

LECTURE 13: STATISTICAL PROPERTIES OF ORDINARY LEAST SQUARES

MECO 7312.
INSTRUCTOR: DR. KHAI CHIONG
NOVEMBER 13, 2019

1. Linear regression models

Let (Y, X, ϵ) be random variables such that:

$$(1) \quad Y = a + bX + \epsilon$$

a and b are unknown parameters, where $\mathbb{E}[\epsilon|X] = 0$. Show that $\mathbb{E}[\epsilon|X] = 0$ implies the following: (i) $\mathbb{E}[X\epsilon] = 0$, (ii) $\text{Cov}(X, \epsilon) = 0$, and (iii) $\mathbb{E}[\epsilon] = 0$.

Suppose n i.i.d random samples: (y_i, x_i, ϵ_i) for $i = 1, \dots, n$ are drawn from the data-generating model, but we only observe $(y_i, x_i)_{i=1}^n$ as our dataset.

Three ways of estimating a and b , all leading to the same estimators!

Method of moments.

$$(2) \quad \mathbb{E}[\epsilon X] = 0$$

$$(3) \quad \mathbb{E}[XY - aX - bX^2] = 0$$

$$(4) \quad \mathbb{E}[\epsilon] = 0$$

$$(5) \quad \mathbb{E}[Y - a - bX] = 0$$

Or we can use Maximum Likelihood Estimator, but we have to additionally assume that $\epsilon \sim \mathcal{N}(0, \sigma^2)$, or equivalently, $Y - a - bX \sim \mathcal{N}(0, \sigma^2)$. Therefore,

$$(6) \quad L(x_1, y_1, \dots, x_n, y_n | a, b, \sigma) = \prod_{i=1}^n \phi\left(\frac{y_i - a - bx_i}{\sigma}\right)$$

$$(7) \quad \underset{a, b, \sigma}{\operatorname{argmax}} \sum_{i=1}^n \log \phi\left(\frac{y_i - a - bx_i}{\sigma}\right)$$

Where ϕ is the pdf of the standard Normal.

Or we can minimize the sum of squared errors using calculus: $\min_{a,b} \sum_{i=1}^n (y_i - a - bx_i)^2$.

1.1. Multivariate linear regression

Now consider:

$$(8) \quad Y = \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_K X_K + \epsilon$$

Suppose we observe n i.i.d random samples: $(y_i, x_{i1}, x_{i2}, \dots, x_{iK})$ for $i = 1, \dots, n$.

$$(9) \quad y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{iK} + \epsilon_i$$

We can manipulate this equation using Matrix Algebra.

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$\boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_K]$$

Where \mathbf{X}_k is a $n \times 1$ column vector containing the k -th explanatory variable. Other names for explanatory variable: features (used by computer scientists), covariates, regressors (used by economists).

$$\mathbf{X}_k = \begin{bmatrix} x_{1k} \\ \vdots \\ x_{nk} \end{bmatrix}$$

$(x_{1k}, x_{2k}, \dots, x_{ik}, \dots, x_{nk})$ are called *observations* for the k -th covariate.

$$\mathbf{X} = \begin{bmatrix} x_{11} & \cdots & x_{1K} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nK} \end{bmatrix}$$

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_K \end{bmatrix}$$

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

$\mathbf{X}\boldsymbol{\beta}$ is the matrix product of a $n \times k$ matrix with a $k \times 1$ matrix, resulting in a $n \times 1$ matrix.

Our entire dataset are contained in the data matrix $[\mathbf{y}, \mathbf{X}]$.

2. Ordinary Least Squares (OLS) estimator

How do we estimate the $\boldsymbol{\beta}$? If we were to use Method of Moments, we need at least K number of moments conditions.

The assumption we need is that \mathbf{X} is *exogenous*, also known as the *conditional mean independence assumption*: $\mathbb{E}[\epsilon|X_1] = 0$, $\mathbb{E}[\epsilon|X_2] = 0$, \dots , $\mathbb{E}[\epsilon|X_K] = 0$. The error term is (conditionally mean) independent of each of the K explanatory variable.

The sample moment conditions can be written as: $\sum_{i=1}^n x_{i1}\epsilon_i = 0$, $\sum_{i=1}^n x_{i2}\epsilon_i = 0$, \dots , $\sum_{i=1}^n x_{ik}\epsilon_i = 0$. In matrix notation:

$$(10) \quad \mathbf{X}_1^T \boldsymbol{\epsilon} = 0$$

$$\vdots$$

$$(11) \quad \mathbf{X}_k^T \boldsymbol{\epsilon} = 0$$

Where \mathbf{X}_k is the k -th column of the data matrix \mathbf{X} , it contains all n observations for the covariate X_k . Now \mathbf{X}^T is the matrix transpose of \mathbf{X} , therefore, \mathbf{X}^T is a $1 \times n$ row vector. $\mathbf{X}^T = [x_{1k}, x_{2k}, \dots, x_{nk}]$.

Finally, the moment conditions can be summarized as just:

$$(12) \quad \mathbf{X}^T \boldsymbol{\epsilon} = \mathbf{0}_K$$

Where $\mathbf{0}_K$ is a $K \times 1$ vector of zeros. \mathbf{X}^T is a $K \times n$ matrix, while $\boldsymbol{\epsilon}$ is a $n \times 1$ matrix, therefore their matrix product has dimension $K \times 1$.

Now, we can derive the OLS (Ordinary Least Square) estimators:

$$(13) \quad \mathbf{X}^T \boldsymbol{\epsilon} = \mathbf{0}$$

$$(14) \quad \mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{0}$$

$$(15) \quad \mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X}\boldsymbol{\beta} = \mathbf{0}$$

$$(16) \quad (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}\boldsymbol{\beta} = \mathbf{0}$$

$$(17) \quad (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} - \boldsymbol{\beta} = \mathbf{0}$$

$$(18) \quad \boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Note that the right-hand side of $\boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ consists entirely of the components of the data matrix. Therefore this is a valid estimator.

2.1. OLS simulation

It is very instructive to implement OLS estimators in a programming language of your choice. For this section, please refer to the R Markdown (lecture13.rmd).

Let the (true) data-generating process be $Y_i = 2 - 3X_{i1} + 0.5X_{i2} + \epsilon_i$ for $i = 1, \dots, 1000$, where $\epsilon_i \sim \text{i.i.d } \mathcal{N}(0, 2)$, $x_{i1} \sim \text{i.i.d Exponential}(0.5)$, and $x_{i2} \sim \text{i.i.d } \mathcal{N}(-1, 1)$. The true coefficient/parameters are therefore $\boldsymbol{\beta} = [2, -3, 0.5]^T$. Since y_i is related to the other variables, we generate y_i through $y_i = 2 - 3x_{i1} + 0.5x_{i2} + \epsilon_i$. After we generated these numbers, we then ignore ϵ_i and stack them according to the data matrix $[\mathbf{y}, \mathbf{X}]$. We compute the OLS estimator $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$, and compare it to the true value.

Are the covariates \mathbf{X} exogenous here? We now consider another data-generating process where the exogeneity assumption is violated.

Suppose we only regress y_i on 1 and x_{i1} , what happens to the OLS estimates? In another words, we now try to estimate the following *misspecified* model: $Y_i = \beta_1 + \beta_2 X_{i1} + v_i$, while the true model is $Y_i = \beta_1 + \beta_2 X_{i1} + \beta_3 X_{i2} + \epsilon_i$. As such, $v_i = \beta_3 X_{i2} + \epsilon_i$.

In order to correctly estimate the coefficient β_2 from $Y_i = \beta_1 + \beta_2 X_{i1} + v_i$, OLS estimator relies on the assumption that $\mathbb{E}[v_i | X_{i1}] = 0$. But because $v_i = \beta_3 X_{i2} + \epsilon_i$, we have: $\mathbb{E}[v_i | X_{i1}] = \mathbb{E}[\beta_3 X_{i2} + \epsilon_i | X_{i1}] = \beta_3 \mathbb{E}[X_{i2} | X_{i1}]$. Therefore, we can still get consistent estimate of β_2 if $\mathbb{E}[X_{i2} | X_{i1}] = 0$, e.g. when the two covariates

are independent. We verify through the simulation that the OLS estimates are inconsistent when $(X_{1i}, X_{2i}) \sim$ i.i.d multivariate Normal with a covariance matrix of $\begin{bmatrix} 1 & 1.5 \\ 1.5 & 3 \end{bmatrix}$.

How do we estimate the variance of ϵ_i ? Sum of squares residuals.

$$(19) \quad \mathbb{E}[\epsilon^2] \approx \frac{1}{n}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

3. Multicollinearity

The OLS estimator is $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$.

The matrix $(\mathbf{X}^T \mathbf{X})$ needs to be invertible. A square matrix that is not invertible is called singular. A square matrix is singular if and only if its determinant is 0.

The matrix $(\mathbf{X}^T \mathbf{X})$ has an inverse if the columns of \mathbf{X} are linearly independent. ($(\mathbf{X}^T \mathbf{X})$ has a full column rank).

Suppose a particular column of \mathbf{X} can be written as a linear function of some other columns of \mathbf{X} (for example, $\mathbf{X}_k = \lambda_1 \mathbf{X}_1 + \lambda_2 \mathbf{X}_2$), then we say that there is a **perfect multicollinearity**. The regressors are linearly dependent. $(\mathbf{X}^T \mathbf{X})$ does not have an inverse – OLS estimator is ill-defined. When one of the regressors are too similar to another regressor, we cannot separately identify their respective coefficients.

In general, even when there is no exact linear relationship between the regressors, OLS estimator will run into problem when one of the regressors are highly correlated with another regressor. This is the **multicollinearity** problem. The inverse $(\mathbf{X}^T \mathbf{X})$ is *almost singular*. Computation of the inverse of an almost singular matrix is highly unstable and numerically imprecise.

Consider the simulation exercise before. Let the (true) data-generating process be $Y_i = 2 - 4X_{i1} + 0.5X_{i2} + \epsilon_i$ for $i = 1, \dots, 1000$, where $\epsilon_i \sim$ i.i.d $\mathcal{N}(0, 2)$, $X_{i1} \sim$ i.i.d Exponential(0.5), and $X_{i2} = 5 - 2X_{i1}$.

Now let $X_{i2} = 5 - 2X_{i1} + v_i$, where $v_i \sim \mathcal{N}(0, 0.1)$.

Multicollinearity can be detected by calculating the condition number of the matrix $(\mathbf{X}^T \mathbf{X})$. When the condition number is high, the matrix is ill-conditioned and almost singular.¹

¹The condition number is computed by finding the square root of the maximum eigenvalue divided by the minimum eigenvalue of the matrix. If the condition number is above 30, the

4. Unbiasedness of OLS estimators

What does unbiasedness mean here? Recall the simulation exercise before – we get different OLS estimates in different simulation when we draw a different random sample from the DGP. What is the average of those OLS estimates over infinitely many simulations?

Let the true data-generating process be: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}_0 + \boldsymbol{\epsilon}$. Here, we do not have to specify the the distribution of $\boldsymbol{\epsilon}$ or \mathbf{X} . Let $\hat{\boldsymbol{\beta}}$ be the OLS estimator.

$$\begin{aligned}
 (20) \quad \mathbb{E}[\hat{\boldsymbol{\beta}}] &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}] \\
 (21) \quad &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\boldsymbol{\beta}_0 + \boldsymbol{\epsilon})] \\
 (22) \quad &= \mathbb{E}[\boldsymbol{\beta}_0 + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}] \\
 (23) \quad &= \boldsymbol{\beta}_0 + \mathbb{E}[\mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon} | \mathbf{X}]] \\
 (24) \quad &= \boldsymbol{\beta}_0 + \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\boldsymbol{\epsilon} | \mathbf{X}]]
 \end{aligned}$$

Expectations above are taken row-wise (for every observation $i = 1, \dots, n$). The Law of Iterated Expectation is applied in the last two equations. It is clear that a sufficient condition for the unbiasedness of OLS estimator is that $\mathbb{E}[\boldsymbol{\epsilon} | \mathbf{X}] = \mathbf{0}$. This expression means that $\mathbb{E}[\epsilon_i | \mathbf{X}] = 0$ for all $i = 1, \dots, n$.

There are two possible ways to satisfy $\mathbb{E}[\epsilon_i | \mathbf{X}] = 0$. EITHER, $\mathbb{E}[\epsilon_i | \mathbf{X}_i] = 0$ and $(\epsilon_i, \mathbf{X}_i)$ are i.i.d across i from some probability distributions. OR, $(\epsilon_i, \mathbf{X}_i)$ are *not necessarily* i.i.d across i , but $\mathbb{E}[\epsilon_i | \mathbf{X}_1, \dots, \mathbf{X}_n] = 0$ for each i . Therefore, in the context of time-series where i.i.d does not hold true, OLS can still be unbiased. Either of these are known as the *exogeneity* condition.

5. OLS covariance matrix

Having established that OLS is unbiased under the *exogeneity* condition, what is $\text{Var}(\hat{\boldsymbol{\beta}})$? When we take the variance of a k -dimensional vector, we mean the k -by- k variance covariance matrix.

Recall that

regression may have significant multicollinearity. The condition number of a matrix indicates the potential sensitivity of the computed inverse to small changes in the original matrix.

$$(25) \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \boldsymbol{\beta}_0 + \boldsymbol{\epsilon})$$

$$(26) \quad = \boldsymbol{\beta}_0 + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}$$

$$(27) \quad \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0 = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}$$

The variance-covariance matrix is given by:

$$(28) \quad \mathbb{E}[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^T] = \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon} ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon})^T]$$

$$(29) \quad = \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}]$$

$$(30) \quad = \mathbb{E}[\mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} | \mathbf{X}]]]$$

$$(31) \quad = \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T | \mathbf{X}] \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}]$$

$$(32) \quad = \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}]$$

$$(33) \quad = \sigma^2 \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1}]$$

An estimate of the variance-covariance matrix is then

$$(34) \quad \text{Var}(\hat{\boldsymbol{\beta}}) = \hat{\sigma}^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

The precision of the OLS estimator is defined as the inverse of the variance-covariance matrix. We see that the precision of the OLS estimator decreases as σ^2 increases.

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