

Review of Ordinary Least Square

February 4, 2020

Teaching Assistant: Jeffrey Je-uei Kuo Email: jeffkuo@gwu.edu

1. The Causality, Selection Bias, and Randomized Control Trial (RCT)

The example of the Randomized Control Trial is excerpted from Chapter 1 to 3 from *Mostly Harmless Econometrics*, authored by J. Angrist, and J-S. Pischke.

In social sciences, the main goal of the research is to identify the causality between the events. However, due to the essence of social science, it is hard (or even impossible) to conduct experiments of interests by running the standard laboratory trial. For example, if we would like to see whether a new medicine is useful, we can not have the same patient received two different trials. Once the patients receive a specific medical treatment, we are not able to reverse the process and rerun the other treatment on the same person. In some cases, it might be possible to do it without involving ethical problems, but it usually costs a lot.

Hence, most of the economic data we have seen are not laboratory data. Instead, the data represents the observed outcomes after the treatment, i.g. The income survey after one enrolled in a particular training program or the income survey after the students have finished their degree.

Here is the example that Angrist and Pischke use in the text. There is an indicator representing the health status of the patients. One is the worst and five is the best. The researcher would like to know if the admission in the hospital ("hospitalized") improves the health status.

The following table shows the result after surveying the people.

Group	Sample Size	Mean Health Status	Std. Error
Hospital	7,774	3.21	0.014
No Hospital	90,049	3.93	0.003

Does "Being Hospitalized" jeopardize the health?

$$E(Y_i | D_i=1) - E(Y_i | D_i=0) = 3.21 - 3.93 = -0.72$$

Assume there are two potential outcomes for the health status for each individual i , Y_{1i} is the health status if they had been hospitalized ($D_i = 1$), Y_{0i} is the health status if they had not been hospitalized ($D_i = 0$).

$$\text{Potential Outcome for } i = \begin{cases} Y_{1i} & \text{if } D_i = 1 \\ Y_{0i} & \text{if } D_i = 0 \end{cases}$$

We have seen is so-called the Average Treatment Effect, which is $E(Y_{1i} | D_i = 1) - E(Y_{0i} | D_i = 0)$
 However, the term capturing the causal effect of the hospitalization should be Average Treatment on the Treated, which is $E(Y_{1i} | D_i = 1) - E(Y_{0i} | D_i = 1)$, or $E(Y_{1i} - Y_{0i} | D_i = 1)$
 What are the relation between the two terms?

Observed Treatment Effects (or Average Treatment Effect).

$$= E(Y_i | D_i=1) - E(Y_i | D_i=0)$$

$$= E(Y_{1i} | D_i=1) - E(Y_{0i} | D_i=0) \leftarrow \text{because treatment depends on outcome.}$$

$$= \underbrace{E(Y_{1i} | D_i=1) - E(Y_{0i} | D_i=1)}_{\text{Average Treatment Effect on the Treated Causality!}} + \underbrace{E(Y_{0i} | D_i=1) - E(Y_{0i} | D_i=0)}_{\text{selection Bias}} \leftarrow \text{mathematical manipulation.}$$

* Average Treatment Effect on the Treated Causality!

in the case, this term is negative.

means, people get treated have worse health status than people do NOT get treated.

So if treatment is independent with outcome, randomly admit patients

$$E(Y_i | D_i=1) - E(Y_i | D_i=0) = E(Y_{1i} | D_i=1) - E(Y_{0i} | D_i=0) \xrightarrow{\text{RCT}} E(Y_{1i} | D_i=1) - E(Y_{0i} | D_i=1)$$

2. Review of Simple Linear Regression

2.1 Definition

“Regression” is a functional relationship between two or more correlated variables that is often empirically determined from data and is used to predict values of one variable when given values of the others. Two words can be used to describe regression analysis: “empirical” and “informative.”

2.2 Simple Linear Regression

There is only one independent variable / regressors/ predictor variable X in the model. The proposed (population) regression model is

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

where β_0 is the intercept coefficient, β_1 is the slope coefficient, and ϵ_i is the model random error term, which takes into account all unpredictable and unknown factors that are not included in the model.

2.3 Estimation with Ordinary Least Square (OLS)

Firstly, let's define the model residual $e_i = y_i - \hat{y}_i$, where $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ is the predicted value of y_i for a given value of x_i . The OLS chooses the prediction equation that minimizes the residual sum of squares for all sample observations, that is:

$$\sum_{i=1}^n \hat{e}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

The goal of OLS is to find the estimator b_0, b_1 , such that minimize the objective function.

$$b_0, b_1 = \arg \min S = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

$$\rightarrow \sum y_i - n\hat{\beta}_0 - \hat{\beta}_1 \sum x_i = 0 \Rightarrow n\bar{y} - n\hat{\beta}_0 - \hat{\beta}_1 n\bar{x} = 0$$

$$\text{FOC } \frac{\partial S}{\partial \hat{\beta}_0} = 0 \Rightarrow 2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) (-1) = 0 \Rightarrow \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

$$\frac{\partial S}{\partial \hat{\beta}_1} = 0 \Rightarrow \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) (-x_i) = 0$$

3. The Role of R^2 . The Goodness of Fit of the Model.

$$\text{SST: Sum of Squares Total} = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$\text{SSR: Sum of Squared Regression} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$\text{SSE: Sum of Squared Errors} = \sum_{i=1}^n \hat{e}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 = \sum_{i=1}^n ([y_i - \hat{y}_i]^2 + 2(y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) + (\hat{y}_i - \bar{y})^2)$$

$$\text{Coefficient of Determinants: } R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} \quad [\text{SST} = \text{SSR} + \text{SSE, but why?}]$$

Note The names of the terms are confusing here. Some Statistics texts refer Sum of Squared Regression as Explained Sum of Squares (ESS), and refer Sum of Square Errors as Residual Sum of Squares (RSS). And often Sum of Squares Total will be called as Total Sum of Squares (TSS). One should pay more attention to the conventions of notation usages before reading through the text.

★ Moment Conditions (FOCs) of OLS estimation

$$\min S = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$

$$\frac{\partial S}{\partial \hat{\beta}_0} = 0 \Rightarrow \sum_{i=1}^n 2(Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)(-1) = 0 \Rightarrow \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) = 0 \quad \text{--- (1)}$$

$$\frac{\partial S}{\partial \hat{\beta}_1} = 0 \Rightarrow \sum_{i=1}^n 2(Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)(-X_i) = 0 \Rightarrow \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) X_i = 0 \quad \text{--- (2)}$$

$$\textcircled{1} \sum_{i=1}^n Y_i - n\hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^n X_i = 0 \quad (\text{using } \sum_{i=1}^n Y_i = n\bar{y}, \sum_{i=1}^n X_i = n\bar{x})$$

$$n\bar{y} - n\hat{\beta}_0 - \hat{\beta}_1 n\bar{x} = 0 \Rightarrow \bar{y} - \hat{\beta}_0 = \hat{\beta}_1 \bar{x} \Rightarrow \bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x} \quad \text{--- (1')}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \rightarrow$$

$$\textcircled{2} \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) X_i = 0$$

$$\Rightarrow \sum_{i=1}^n X_i Y_i - \hat{\beta}_0 \sum_{i=1}^n X_i - \hat{\beta}_1 \sum_{i=1}^n X_i^2 = 0 \quad \text{plug in (1')}$$

$$\sum_{i=1}^n X_i Y_i - (\bar{y} - \hat{\beta}_1 \bar{x}) \sum_{i=1}^n X_i - \hat{\beta}_1 \sum_{i=1}^n X_i^2 = 0$$

$$\sum_{i=1}^n X_i Y_i - \bar{y} \sum_{i=1}^n X_i + \hat{\beta}_1 \frac{\sum_{i=1}^n X_i \sum_{i=1}^n X_i}{n} - \hat{\beta}_1 \sum_{i=1}^n X_i^2 = 0 \quad \text{sub } \sum_{i=1}^n X_i = n\bar{x}$$

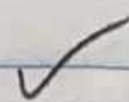
$$\sum_{i=1}^n Y_i = n\bar{y}$$

$$\sum_{i=1}^n X_i Y_i - \bar{y} \cdot n\bar{x} + \hat{\beta}_1 (n\bar{x} \cdot n\bar{x} - \sum_{i=1}^n X_i^2) = 0$$

$$\sum_{i=1}^n X_i Y_i - n\bar{x}\bar{y} = \hat{\beta}_1 \left[\sum_{i=1}^n X_i^2 - (n\bar{x})^2 \right]$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n X_i Y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n X_i^2 - n\bar{x}^2} = \frac{\frac{1}{n} \sum_{i=1}^n X_i Y_i - \bar{x}\bar{y}}{\frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{x}^2} \quad \text{if } n \rightarrow \infty$$

$$\rightarrow \frac{E(X_i Y_i) - E(X_i)E(Y_i)}{E(X_i^2) - [E(X_i)]^2} = \frac{\text{Cov}(X_i, Y_i)}{\text{Var}(X_i)}$$



* Fit of the Model. Role of the R^2 .

$$SST = SSR + SSE$$

$$\sum_{i=1}^n (y_i - \bar{y})^2$$

$$= \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2$$

substitute $y_i - \hat{y}_i = \hat{e}_i$

$$= \sum_{i=1}^n (\hat{e}_i + \hat{y}_i - \bar{y})^2$$

$$= \sum_{i=1}^n (\hat{e}_i^2 + 2 \cdot \hat{e}_i (\hat{y}_i - \bar{y}) + (\hat{y}_i - \bar{y})^2)$$

$$= \sum_{i=1}^n \hat{e}_i^2 + 2 \sum_{i=1}^n \hat{e}_i (\hat{y}_i - \bar{y}) + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$= SSE + 2 \sum_{i=1}^n \hat{e}_i \hat{y}_i - 2 \sum_{i=1}^n \hat{e}_i \bar{y} + SSR$$

b/c ②

prediction & residual
are orthogonal.

$\sum_{i=1}^n (\hat{e}_i) = 0$ sum of residual = 0

So $SST = SSE + SSR$ is shown. **

$$GDP_t = \beta_0 + \beta_1 RainFall_t + \epsilon_t$$

$$t = 1989(1), 1990(2), 1991(3), \dots, 2018(20)$$

Appendix: Taxonomy of Data, Cross-sectional v.s. Panel

Econometrics textbooks usually start with the introduction of cross-sectional data since the methodology, estimation, and analysis could be straight-forwardly constructed based on the Gauss-Markov assumptions. In the previous lectures, we had already gone through the main analysis methods, Ordinary Least Square (OLS) and learned how to make the estimators become unbiased and consistent. The specification (i.e. econometric model) is usually written as following, based on the assumptions that the data was randomly selected from the population at the same time. $i = 1, 2, \dots, N$ could be states, cities, countries, individuals, etc.

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + e_i$$

However, most of the institutions in charge of collecting the macroeconomic data usually keep track of their observations. For example, the World Bank would like to know each country's annual GDP, CPI, or currency exchange rate; and the Federal Reserve pays attention to the indices of the financial market such as interest rates and Yields to the Treasury Bonds day-by-day. Hence, following an identical group of sample creates another dimension of the existing variables in the OLS model, which is the time. Furthermore, usually, we will put another subscript t underneath the variables in the original OLS equation, and the finite time index is usually denoted by $t = 1, 2, \dots, T$.

$$Y_{it} = \beta_0 + \beta_1 X_{1it} + \beta_2 X_{2it} + \dots + \beta_k X_{kit} + e_{it}$$

However, what is wrong with the OLS? Aren't they just different in notations? Aren't they the same mathematically? What will happen if we run the Ordinary Least Square on panel data or applied the Weighted Least Square by using White's heteroskedasticity robustness correction?

The answer to the above is one of the key Gauss-Markov assumptions is no longer hold. Specifically, because the nature of panel data, the condition $\text{Cov}(e_i, e_j) = 0, i \neq j$ usually will be violated. If we use the notation in the panel data equation, $\text{Cov}(e_{it}, e_{is}) \neq 0$, and it is also the fundamental reason why we need to learn another method as well as some additional assumption on the error terms for the panel data or time-series data.

Let me use a simple reduced example to walk you through the idea. Assume we would like to know should the flash flood affects the economic activities in the DC area. For some reason, we could only get access to the GDP and RainFall data of the State of Virginia since 1989. Now, we should be able to construct a model that is similar to the following equation. $i = 1, t = 1989, 1990, \dots, 2019$ ($N = 1, T = 20$). In other words, for the State of Virginia,

$$GDP_t = \beta_0 + \beta_1 \text{RainFall}_t + e_t$$

where

$$t = 1989(1), 1990(2), 1991(3), \dots, 2018(20)$$

Then we can think about it. RainFall_t for each year should be independent with each other, but what about the GDP_t ? We have learned Real Business Cycle theory in Macroeconomics, so we should presume $\text{Cov}(GDP_t, GDP_s) \neq 0$, for $t \neq s$ or $t \rightarrow s$. Then if we plug in two data points of dependent variables $GDP_t = \beta_0 + \beta_1 \text{RainFall}_t + e_t$ and $GDP_s = \beta_0 + \beta_1 \text{RainFall}_s + e_s$.

$$\text{Cov}(GDP_t, GDP_s) \neq 0$$

$$\Rightarrow \text{Cov}(\beta_0 + \beta_1 \text{RainFall}_t + e_t, \beta_0 + \beta_1 \text{RainFall}_s + e_s) \neq 0$$

$$\Rightarrow \text{Cov}(e_t, \beta_0 + \beta_1 \text{RainFall}_s + e_s) \neq 0$$

$$\Rightarrow \text{Cov}(e_t, e_s) \neq 0$$

So now, you should be able to see what happened if we apply OLS to this data set. Since one of the important assumption of OLS is violated, the estimator of the coefficient will NOT be unbiased toward to the parameters, no matter how large the sample size is. [i.e. $E(b_1) \neq \beta_1$ $E(b_0) \neq \beta_0$] The whole estimation process becomes problematic and meaningless.