

LECTURE 15: STATISTICAL PROPERTIES OF ORDINARY LEAST SQUARES

MECO 7312.
INSTRUCTOR: DR. KHAI CHIONG

1. OLS covariance matrix

Recall the OLS estimator $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$.

What is the sampling distribution of $\hat{\beta}$? We must start from the data-generating process. The sampling distribution is induced by the randomness in $\mathbf{y} = \mathbf{X}\beta_0 + \epsilon$, where either \mathbf{X} is nonstochastic (fixed) and ϵ is random, or (\mathbf{X}, ϵ) are both random. In scientific experiments, we can treat \mathbf{X} as being fixed. In either case, all OLS-related derivations are similar.

Previously, we have seen that under the *exogeneity* condition $\mathbb{E}[\epsilon|\mathbf{X}] = \mathbf{0}$, the OLS estimator is unbiased, i.e. $\mathbb{E}[\hat{\beta}] = \beta_0$.

Now we want to examine other features of the sampling distribution, such as the precision of the estimator – what is $\text{Var}(\hat{\beta})$? When we take the variance of a k -dimensional vector, we mean the k -by- k variance covariance matrix.

Recall that

$$(1) \quad \hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\beta_0 + \epsilon)$$

$$(2) \quad = \beta_0 + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

$$(3) \quad \hat{\beta} - \beta_0 = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

The variance-covariance matrix is given by:

$$(4) \quad \mathbb{E}[(\hat{\beta} - \beta_0)(\hat{\beta} - \beta_0)^T] = \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon^T]$$

$$(5) \quad = \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}]$$

$$(6) \quad = \mathbb{E}[\mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} | \mathbf{X}]]$$

$$(7) \quad = \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\epsilon \epsilon^T | \mathbf{X}] \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}]$$

Now we need to assume something about $\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T|\mathbf{X}]$, in particular, we assume that

$$(8) \quad \mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T|\mathbf{X}] = \sigma_0^2\mathbf{I}$$

This assumption means: $\text{Var}[\epsilon_i] = \sigma_0^2$ for all $i = 1, \dots, n$, and that $\mathbb{E}[\epsilon_i\epsilon_j] = 0$ for all $i \neq j$. In words, the error term across all observations have the same variance σ_0^2 , and the covariance of the error term across different observations is zero. That the observations are i.i.d. would imply $\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T|\mathbf{X}] = \sigma_0^2\mathbf{I}$, but i.i.d is a stronger requirement.

When the error terms have identical variance across observations, we are said to be imposing the *homoskedasticity* assumption (as opposed to heteroskedasticity).

$$(9) \quad \mathbb{E}[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^T] = \mathbb{E}[(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\sigma_0^2\mathbf{I}\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}]$$

$$(10) \quad = \sigma_0^2\mathbb{E}[(\mathbf{X}^T\mathbf{X})^{-1}]$$

In the case where \mathbf{X} is non-stochastic, then $\text{Var}(\hat{\boldsymbol{\beta}}) = \sigma_0^2(\mathbf{X}^T\mathbf{X})^{-1}$. When \mathbf{X} is stochastic, we simply estimate $\mathbb{E}[(\mathbf{X}^T\mathbf{X})^{-1}]$ as $(\mathbf{X}^T\mathbf{X})^{-1}$. Alternatively, we are not interested in the stochastic process governing \mathbf{X} , and so we calculate the variance-covariance matrix conditioning on \mathbf{X} . Therefore,

$$(11) \quad \text{Var}(\hat{\boldsymbol{\beta}}|\mathbf{X}) = \sigma_0^2(\mathbf{X}^T\mathbf{X})^{-1}$$

The precision of the OLS estimator is defined as the inverse of $\text{Var}(\hat{\boldsymbol{\beta}})$.

Check using simulation that the variance of OLS estimators decrease (OLS estimators become more precise) when: (i) sample size n increases, and (ii) when the collinearity between regressors decreases.

We can demonstrate this more rigorously using the FWL theorem.

Consider the regression model $\mathbf{y} = \mathbf{x}_1\beta_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u}$. By the FWL theorem, the OLS estimator of β_1 in this regression model is equivalent to the corresponding OLS estimator in $\mathbf{M}_2\mathbf{y} = \mathbf{M}_2\mathbf{x}_1\beta_1 + \boldsymbol{\epsilon}$, where $\mathbf{M}_2 = \mathbf{I} - \mathbf{P}_2$, and $\mathbf{P}_2 = \mathbf{X}_2(\mathbf{X}_2^T\mathbf{X}_2)^{-1}\mathbf{X}_2^T$.

Therefore,

$$(12) \quad \text{Var}(\hat{\beta}_1) = \sigma_0^2((\mathbf{M}_2\mathbf{x}_1)^T(\mathbf{M}_2\mathbf{x}_1))^{-1}$$

$$(13) \quad = \frac{\sigma_0^2}{(\mathbf{M}_2\mathbf{x}_1)^T(\mathbf{M}_2\mathbf{x}_1)}$$

Therefore, $\hat{\beta}_1$ becomes more precise when the squared Euclidean length of the vector $\mathbf{M}_2\mathbf{x}_1$ is large.

The squared Euclidean length of the vector $\mathbf{M}_2\mathbf{x}_1$ is just the sum of squared residuals from the regression:

$$(14) \quad \mathbf{x}_1 = \mathbf{X}_2\boldsymbol{\alpha} + \text{residuals}$$

Thus, when \mathbf{x}_1 can be explained by \mathbf{X}_2 , $(\mathbf{M}_2\mathbf{x}_1)^T(\mathbf{M}_2\mathbf{x}_1)$ becomes small, and consequently, $\hat{\beta}_1$ becomes less precise. The intuition is that, we cannot estimate the effect of a regressor on the dependent variable well if that regressor can be explained by the other regressors – this regressor does not add new orthogonal information. For instance, if the dependent variable is sales across customers, and the explanatory variables are demographics regressors, we may find that it is very hard to disentangle the effects between two highly correlated regressors, say income and education level.

2. Estimating the variance of the error terms

How do we estimate σ_0^2 ? Using the method of moments, and the moment condition that $\mathbb{E}[\epsilon_i^2] = \sigma_0^2$, one estimator is $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2$, where $\hat{\epsilon} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ is the fitted residual.

It can be shown that $\hat{\sigma}^2$ is a consistent estimator of σ_0 , however $\hat{\sigma}^2$ is a biased estimator. In particular, $\mathbb{E}[\hat{\sigma}^2] = \frac{n-k}{n}\sigma_0^2$, and therefore, an unbiased estimator of σ_0 would be:

$$(15) \quad s^2 = \frac{1}{n-k} \sum_{i=1}^n \hat{\epsilon}_i^2$$

To see this, note that:

$$(16) \quad \hat{\epsilon} = \mathbf{M}\mathbf{y}$$

$$(17) \quad = \mathbf{M}(\mathbf{X}\boldsymbol{\beta}_0 + \boldsymbol{\epsilon})$$

$$(18) \quad = \mathbf{M}\boldsymbol{\epsilon}$$

Now consider the estimator $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2 = \frac{1}{n} \text{Trace}(\hat{\boldsymbol{\epsilon}}\hat{\boldsymbol{\epsilon}}^T)$, where the Trace of a matrix is just the sum of the diagonal elements. Now we want to compute $\mathbb{E}[\hat{\sigma}^2]$. Assume that \mathbf{X} is non-stochastic (or we implicitly condition on \mathbf{X} , i.e. $\mathbb{E}[\hat{\sigma}^2|\mathbf{X}]$).

$$\begin{aligned}
(19) \quad & \mathbb{E}[\text{Trace}(\hat{\boldsymbol{\epsilon}}\hat{\boldsymbol{\epsilon}}^T)] = \mathbb{E}[\text{Trace}(\mathbf{M}\boldsymbol{\epsilon}(\mathbf{M}\boldsymbol{\epsilon})^T)] \\
(20) \quad & = \text{Trace}(\mathbb{E}[\mathbf{M}\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T\mathbf{M}]) \\
(21) \quad & = \text{Trace}(\mathbf{M}\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T]\mathbf{M}) \\
(22) \quad & = \sigma_0^2\text{Trace}(\mathbf{M}\mathbf{M}) \\
(23) \quad & = \sigma_0^2\text{Trace}(\mathbf{M})
\end{aligned}$$

Now:

$$\begin{aligned}
(24) \quad & \text{Trace}(\mathbf{M}) = \text{Trace}(\mathbf{I} - \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T) \\
(25) \quad & = \text{Trace}(\mathbf{I}) - \text{Trace}(\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T) \\
(26) \quad & = n - \text{Trace}(\mathbf{X}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}) \\
(27) \quad & = n - k
\end{aligned}$$

Therefore $\mathbb{E}[\hat{\sigma}^2] = \mathbb{E}[\frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2] = \frac{n-k}{n} \sigma_0^2$. An unbiased estimator of σ_0^2 is:

$$(28) \quad s^2 = \frac{1}{n-k} \sum_{i=1}^n \hat{\epsilon}_i^2$$

To see the intuition behind the factor $\frac{1}{n-k}$. Suppose there are as many regressors as there are observations, then the residuals are always zero, and we would not have obtained any information on σ_0^2 . Of course this is an extreme case. In practice we confine ourselves to the case $k < n$. The very fact that we choose $\hat{\boldsymbol{\beta}}$ in such a way that the sum of squared residuals is minimized is the cause of the fact that the squared residuals $\hat{\boldsymbol{\epsilon}}$ are smaller (on average) than the squared error terms $\boldsymbol{\epsilon}$.

2.1. Leverage

If we denote the i -th diagonal of $\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ as h_i , then h_i is also known as the *leverage* of the observation i .

$$\begin{aligned}
(29) \quad & \text{Var}(\hat{\epsilon}_i) = \mathbb{E}[\hat{\epsilon}_i^2] \\
(30) \quad & = \mathbb{E}[(\hat{\boldsymbol{\epsilon}}\hat{\boldsymbol{\epsilon}}^T)_{ii}] \\
(31) \quad & = (\mathbf{M}\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T]\mathbf{M})_{ii} \\
(32) \quad & = \sigma_0^2\mathbf{M}_{ii} \\
(33) \quad & = \sigma_0^2(1 - h_i)
\end{aligned}$$

Therefore, when an observation has a high leverage, $1 - h_i$ is small, and the residual $\hat{\epsilon}_i$ is close to zero. Even though it must be the case that $\sum_i \hat{\epsilon}_i = 0$, and $\mathbb{E}[\hat{\epsilon}_i] = 0$, not all residuals are equally small – some observations have smaller residuals – OLS prioritize these observations (those with high leverage). By the way, can you show that $\sum_i \hat{\epsilon}_i = 0$, and $\mathbb{E}[\hat{\epsilon}_i] = 0$?

Leverage is a measure of how influential an observation is – how much the OLS estimate changes when the point is removed. Data points with high leverage or influence force the regression line to be close to the point. Leverage is not to be confused with the notion of outliers.

3. Heteroskedasticity-consistent covariance matrix estimator

Recall that the variance-covariance matrix of the OLS estimator $\hat{\beta}$ is:

$$(34) \quad \mathbb{E}[(\hat{\beta} - \beta_0)(\hat{\beta} - \beta_0)^T] = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T] \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$

(Either treat \mathbf{X} to be fixed, or implicitly condition on \mathbf{X}).

Heteroskedasticity means:

$$\mathbb{E}[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T] = \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{bmatrix}$$

We have reasons to believe that some observations have more noise than others – $\text{Var}(\epsilon_i)$ differs across i .

Heteroskedasticity does **not** cause OLS to be biased, but the estimator for the variance-covariance matrix of OLS would be biased and wrong. Therefore, we still get the same estimate regardless of whether we assume heteroskedasticity or not, but our inference (hypothesis test, confidence interval, etc) would be wrong.

In one of the most cited paper of all time in Economics, Halber White (1980) introduced the following estimator:

$$(35) \quad (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \hat{\Sigma} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$

Where:

$$\hat{\Sigma} = \begin{bmatrix} \hat{\epsilon}_1^2 & 0 & \cdots & 0 \\ 0 & \hat{\epsilon}_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \hat{\epsilon}_n^2 \end{bmatrix}$$

White shows that this is a consistent estimator for the variance-covariance matrix of OLS estimators, i.e. $\mathbb{E}[(\hat{\beta} - \beta_0)(\hat{\beta} - \beta_0)^T]$. This is the estimator that Stata implements by default with the command `reg y x, robust`.

The intuition behind this estimator can be seen by recalling Equation 33: $\text{Var}(\epsilon_i) = \mathbb{E}[\hat{\epsilon}_i^2] = \sigma_0^2(1 - h_i)$. Asymptotically when n is large, we can ignore the leverage factor. There are other heteroskedastic-consistent variance-covariance estimator, such as using:

$$\hat{\Sigma} = \begin{bmatrix} \hat{\epsilon}_1^2(1 - h_1) & 0 & \cdots & 0 \\ 0 & \hat{\epsilon}_2^2(1 - h_2) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \hat{\epsilon}_n^2(1 - h_n) \end{bmatrix}$$

White's estimator is the standard estimator. None of these estimators have any finite-sample guarantee (unbiasedness), but asymptotically they are all equivalent and consistent. Therefore, there is no formal mathematical reason to think that the latter estimator is better.

3.1. Serial correlation

In the presence of serial correlation,

$$(36) \quad \mathbb{E}[\epsilon\epsilon^T] = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1n} \\ \sigma_{12} & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{bmatrix}$$

The off-diagonals are non-zero. The error terms are correlated across observations. This is quite common in time-series (but not in cross-sectional data). In Equation 36 above, we have both serial correlation and heteroskedasticity.

Serial correlation does **not** affect the unbiasedness of OLS estimators. Similar to heteroskedasticity, serial correlation results in incorrect confidence intervals and hypothesis tests.

Serial correlation is usually corrected by assuming that the serial correlation follows a specific form: $\epsilon_t = \rho\epsilon_{t-1} + u_t$. This is known as the AutoRegressive(1) errors. We can test for the presence of this kind of serial correlation using the Durbin-Watson test. Correcting for serial correlation involves differencing: we regress the difference $y_t - \rho y_{t-1}$ on the difference $x_t - \rho x_{t-1}$. The parameter ρ can be estimated consistently using OLS residuals, by regressing $\hat{\epsilon}_t$ on $\hat{\epsilon}_{t-1}$.

4. Hypothesis testing and confidence interval involving OLS estimators

Suppose that $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}_0 + \mathbf{u}$. Assume the following: (1) Exogeneity, $\mathbb{E}[\mathbf{u}|\mathbf{X}] = 0$, (2) No perfect multicollinearity, $(\mathbf{X}^T \mathbf{X})^{-1}$ exists, (3) Homoskedastic and no serial correlation, $\mathbb{E}[\mathbf{u}\mathbf{u}^T] = \sigma_0^2 \mathbf{I}$, (4) Normality, $\mathbf{u}|\mathbf{X} \sim \mathcal{N}(0, \sigma_0^2 \mathbf{I})$.

These assumptions are collectively called the Classical Linear Regression Model.

Then the OLS estimator $\hat{\boldsymbol{\beta}}$ satisfies:

$$(37) \quad \hat{\boldsymbol{\beta}}|\mathbf{X} \sim \mathcal{N}(\boldsymbol{\beta}_0, \sigma_0^2 (\mathbf{X}^T \mathbf{X})^{-1})$$

This result can be derived by using the fact that a linear combination of Normal random variables is a Normal random variable, and that OLS takes linear combination of the Normal error terms. Specifically if $\mathbf{u} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then $A + B\mathbf{u} \sim \mathcal{N}(A + B\boldsymbol{\mu}, B\boldsymbol{\Sigma}B^T)$.

Therefore, to construct hypothesis tests and confidence intervals for OLS estimates, we can simply applied what we have learned in the last few classes here.

If we are unwilling to assume Normal error terms, then there are two alternative approaches: (1) bootstrapping, (2) asymptotics.

Asymptotic sampling distribution. Under certain condition,

$$(38) \quad \sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) \xrightarrow{d} \mathcal{N}\left(\mathbf{0}, \sigma_0^2 \left(\frac{\mathbf{X}^T \mathbf{X}}{n}\right)^{-1}\right)$$

As such, the sampling distribution of $\hat{\boldsymbol{\beta}}$ for large n can be approximated as:

$$(39) \quad \hat{\boldsymbol{\beta}} \sim \mathcal{N}(\boldsymbol{\beta}_0, \sigma_0^2 (\mathbf{X}^T \mathbf{X})^{-1})$$

5. Summary

- (i) Exogeneity alone guarantees unbiasedness. Exogeneity can be violated under many circumstances – whenever the regressor is correlated with the error term.

For instance, if we regress weekly demand on weekly price, we expect high price to cause lower demand, so the coefficient is negative. But often this regression gives us positive price coefficient! This is because price is an endogenous variable. If price is set randomly (experimentation or A/B testing), then it would be exogenous. But price is set strategically by firms. When an unexpected shock drives demand higher, the firm responds by changing price, therefore, there is a correlation between price and the error term (demand shocks).

This is the foremost concern in any empirical research. The branch of statistics/econometrics dealing with this concern is called *causal inference*. Tools that fall under causal inference include (1) Instrumental Variable approach, (2) Difference-in-difference, (3) Regression discontinuity, (4) Propensity score matching, (5) experimentation and A/B testing.

- (ii) Heteroskedasticity and serial correlation causes incorrect statistical inference (wrong formula for calculating the variance-covariance matrix of OLS estimator).
- (iii) Multicollinearity increases the variance of OLS.
- (iv) Under-specification (omission of relevant variables) causes bias since the exogeneity condition is violated. However over-specification (inclusion of irrelevant variables) does *not* cause bias. However it does increase the variance of OLS. Specifically, over-specification means the true data-generating process is $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$, but we estimate $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\alpha} + \boldsymbol{\epsilon}$. It is straightforward to see why OLS is still unbiased.