Based on the equation above, we have:

$$Ay = (X'X)^{-1}X'y + Cy$$

$$Cy = Ay - (X'X)^{-1}X'y$$

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

The Gauss-Markov Theorem:

If we can assume that $E(u|X) = 0$ (exogeneity) and $E(uu'|X) = \sigma^2 I$ (homoskedasticity) in the linear regression model, then $\hat{\beta}$ is **BLUE** (best linear unbiased estimator). That is, the OLS estimator is more efficient than any other linear unbiased estimator. Note that the OLS estimator may not be more efficient compared to other biased estimators or nonlinear estimators and it only has the BLUE property when the assumptions (e.g., exogeneity and homoskedasticity) are met.

We can prove that if:

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

And

$$\tilde{\beta} = Ay = (X'X)^{-1}X'y + Cy$$

If we assume $\tilde{\beta}$ is unbiased, we will have a strong constraint on C , which makes $Cov(\hat{\beta}, Cy) = 0$, (see more details in Davison and Mackinnon, 2003). Thus, we have:

$$\begin{aligned}
Var(\tilde{\beta}) &= Var(\hat{\beta} + (\tilde{\beta} - \hat{\beta})) \\
&= Var(\hat{\beta} + Cy) \\
&= Var(\hat{\beta}) + Var(Cy) - Cov(\hat{\beta}, Cy) \\
&= Var(\hat{\beta}) + Var(Cy)
\end{aligned}$$

Therefore, we have:

$$Var(\tilde{\beta}) - Var(\hat{\beta}) = Var(Cy)$$

Since $Var(Cy)$ is a covariance matrix, $Var(\tilde{\beta}) - Var(\hat{\beta})$ must be positive semidefinite. This proves the **Gauss-Markov Theorem**.

Residual and the error term

Once we have the parameter estimate $\hat{\beta}$, we can calculate the estimate of error term, e.g., \hat{u} :

$$\hat{u} \equiv y - X\hat{\beta}$$

According to the numerical property of OLS, \hat{u} is orthogonal to $X\hat{\beta}$ and every vector in the $S(X)$. We need \hat{u} to estimate the error variance σ^2 so that we can calculate the variance matrix of $\hat{\beta}$:

$$Var(\hat{\beta}) = \frac{\sigma_0^2}{n} \left(\frac{1}{n} X'X \right)^{-1}$$

We know that $\hat{\beta}$ is consistent, thus, $\hat{u} \rightarrow u$ when $n \rightarrow \infty$. however, the finite property of \hat{u} differs from those of u .

If the model is correctly specified, we have:

$$\begin{aligned}
\hat{u} &= M_X y = M_X X \beta_0 + M_X u \\
\hat{u} &= M_X y = M_X u \\
\hat{u} &= M_X u
\end{aligned}$$

Each residual in \hat{u} is a linear combination of every element in u . Consider a row of the matrix product of $\hat{u} = M_X u$,

$$\begin{aligned}
\hat{u}_t &= u_t - X_t(X'X)^{-1}X'u \\
\hat{u}_t &= u_t - \sum_{s=1}^n X_t(X'X)^{-1}X'_s u_s
\end{aligned}$$

(3.42)

Thus, even we may assume that u_t is independent of all the other error terms (e.g., u_s), \hat{u}_t is not. There is some dependency between each pair of \hat{u}_t (e.g., in (3.42), all the \hat{u}_t depend on $\sum_{s=1}^n \mathbf{X}_t(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}_s'u_s$) though this dependency diminishes as sample size n increase.

If we assume exogeneity, e.g., $E(\mathbf{u}|\mathbf{X}) = 0$. This assumption suggests that $E(u_t|\mathbf{X}) = 0$. Also, we know from (3.42) that \hat{u}_t is a linear combination of all the u_t . Therefore, this means the $E(\hat{u}_t|\mathbf{X}) = 0$. In this respect, \hat{u}_t (or $\hat{\mathbf{u}}$) just behaves like u_t (or \mathbf{u}).

In other respects, \hat{u}_t (or $\hat{\mathbf{u}}$) does not behave like u_t (or \mathbf{u}). e.g., the Euclidean length of the vector of the LS residuals $\hat{\mathbf{u}}$ is always smaller than that of the vector of the residuals evaluated at any other value. This means, $\hat{\mathbf{u}}$ must be shorter than the vector of error terms $\mathbf{u} = \mathbf{u}(\beta_0)$

We can calculate the variance of $\hat{\mathbf{u}}$:

$$\begin{aligned} Var(\hat{\mathbf{u}}) &= Var(\mathbf{M}_X\mathbf{u}) = E(\mathbf{M}_X\mathbf{u}\mathbf{u}'\mathbf{M}_X) \\ &= E(\mathbf{M}_X\mathbf{u}\mathbf{u}'\mathbf{M}_X) = \mathbf{M}_X E(\mathbf{u}\mathbf{u}') \mathbf{M}_X = \mathbf{M}_X Var(\mathbf{u}) \mathbf{M}_X \\ &= \mathbf{M}_X(\sigma^2 \mathbf{I}) \mathbf{M}_X = \sigma^2 \mathbf{M}_X \mathbf{M}_X = \sigma^2 \mathbf{M}_X \end{aligned} \quad (3.43)$$

The second equation holds because the expected value of $\mathbf{M}_X\mathbf{u}$ is zero. (3.43) suggests that in general the variance-covariance of $\hat{\mathbf{u}}$ is not a diagonal matrix. This suggests that the off-diagonal elements in the matrix are not zero, e.g., $E(\hat{u}_t\hat{s}_t) \neq 0$. Thus, even we assume the original errors are uncorrelated, the residuals are not uncorrelated.

(3.43) also suggests that the residuals do not have a constant variance, and the residual variance for every \hat{u}_t must be smaller than σ_0^2 (see details in Davison and Mackinnon, 2003).

Estimating the variance of the error terms

The method of LS estimate β but not σ^2 . We can estimate σ^2 using MM based on the moment of the sample.

Ideally, we have:

$$\hat{\sigma}^2 = Var(u_t) = \frac{1}{n} \sum_{t=1}^n u_t^2$$

However, we do not observe u_t but only \hat{u}_t , thus we have:

$$\hat{\sigma}^2 = Var(\hat{u}_t) = \frac{1}{n} \sum_{t=1}^n \hat{u}_t^2 \quad (3.46)$$

It can be shown that (3.46) is consistent. However, because each \hat{u}_t^2 is smaller than σ^2 , (3.46) is biased downward. It can be shown that we can make some adjustment to obtain an unbiased estimator for σ^2 , e.g.,

$$s^2 = \frac{1}{n-k} \sum_{t=1}^n \hat{u}_t^2$$

All regression programs report s^2 . We can therefore update the estimate variance of the OLS parameter from:

$$Var(\hat{\beta}) = \frac{\sigma_0^2}{n} \left(\frac{1}{n} \mathbf{X}'\mathbf{X} \right)^{-1}$$

To:

$$\widehat{Var}(\hat{\beta}) = s^2 (\mathbf{X}'\mathbf{X})^{-1}$$

It can be proven that s^2 is the best quadratic unbiased estimator of σ^2 in the CLR model and is the best unbiased estimator of σ^2 in the CNLR model.

3.7 misspecification of linear regression models

Overspecification

Suppose that the DGP is generated by:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}_0 + \mathbf{u}, \mathbf{u} \sim IID(\mathbf{0}, \sigma_0^2 \mathbf{I}) \quad (3.52)$$

If we have a correctly specified model, e.g.,

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

We will have the OLS estimate as:

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

However, suppose we have an over-specified model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \mathbf{u}, \mathbf{u} \sim IID(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (3.51)$$

We will have the OLS estimate as:

$$\tilde{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{M}_Z\mathbf{X})^{-1}\mathbf{X}'\mathbf{M}_Z\mathbf{y}$$

Overspecification is strictly speaking not a misspecification because (3.52) is still embedded in (3.51). Indeed, if we replace \mathbf{y} by $\mathbf{X}\boldsymbol{\beta}_0 + \mathbf{u}$, we have:

$$\tilde{\boldsymbol{\beta}} = \boldsymbol{\beta}_0 + (\mathbf{X}'\mathbf{M}_Z\mathbf{X})^{-1}\mathbf{X}'\mathbf{M}_Z\mathbf{u}$$

It is obvious that both $\hat{\boldsymbol{\beta}}$ and $\tilde{\boldsymbol{\beta}}$ are unbiased under the assumption of exogeneity (or fixed regressor). However, we could find that $\hat{\boldsymbol{\beta}}$ is more efficient than $\tilde{\boldsymbol{\beta}}$.

Proof:

e.g., we have:

$$\begin{aligned} \text{Var}(\tilde{\boldsymbol{\beta}}) &= E((\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)') \\ &= \sigma_0^2 (\mathbf{X}'\mathbf{M}_Z\mathbf{X})^{-1} \end{aligned}$$

$$\text{Var}(\hat{\boldsymbol{\beta}}) = \sigma_0^2 (\mathbf{X}'\mathbf{X})^{-1}$$

Thus, we have:

$$\text{Var}(\hat{\boldsymbol{\beta}})^{-1} - \text{Var}(\tilde{\boldsymbol{\beta}})^{-1} = \frac{1}{\sigma_0^2} (\mathbf{X}'\mathbf{X} - \mathbf{X}'\mathbf{M}_Z\mathbf{X}) = \frac{1}{\sigma_0^2} (\mathbf{P}_Z\mathbf{X})'\mathbf{P}_Z\mathbf{X} \quad (3.58)$$

(3.58) is a matrix of the form $\mathbf{B}'\mathbf{B}$, it must be a positive semidefinite matrix. Thus, this is equivalent to saying that: $\text{Var}(\tilde{\boldsymbol{\beta}}) - \text{Var}(\hat{\boldsymbol{\beta}})$ is a positive semidefinite matrix.

(3.58) suggests that $\hat{\boldsymbol{\beta}}$ is more efficient than $\tilde{\boldsymbol{\beta}}$, and $\hat{\boldsymbol{\beta}}$ and $\tilde{\boldsymbol{\beta}}$ are equivalent only when $\mathbf{P}_Z\mathbf{X} = \mathbf{0}$ which means \mathbf{X} and \mathbf{Z} are mutually orthogonal. In general, including regressors that do not belong in a model increase the variance of the estimates of coefficients on the regressors that do belong, and the increase can be very great in many cases.

Under-specification

Suppose we have the following DGP:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}_0 + \mathbf{Z}\boldsymbol{\gamma}_0 + \mathbf{u}, \mathbf{u} \sim IID(\mathbf{0}, \sigma_0^2 \mathbf{I}) \quad (3.59)$$

And we have an underspecified model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \mathbf{u} \sim IID(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (3.55)$$

Under-specification is a misspecification because (3.59) is NOT embedded in (3.51). Now we need to realize that the OLS estimate based on (3.55) is biased in general. e.g.,

$$\begin{aligned} E(\hat{\boldsymbol{\beta}}) &= E((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}) \\ E(\hat{\boldsymbol{\beta}}) &= E((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\boldsymbol{\beta}_0 + \mathbf{Z}\boldsymbol{\gamma}_0 + \mathbf{u})) \\ &= \boldsymbol{\beta}_0 + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Z}\boldsymbol{\gamma}_0 \end{aligned} \quad (3.60)$$

(3.60) equals to zero only when $\mathbf{X}'\mathbf{Z} = 0$ (e.g., \mathbf{X} and \mathbf{Z} are mutually orthogonal). The magnitude of the bias depends on \mathbf{X} , \mathbf{Z} , and $\boldsymbol{\gamma}_0$. Also, because this bias does not vanish as n increase, thus $\hat{\boldsymbol{\beta}}$ is not consistent.

Since $\hat{\boldsymbol{\beta}}$ is biased, we cannot calculate the variance-covariance matrix to evaluate its accuracy. Instead, we use the mean square error (**MSE**) matrix:

$$MSE(\hat{\boldsymbol{\beta}}) = E((\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)')$$

Note that $MSE(\hat{\boldsymbol{\beta}})$ looks similar to $Var(\hat{\boldsymbol{\beta}})$ but is different because the expected value of $\hat{\boldsymbol{\beta}}$ is not $\boldsymbol{\beta}_0$ (see (3.60)). The MSE matrix is equivalent to $Var(\hat{\boldsymbol{\beta}})$ when $\hat{\boldsymbol{\beta}}$ is unbiased, though not otherwise. Since we know that:

$$\begin{aligned} \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0 &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\boldsymbol{\beta}_0 + \mathbf{Z}\boldsymbol{\gamma}_0 + \mathbf{u}) - \boldsymbol{\beta}_0 \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Z}\boldsymbol{\gamma}_0 - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u} \end{aligned}$$

Thus, we can derive the MSE matrix as:

$$\begin{aligned} MSE(\hat{\boldsymbol{\beta}}) &= E((\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)') \\ MSE(\hat{\boldsymbol{\beta}}) &= \sigma_0^2(\mathbf{X}'\mathbf{X})^{-1} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Z}\boldsymbol{\gamma}_0\boldsymbol{\gamma}_0'\mathbf{Z}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \end{aligned} \quad (3.62)$$

Suppose that $\tilde{\boldsymbol{\beta}}$ is the OLS estimate based on the correctly specified model. Note that:

$$\begin{aligned} Var(\tilde{\boldsymbol{\beta}}) &= E((\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)') \\ &= \sigma_0^2(\mathbf{X}'\mathbf{M}_Z\mathbf{X})^{-1} \end{aligned}$$

Thus, we would like to compare $MSE(\hat{\boldsymbol{\beta}})$ with $Var(\tilde{\boldsymbol{\beta}})$ (which is what it should be if the model is correctly specified). However, no unambiguous comparison is possible. It is possible that

some parameters may be estimated more efficiently by $\hat{\beta}$ and others more efficiently by $\tilde{\beta}$. Nevertheless, the covariance matrix of $\hat{\beta}$ calculated by programs will be incorrect as it only calculates the first component of (3.62).

Therefore, we conclude that under-specification will lead the OLS estimate to be **biased** and **inconsistent** and the variance-covariance matrix of the parameters to be **misleading**.