

Lecture Note 1

2014-03-18

1.1. Classical Linear Regression

Definition \mathbf{y} is a dependent variable.

$$\mathbf{y} := \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad \mathbf{y} \in \mathbb{R}^n, \quad y_i \in \mathbb{R} \quad \forall i \in \{1, 2, \dots, n\}$$

Definition \mathbf{X} is an independent variable.

$$\begin{aligned} \mathbf{X} &:= \begin{pmatrix} 1 & x_{11} & x_{21} & \cdots & x_{k1} \\ 1 & x_{12} & x_{22} & \cdots & x_{k1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & x_{2n} & \cdots & x_{kn} \end{pmatrix} & \mathbf{X} \in \mathbb{R}^{n \times k}, \quad x_{ij} \in \mathbb{R} \quad \forall j \in \{1, 2, \dots, k\} \\ &= (\mathbf{1} \ \mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_k) \\ &= \begin{pmatrix} \mathbf{x}_{,1} \\ \mathbf{x}_{,2} \\ \vdots \\ \mathbf{x}_{,n} \end{pmatrix} \end{aligned}$$

Definition $\boldsymbol{\beta}$ is a parameter.

$$\boldsymbol{\beta} := \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} \quad \boldsymbol{\beta} \in \mathbb{R}^{k+1}$$

Definition $\tilde{\boldsymbol{\varepsilon}}$ is a residual.

$$\begin{aligned} \tilde{\boldsymbol{\varepsilon}} &:= \begin{pmatrix} \tilde{\varepsilon}_1 \\ \tilde{\varepsilon}_2 \\ \vdots \\ \tilde{\varepsilon}_n \end{pmatrix} \\ \Rightarrow \tilde{\mathbf{y}} &= \mathbf{X}\boldsymbol{\beta} + \tilde{\boldsymbol{\varepsilon}} \end{aligned}$$

Since $\tilde{\boldsymbol{\varepsilon}}$ is stochastic, \mathbf{y} is also stochastic. However, \mathbf{X} and $\boldsymbol{\beta}$ are not random.

1.1.1. Least Squares

However, in reality, \mathbf{y} and \mathbf{X} are given, while $\boldsymbol{\beta}$ is unknown. Therefore, we should estimate it based on several assumptions; *Least Squares Estimation* is one way to estimate it.

Objective function of the estimation can be denoted as below. By minimizing the function, we can estimate β , which is indeed unknown in reality.

$$\begin{aligned} & \min_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}_i' \tilde{\beta})^2 \\ &= \min_{\beta} \tilde{\mathbf{e}}' \tilde{\mathbf{e}} \end{aligned}$$

Classical linear regression equation can be depicted as below.

$$\begin{aligned} \tilde{\mathbf{y}} &= \mathbf{X}\beta + \tilde{\mathbf{e}} \\ \Rightarrow E(\tilde{\mathbf{y}}|\mathbf{X}) &= \mathbf{X}\hat{\beta} \quad (\text{conditional regression } \mathbf{y} \text{ on } \mathbf{X}) \end{aligned}$$

$$\text{Assumption 1} \quad E(\tilde{\mathbf{e}}|\mathbf{X}) = \mathbf{0} \Leftrightarrow E(\tilde{\varepsilon}_i|\mathbf{X}) = 0 \quad \forall i$$

$$\text{Assumption 2} \quad E(\tilde{\mathbf{e}}\tilde{\mathbf{e}}'|\mathbf{X}) = \sigma^2 \mathbf{I}_n \Leftrightarrow E(\tilde{\varepsilon}_i^2|\mathbf{X}) = \sigma^2 \text{ and } E(\tilde{\varepsilon}_i\tilde{\varepsilon}_j|\mathbf{X}) = 0 \quad \forall i \neq j$$

$$E(\tilde{\mathbf{e}}\tilde{\mathbf{e}}'|\mathbf{X}) = \begin{pmatrix} \sigma^2 & 0 & \cdots & 0 \\ 0 & \sigma^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma^2 \end{pmatrix}$$

Assumption 2 is called as *Homoskedasticity* assumption, which is frequently violated in reality. Hypothesis testing based on OLS estimator will not be valid if this assumption is violated.

$$\text{Assumption 3} \quad \text{rank}(\mathbf{X}) = k + 1$$

In terms of econometrics, following statements are equivalent. \mathbf{X} is full rank; determinants of \mathbf{X} are not zero; \mathbf{X} is non-singular; $\exists \mathbf{X}^{-1}$; columns of \mathbf{X} are linearly independent.

$$\text{Assumption 4} \quad \tilde{\mathbf{e}}|\mathbf{X} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

Assumption 4 is called as *Normality* assumption, which is not critical, but convenient; hence, we will use this assumption.

1.1.2. Ordinary Least Squares Estimation

Then, we can estimate β by solving the minimization problem above.

$$\hat{\beta} := \underset{\beta}{\operatorname{argmin}} \tilde{\mathbf{e}}' \tilde{\mathbf{e}}$$

The problem can be solved by using differentiation.

$$\tilde{\mathbf{e}}' \tilde{\mathbf{e}} = (\mathbf{y} - \mathbf{X}\tilde{\beta})' (\mathbf{y} - \mathbf{X}\tilde{\beta})$$

$$\begin{aligned}
&= \mathbf{y}'\mathbf{y} - \tilde{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{y} - \mathbf{y}'\mathbf{X}\tilde{\boldsymbol{\beta}} + \tilde{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\tilde{\boldsymbol{\beta}} \\
&= \mathbf{y}'\mathbf{y} - 2\tilde{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{y} + \tilde{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\tilde{\boldsymbol{\beta}} \\
\Rightarrow \frac{\partial \tilde{\boldsymbol{\varepsilon}}'\tilde{\boldsymbol{\varepsilon}}}{\partial \tilde{\boldsymbol{\beta}}} &= -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} \\
&= \mathbf{0} \\
\Rightarrow \mathbf{X}'\mathbf{y} &= \mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} \\
\therefore \hat{\boldsymbol{\beta}} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \\
&= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\tilde{\mathbf{y}} \\
&= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\boldsymbol{\beta} + \tilde{\boldsymbol{\varepsilon}}) \\
&= \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\tilde{\boldsymbol{\varepsilon}}
\end{aligned}$$

Note that the term *estimator* means function, while the term *estimate* means numeric.

1.1.3. Properties of OLS Estimator

Property 1 $E(\hat{\boldsymbol{\beta}}|\mathbf{X}) = \boldsymbol{\beta}$; $\hat{\boldsymbol{\beta}}$ is unbiased.

Proof

$$\begin{aligned}
\hat{\boldsymbol{\beta}} &= \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\tilde{\boldsymbol{\varepsilon}} \\
\therefore E(\hat{\boldsymbol{\beta}}|\mathbf{X}) &= E[\boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\tilde{\boldsymbol{\varepsilon}}|\mathbf{X}] \\
&= E(\boldsymbol{\beta}|\mathbf{X}) + E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\tilde{\boldsymbol{\varepsilon}}|\mathbf{X}] \\
&= \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\tilde{\boldsymbol{\varepsilon}}|\mathbf{X}) \quad \text{Assumption 1} \\
&= \boldsymbol{\beta} \quad \blacksquare
\end{aligned}$$

Property 2 $\hat{\boldsymbol{\beta}}$ is efficient.

In terms of econometrics, following statements are equivalent. $\hat{\boldsymbol{\beta}}$ is the *best* estimator; $\hat{\boldsymbol{\beta}}$ has the smallest variance; $\hat{\boldsymbol{\beta}}$ is the most efficient estimator.

If an estimator has large variance, we cannot say whether the estimator is significant or not. With Assumption 4, it is well-known that the OLS estimator $\hat{\boldsymbol{\beta}}$ is equal to the ML estimator $\hat{\boldsymbol{\beta}}_{\text{ML}}$, which has the smallest variance. Thus, $\hat{\boldsymbol{\beta}}$ is efficient.

Since the proof requires too tedious mathematics, we will skip the proof.

$$\begin{aligned}
\text{Var}(\hat{\boldsymbol{\beta}}|\mathbf{X}) &= E[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})'|\mathbf{X}] \\
&= E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\tilde{\boldsymbol{\varepsilon}}\tilde{\boldsymbol{\varepsilon}}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}|\mathbf{X}] \quad \text{Assumption 2} \\
&= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\
&= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}
\end{aligned}$$

Property 3 $\text{plim}\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}$; $\hat{\boldsymbol{\beta}}$ is consistent.

Proof

$$\begin{aligned}\frac{\mathbf{X}'\mathbf{X}}{n} &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \\ &= \hat{E}(\mathbf{x}_i \mathbf{x}_i') \quad \text{as } n \rightarrow \infty\end{aligned}$$

$$\begin{aligned}\frac{\mathbf{X}'\tilde{\boldsymbol{\varepsilon}}}{n} &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \tilde{\varepsilon}_i \\ &= \hat{E}(\mathbf{x}_i \tilde{\varepsilon}_i) \quad \text{as } n \rightarrow \infty \\ &= \mathbf{x}_i \hat{E}(\tilde{\varepsilon}_i) \\ &= \mathbf{0}\end{aligned}$$

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\tilde{\boldsymbol{\varepsilon}} \\ &= \boldsymbol{\beta} + \left(\frac{\mathbf{X}'\mathbf{X}}{n}\right)^{-1} \frac{\mathbf{X}'\tilde{\boldsymbol{\varepsilon}}}{n} \\ &= \boldsymbol{\beta} \quad \text{as } n \rightarrow \infty\end{aligned}$$

$$\begin{aligned}\text{Var}(\hat{\boldsymbol{\beta}}|\mathbf{X}) &= \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \\ &= \frac{\sigma^2}{n} \left(\frac{\mathbf{X}'\mathbf{X}}{n}\right)^{-1} \\ &= 0 \times \hat{E}(\mathbf{x}_i \mathbf{x}_i') \quad \text{as } n \rightarrow \infty \\ &= \mathbf{0} \quad \blacksquare\end{aligned}$$

Based on these properties, OLS estimator is BLUE; *Best, Linear, Unbiased* and *Efficient*. By Gauss-Markov Theorem, it can be proved that OLS estimator is BLUE. Note that Lucas, one of the prominent economist, said, “OLS is everything!”

1.2. Hypothesis Testing

1.2.1. Test Statistic

We can conduct *Hypothesis Testing* with the estimator already estimated above. The null hypothesis H_0 can be denoted as below.

$$H_0: \mathbf{R}\boldsymbol{\beta} = \mathbf{c} \quad \mathbf{R} \in \mathbb{R}^{q \times (k+1)}, \mathbf{c} \in \mathbb{R}^q$$

In above expression, q is the number of parameters to be tested. Then, the test statistic for the null hypothesis can be written as below.

$$t = \frac{\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}}{\sqrt{\text{Var}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c})}}$$

Above t will be utilized in the process of hypothesis testing. Under the null hypothesis,

above t statistic will follow t distribution, which is similar to Normal distribution. Thus, if the absolute value of calculated t statistic is large, then we can reject the null hypothesis under a given significance level; we can reject the null hypothesis if the absolute value of $\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}$ is significantly big or if $\sqrt{\text{Var}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c})}$ is significantly small.

$$\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c} \sim N[0, \text{Var}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c})] \quad \text{under } H_0$$

$$\begin{aligned} \because \mathbf{R}\boldsymbol{\beta} &= \mathbf{c} \\ \hat{\boldsymbol{\beta}} &= \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\tilde{\boldsymbol{\varepsilon}} \\ \Rightarrow \hat{\boldsymbol{\beta}} &\sim N[\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}] \quad \text{under } H_0 \end{aligned}$$

$$\begin{aligned} \therefore \text{Var}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}) &= \text{Var}[\mathbf{R}(\boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\tilde{\boldsymbol{\varepsilon}}) - \mathbf{c}] \\ &= \text{Var}[\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\tilde{\boldsymbol{\varepsilon}}] \\ &= \sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}' \quad \text{under } H_0 \end{aligned}$$

The problem is that σ^2 is unknown in reality. Instead of σ^2 , we adopt $\hat{\sigma}^2$ as an alternative.

1.2.2. Estimation of Residual Variance

$$\begin{aligned} \hat{\boldsymbol{\varepsilon}}'\hat{\boldsymbol{\varepsilon}} &= (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \\ \therefore \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} &= \mathbf{M}_X\mathbf{y} \quad \text{where } \mathbf{M}_X = [\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] \\ &= \mathbf{M}_X\tilde{\boldsymbol{\varepsilon}} \end{aligned}$$

By Spectral Decomposition, \mathbf{M}_X can be decomposed as below.

$$\begin{aligned} \mathbf{M}_X &= \mathbf{V}\mathbf{\Lambda}\mathbf{V}' \quad \ni \quad \mathbf{V}'\mathbf{V} = \mathbf{V}\mathbf{V}' = \mathbf{I}_n \\ \mathbf{V} &= \begin{pmatrix} v_1 & 0 & \cdots & 0 \\ 0 & v_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & v_n \end{pmatrix} \\ \mathbf{\Lambda} &= \begin{pmatrix} 1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 1 & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{I}_{n-(k+1)} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \end{aligned}$$

Therefore, we can find a property of its *Eigenvalue*.

$$\begin{aligned} \lambda \mathbf{V} &= \mathbf{M}_X\mathbf{V} \\ &= \mathbf{M}_X\mathbf{M}_X\mathbf{V} \\ &= \mathbf{M}_X(\mathbf{M}_X\mathbf{V}) \end{aligned}$$

$$\begin{aligned}
&= \mathbf{M}_X(\lambda \mathbf{V}) \\
&= \lambda \mathbf{M}_X \mathbf{V} \\
&= \lambda^2 \mathbf{V}
\end{aligned}$$

$$\therefore \lambda = 0 \text{ or } 1$$

Hereafter, we will examine the distributional characteristics of $\hat{\boldsymbol{\varepsilon}}$ with the spectral decomposition of \mathbf{M}_X mentioned above.

$$\begin{aligned}
(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) &= \hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}} \\
&= (\mathbf{M}_X \tilde{\boldsymbol{\varepsilon}})' \mathbf{M}_X \tilde{\boldsymbol{\varepsilon}} \\
&= \tilde{\boldsymbol{\varepsilon}}' \mathbf{M}_X \tilde{\boldsymbol{\varepsilon}} \\
&= \tilde{\boldsymbol{\varepsilon}}' \mathbf{V} \boldsymbol{\Lambda} \mathbf{V}' \tilde{\boldsymbol{\varepsilon}} \\
&= \tilde{\boldsymbol{\varepsilon}}_*' \boldsymbol{\Lambda} \tilde{\boldsymbol{\varepsilon}}_* \quad \text{where } \tilde{\boldsymbol{\varepsilon}}_* = \mathbf{V}' \tilde{\boldsymbol{\varepsilon}}
\end{aligned}$$

$$\begin{aligned}
\therefore E(\tilde{\boldsymbol{\varepsilon}}_*) &= E(\mathbf{V}' \tilde{\boldsymbol{\varepsilon}}) \\
&= \mathbf{0}
\end{aligned}$$

$$\begin{aligned}
\therefore \text{Var}(\tilde{\boldsymbol{\varepsilon}}_*) &= E(\tilde{\boldsymbol{\varepsilon}}_* \tilde{\boldsymbol{\varepsilon}}_*') \\
&= E(\mathbf{V}' \tilde{\boldsymbol{\varepsilon}} \tilde{\boldsymbol{\varepsilon}}' \mathbf{V}) \\
&= \sigma^2 \mathbf{V}' \mathbf{V} \\
&= \sigma^2 \mathbf{I}_n
\end{aligned}$$

$$\therefore \tilde{\boldsymbol{\varepsilon}}_* \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\Leftrightarrow \frac{\tilde{\boldsymbol{\varepsilon}}_*}{\sigma} \sim N(\mathbf{0}, \mathbf{I}_n)$$

$$\Leftrightarrow \frac{\tilde{\boldsymbol{\varepsilon}}_*}{\sigma} \sim N(0, 1) \quad \forall i$$

$$\therefore \hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}} = \tilde{\boldsymbol{\varepsilon}}_*' \boldsymbol{\Lambda} \tilde{\boldsymbol{\varepsilon}}_*$$

$$\begin{aligned}
\Rightarrow \frac{\hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}}}{\sigma^2} &= \frac{\tilde{\boldsymbol{\varepsilon}}_*'}{\sigma} \boldsymbol{\Lambda} \frac{\tilde{\boldsymbol{\varepsilon}}_*}{\sigma} \\
&= \frac{1}{\sigma^2} \tilde{\boldsymbol{\varepsilon}}_*' \boldsymbol{\Lambda} \tilde{\boldsymbol{\varepsilon}}_* \\
&= \frac{1}{\sigma^2} (\tilde{\varepsilon}_{*1}^2 + \cdots + \tilde{\varepsilon}_{*n-k-1}^2) \\
&= \sum_{j=1}^{n-k-1} \left(\frac{\tilde{\varepsilon}_{*j}}{\sigma} \right)^2 \sim \chi^2(n-k-1)
\end{aligned}$$

Note that if $z_i \sim N(0,1)$, then $\sum_{i=1}^n z_i^2 \sim \chi^2(n)$.

$$\therefore \hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}} = \sigma^2 \chi^2(n-k-1)$$

Therefore, if $\hat{\sigma}^2 = \hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}}$, then $E(\hat{\sigma}^2) = (n-k-1)\sigma^2 \neq \sigma^2$ and hence biased. Furthermore, if $z_i \sim \chi^2(n)$, then $E(z_i) = n$ and $Var(z_i) = n^2$

$$\begin{aligned} \therefore \hat{\sigma}^2 &= \frac{\hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}}}{n-k-1} \\ &\sim \frac{\sigma^2}{n-k-1} \chi^2(n-k-1) \end{aligned}$$

$$\Rightarrow \frac{n-k-1}{\sigma^2} \hat{\sigma}^2 \sim \chi^2(n-k-1)$$

1.2.3. Construction of Test Statistic

With these distributional characteristics, recall the test statistic for the null hypothesis mentioned above.

$$\begin{aligned} t &= \frac{\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}}{\sqrt{Var(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c})}} \\ &= \frac{\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}}{\sigma \sqrt{\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'}} \\ &\sim N(0,1) \quad \text{under } H_0 \end{aligned}$$

Note that if \tilde{z} and \tilde{x} are independent and follow Standard Normal distribution and Chi-square distribution respectively, then the statement below can be derived.

$$\frac{\tilde{z}}{\sqrt{\tilde{x}/n}} \sim t(n)$$

Since we can construct one Normal variable and the other Chi-square variable under the null hypothesis, we can also derive another variable which follows t distribution.

$$\begin{aligned} \frac{\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}}{\sigma \sqrt{\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'}} &\sim N(0,1) \\ \frac{\sqrt{\hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}} / \sigma^2}}{n-k-1} &\sim \chi^2(n-k-1) \end{aligned}$$

$$\begin{aligned}
\therefore \frac{\frac{\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}}{\sigma \sqrt{\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'}}}{\frac{\sqrt{\hat{\boldsymbol{\varepsilon}}'\hat{\boldsymbol{\varepsilon}}/\sigma^2}}{n-k-1}} &= \frac{\frac{\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}}{\sqrt{\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'}}}{\sqrt{\frac{\hat{\boldsymbol{\varepsilon}}'\hat{\boldsymbol{\varepsilon}}}{n-k-1}}} \\
&= \frac{\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}}{\hat{\sigma} \sqrt{\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'}} \\
&\sim t(n-k-1)
\end{aligned}$$

Through the t statistic above, we can test the null hypothesis correctly if the assumptions are not violated. And this process can also be used in many cases by manipulating \mathbf{R} and \mathbf{c} of the null hypothesis.

For instance, in the given regression equation $y = \beta_0 + \beta_1 x + \varepsilon$ and the null hypothesis is stated as $H_0: \beta_1 = 0$, then \mathbf{R} and \mathbf{c} for the test can be designed respectively.

$$\begin{aligned}
\mathbf{R} &= (0 \ 1) \\
\mathbf{c} &= 0
\end{aligned}$$

$$\Rightarrow \mathbf{R}\boldsymbol{\beta} = \mathbf{c}$$

$$\therefore (0 \ 1) \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} = 0$$

Lecture Note 2

2014-03-25

2.1. Wald Test

2.1.1. Construction of Test Statistic

$$H_0: \mathbf{R}\boldsymbol{\beta} = \mathbf{c}$$

Note Wald criterion

$$\mathbf{u} \sim N(\mathbf{0}, \boldsymbol{\Sigma})$$

$$\Rightarrow \boldsymbol{\Sigma}^{-1/2} \mathbf{u} \sim N(\mathbf{0}, \mathbf{I}_n)$$

$$\therefore \mathbf{u}' \boldsymbol{\Sigma}^{-1} \mathbf{u} \sim \chi^2(q)$$

Or equivalently,

$$u_i \stackrel{i.i.d.}{\sim} N(0, 1)$$

$$\Rightarrow \sum_{i=1}^n u_i^2 \sim \chi^2(n)$$

For $\hat{\boldsymbol{\beta}}_{\text{OLS}}$, hereafter $\hat{\boldsymbol{\beta}}$, we already know the fact that

$$\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c} \sim N(\mathbf{0}, \sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}') \quad \text{under } H_0$$

$$\therefore (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c})' [\sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}) \sim \chi^2(q)$$

Recall that

$$\frac{(n-k-1)\hat{\sigma}^2}{\sigma^2} \sim \chi^2(n-k-1)$$

Remark For $x_1 \sim \chi^2(n_1)$ and $x_2 \sim \chi^2(n_2)$

$$\Rightarrow \frac{x_1/n_1}{x_2/n_2} \sim F(n_1, n_2)$$

$$\frac{(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c})' [\sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c})/q}{(n-k-1)\hat{\sigma}^2 / [(n-k-1)\sigma^2]} = \frac{(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c})' [\sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c})/q}{\hat{\sigma}^2 / \sigma^2}$$

$$\sim F(q, n - k - 1)$$

Therefore,

$$\frac{1}{\hat{\sigma}^2} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c})' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}]^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}) \sim F(q, n - k - 1) \quad \text{under } H_0$$

This is Wald test statistic.

2.1.2. Relation between t Test and F Test

Note The mathematical relationship between t distribution with n degrees of freedom and F distribution with 1 and n degrees of freedom can be depicted as below.

$$\begin{aligned} x &\sim t(n) \\ \Rightarrow x^2 &\sim F(1, n) \end{aligned}$$

Not rigorously but roughly, if $q = 1$, then $t^2 = F$; this distributional characteristic is similar to the relationship between Standard Normal distribution and Chi-square distribution.

$$\begin{aligned} y &\sim N(0, 1) \\ \Rightarrow y^2 &\sim \chi^2(1) \end{aligned}$$

2.2. Heteroskedasticity

Recall t test statistic above.

$$t = \frac{\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}}{\hat{\sigma} \sqrt{\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}}}$$

Traditional Econometrics course often teaches the way to test the null hypothesis based on several assumptions, such as Homoskedasticity assumption; unfortunately, these binding assumptions are frequently violated in reality and hence Econometric course also teaches various helpful remedies whereby the test can be correctly conducted. This is a brief sketch of an undergraduate level Econometrics course.

Heteroskedasticity, which is the violation of Assumption 2, can be sketched as below; for a classical linear regression equation below,

$$\tilde{\mathbf{y}} = \mathbf{X}\boldsymbol{\beta} + \tilde{\boldsymbol{\varepsilon}}$$

Heteroskedasticity literally means that the size of variance for each residual term is heterogeneous, which could mathematically be denoted as below.

$$E(\tilde{\varepsilon}_i^2|\mathbf{X}) = \sigma_i^2 \quad \sigma_i^2 \neq \sigma_j^2 \quad \forall i \neq j$$

Or compactly,

$$\begin{aligned} E(\tilde{\varepsilon}\tilde{\varepsilon}'|\mathbf{X}) &= \mathbf{\Omega} \\ &= \begin{pmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{pmatrix} \\ &= \sigma^2 \underbrace{\begin{pmatrix} \Sigma_1 & 0 & \cdots & 0 \\ 0 & \Sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \Sigma_n \end{pmatrix}}_{\mathbf{\Sigma}} \quad \Sigma \in \mathbb{R}^{n \times n}, \quad \Sigma_i \in \mathbb{R}, \quad \sigma_i^2 = \sigma^2 \Sigma_i \\ &= \sigma^2 \mathbf{\Sigma} \end{aligned}$$

Although there exists Heteroskedasticity problem in the residual, $\hat{\beta}$ is still unbiased since Assumption 1 is not violated; i.e. $E(\tilde{\varepsilon}|\mathbf{X}) = \mathbf{0}$. However, $\hat{\beta}$ is not efficient anymore.

$$\begin{aligned} Var(\hat{\beta}|\mathbf{X}) &= Var(\hat{\beta} - \beta|\mathbf{X}) \\ &= Var[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\tilde{\varepsilon}|\mathbf{X}] \\ &= E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\tilde{\varepsilon}\tilde{\varepsilon}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}|\mathbf{X}] \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\tilde{\varepsilon}\tilde{\varepsilon}'|\mathbf{X})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{\Omega}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\ &\neq \sigma^2(\mathbf{X}'\mathbf{X})^{-1} \end{aligned}$$

Because of this inefficiency, other econometric techniques to calculate the size of variance are required to test the null hypothesis precisely.

Case 1 If $\mathbf{\Sigma}$ is known and diagonal, then we will use *Weighted Least Squares*. We don't know σ^2 in this case. Hence, we will again use Spectral Decomposition, which is also known as Cholesky Decomposition or Eigenvalue Decomposition.

Theorem Positive semi-definite matrix can always be decomposed as below.

$$\mathbf{\Sigma}^{-1} = \mathbf{P}\mathbf{P}'$$

Instead of introducing a tedious proof of above Spectral Theorem, we will just use this theorem hereafter.

$$\begin{aligned} \Rightarrow \mathbf{\Sigma} &= (\mathbf{P}\mathbf{P}')^{-1} \\ &= \mathbf{P}'^{-1}\mathbf{P}^{-1} \\ \Rightarrow \mathbf{P}'\mathbf{\Sigma}\mathbf{P} &= \mathbf{I}_n \end{aligned}$$

$$\Rightarrow \mathbf{P} := \begin{pmatrix} \Sigma_1^{-1/2} & 0 & \cdots & 0 \\ 0 & \Sigma_2^{-1/2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \Sigma_n^{-1/2} \end{pmatrix}$$

$$\therefore \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$$\Rightarrow \mathbf{P}'\mathbf{y} = \mathbf{P}'\mathbf{X}\boldsymbol{\beta} + \mathbf{P}'\boldsymbol{\varepsilon}$$

$$\Rightarrow \mathbf{y}_* = \mathbf{X}_*\boldsymbol{\beta} + \boldsymbol{\varepsilon}_*$$

$$\begin{aligned} \Rightarrow \text{Var}(\boldsymbol{\varepsilon}_*|\mathbf{X}) &= \text{Var}(\mathbf{P}'\boldsymbol{\varepsilon}|\mathbf{X}) \\ &= E(\mathbf{P}'\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'\mathbf{P}|\mathbf{X}) \\ &= \sigma^2\mathbf{P}'\boldsymbol{\Sigma}\mathbf{P} \\ &= \sigma^2\mathbf{I}_n \end{aligned}$$

$$\begin{aligned} \therefore \hat{\boldsymbol{\beta}}_{\text{WLS}} &= (\mathbf{X}_*'\mathbf{X}_*)^{-1}\mathbf{X}_*'\mathbf{y}_* \\ &= (\mathbf{X}'\mathbf{P}\mathbf{P}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{P}\mathbf{P}'\mathbf{y} \\ &= (\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{y} \end{aligned}$$

Therefore, above $\hat{\boldsymbol{\beta}}_{\text{WLS}}$ is BLUE.

$$\begin{aligned} \therefore \text{Var}(\hat{\boldsymbol{\beta}}_{\text{WLS}}|\mathbf{X}) &= \sigma^2(\mathbf{X}_*'\mathbf{X}_*)^{-1} \\ &= \sigma^2(\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X})^{-1} \end{aligned}$$

However, there is no answer for the question that “How do we know the structure of $\boldsymbol{\Sigma}$ in reality?” If there is no justification for this question, WLS is no more valid; *Inference Manipulation*. Various papers have adopted WLS without any notification. This bad attitude should be avoided.

Case 2 If $\boldsymbol{\Sigma}$ (or $\boldsymbol{\Omega}$ equivalently) is unknown, then we should estimate $\hat{\boldsymbol{\Omega}}$ firstly and then conduct remain procedures secondly; this is called as *Feasible Generalized Least Squares*.

Indeed, it is impossible to estimate $\hat{\boldsymbol{\Omega}}$ since it contains n^2 entries in it and we cannot estimate it based on n observations. Instead, we can successfully estimate $\hat{E}(\mathbf{X}'\boldsymbol{\Omega}\mathbf{X}) = \mathbf{X}'\widehat{\boldsymbol{\Omega}}\mathbf{X}$, which includes k^2 elements in it; usually, n is bigger than k . Note below.

$$\mathbf{X}'\boldsymbol{\Omega}\mathbf{X} = \sum_{i=1}^n \sigma_i^2 \mathbf{x}_i \mathbf{x}_i'$$

$$\therefore \mathbf{x}_{\cdot i} = \begin{pmatrix} 1 \\ x_{1i} \\ x_{2i} \\ \vdots \\ x_{ki} \end{pmatrix}$$

2.2.1. White's Robust Variance

Definition White (1980, *Econometrica*) suggests *Heteroskedasticity Consistent Variance Estimator*, which is also known as *HC Estimator* or *White Standard Error*.

$$\begin{aligned}\hat{E}(\mathbf{X}'\boldsymbol{\Omega}\mathbf{X}) &= \sum_{i=1}^n \hat{\varepsilon}_i^2 \mathbf{x}_i \mathbf{x}_i' \quad \text{for} \quad \sum_{i=1}^n \sigma_i^2 \mathbf{x}_i \mathbf{x}_i' \\ \therefore \text{Var}(\hat{\boldsymbol{\beta}}_{\text{LS}}|\mathbf{X}) &= (\mathbf{X}'\mathbf{X})^{-1} \hat{E}(\mathbf{X}'\boldsymbol{\Omega}\mathbf{X}) (\mathbf{X}'\mathbf{X})^{-1} \\ &= \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\sum_{i=1}^n \hat{\varepsilon}_i^2 \mathbf{x}_i \mathbf{x}_i' \right) \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \\ \hat{\varepsilon}_i &= y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}\end{aligned}$$

Fortunately, White $\text{Var}(\hat{\boldsymbol{\beta}}_{\text{LS}}|\mathbf{X})$ is consistent with Heteroskedasticity.

2.2.2. Newey-West Robust Variance

However, this estimator is still problematic if there exists not only Heteroskedasticity, but also Autocorrelation, simultaneously.

$$E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'|\mathbf{X}) = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1n} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \cdots & \sigma_n^2 \end{pmatrix}$$

Or equivalently,

$$\begin{aligned}\Rightarrow E(\varepsilon_i \varepsilon_j | \mathbf{X}) &= \begin{cases} \sigma_i^2 & \forall i = j \\ \sigma_{ij} & \forall i \neq j \end{cases} \\ \therefore \mathbf{X}' \begin{pmatrix} 0 & \sigma_{12} & \cdots & \sigma_{1n} \\ \sigma_{21} & 0 & \cdots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \cdots & 0 \end{pmatrix} \mathbf{X} &= \sum_{j=1}^{n-1} \sum_{i=j+1}^n E(\varepsilon_i \varepsilon_{i-j} | \mathbf{X}) (\mathbf{x}_i \mathbf{x}_{i-j}' + \mathbf{x}_{i-j} \mathbf{x}_i')\end{aligned}$$

The term j of above equation means the distance between two observations, which should be determined decreasingly in increasing distance. *Bartlett Kernel* is frequently applied as a weighting scheme such that above relationships can be reflected.

$$w_j = \begin{cases} 1 - \frac{j}{J+1} & \forall 1 \leq j \leq J \\ 0 & \forall J < j \leq n \end{cases}$$

Definition Newey and West (1987, *Econometrica*) suggests *Heteroskedasticity Autocorrelation Consistent Variance Estimator*, which is also known as *HAC Estimator* or *Newey-West Standard Error*.

$$\hat{E}(\mathbf{X}'\Omega\mathbf{X}) = \sum_{i=1}^n \hat{\varepsilon}_i^2 \mathbf{x}_i \mathbf{x}_i' + \sum_{j=1}^J w_j \sum_{i=j+1}^n \hat{\varepsilon}_i \hat{\varepsilon}_{i-j} (\mathbf{x}_i \mathbf{x}_{i-j}' + \mathbf{x}_{i-j} \mathbf{x}_i')$$

$$\therefore \text{Var}(\hat{\beta}_{\text{LS}}|\mathbf{X}) = (\mathbf{X}'\mathbf{X})^{-1} \hat{E}(\mathbf{X}'\Omega\mathbf{X}) (\mathbf{X}'\mathbf{X})^{-1}$$

Newey-West $\text{Var}(\hat{\beta}_{\text{LS}}|\mathbf{X})$ is consistent with Heteroskedasticity and Autocorrelation.

2.3. Seemingly Unrelated Regression

Suppose that there are m regression equations.

$$\begin{aligned} \mathbf{y}^1 &= \mathbf{X}^1 \boldsymbol{\beta}^1 + \mathbf{u}^1 \\ \mathbf{y}^2 &= \mathbf{X}^2 \boldsymbol{\beta}^2 + \mathbf{u}^2 \\ &\vdots \\ \mathbf{y}^m &= \mathbf{X}^m \boldsymbol{\beta}^m + \mathbf{u}^m \end{aligned}$$

Note that $\boldsymbol{\beta}^1, \boldsymbol{\beta}^2, \dots, \boldsymbol{\beta}^m$ can have different dimensions, respectively.

$$\begin{aligned} \Rightarrow \begin{pmatrix} \mathbf{y}^1 \\ \mathbf{y}^2 \\ \vdots \\ \mathbf{y}^m \end{pmatrix} &= \begin{pmatrix} \mathbf{X}^1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{X}^2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{X}^m \end{pmatrix} \begin{pmatrix} \boldsymbol{\beta}^1 \\ \boldsymbol{\beta}^2 \\ \vdots \\ \boldsymbol{\beta}^m \end{pmatrix} + \begin{pmatrix} \mathbf{u}^1 \\ \mathbf{u}^2 \\ \vdots \\ \mathbf{u}^m \end{pmatrix} \\ \Rightarrow \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{u} \end{aligned}$$

If the assumption below is satisfied, then there is no problem at all.

$$\begin{aligned} E(\mathbf{u}^i|\mathbf{X}) &= \mathbf{0} \quad \forall i \in \{1, 2, \dots, m\} \\ E(\mathbf{u}^i \mathbf{u}^{i'}|\mathbf{X}) &= \sigma^{i2} \mathbf{I}_n \end{aligned}$$

Above means i.e. each regression equation is Classical Linear Regression. However,

$$E(\mathbf{u}^j \mathbf{u}^{l'}|\mathbf{X}) = \sigma^{jl} \mathbf{I}_n$$

i.e.

$$\begin{aligned} E(u_t^j u_t^l|\mathbf{X}) &= \sigma^{jl} \quad \forall j \neq l \\ E(u_t^j u_s^j|\mathbf{X}) &= 0 \quad \forall t \neq s \end{aligned}$$

This relationship cannot be found if there is only one regression equation; thus, it is

usually viewed ‘seemingly unrelated.’ If we use GLS, we can estimate these estimators, which are more efficient than OLS estimator.

$$\Rightarrow \begin{pmatrix} \mathbf{y}^1 \\ \mathbf{y}^2 \\ \vdots \\ \mathbf{y}^m \end{pmatrix} = \begin{pmatrix} \mathbf{X}^1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{X}^2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{X}^m \end{pmatrix} \begin{pmatrix} \boldsymbol{\beta}^1 \\ \boldsymbol{\beta}^2 \\ \vdots \\ \boldsymbol{\beta}^m \end{pmatrix} + \begin{pmatrix} \mathbf{u}^1 \\ \mathbf{u}^2 \\ \vdots \\ \mathbf{u}^m \end{pmatrix}$$

$$\Rightarrow \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u} \quad \mathbf{y} \in \mathbb{R}^{mn}, \mathbf{X} \in \mathbb{R}^{mn \times \sum k_i}, \boldsymbol{\beta} \in \mathbb{R}^{\sum k_i}, \mathbf{u} \in \mathbb{R}^{mn}$$

Note that if \mathbf{y}^i have same dimensions for all $i \in \{1, 2, \dots, m\}$, the data can be called as a balanced panel. If not, it can be called as an unbalanced panel and using an unbalanced panel requires careful deliberation.

SAS is not good to analyze these *Seemingly Unrelated Regressions* with unbalanced panel data because it truncates its observation to make the data balanced. Thus, SAS can be recommended only if the data is balanced in terms of panel analysis. On the other hand, STATA provides more attractive way to analyze these regressions. In addition, EViews can do this, but it is not as powerful as STATA.

2.3.1. Generalized Least Squares Estimation for SUR

These non-standard error structure can be written as below.

$$E(\mathbf{u}\mathbf{u}'|\mathbf{X}) \neq \sigma^2 \mathbf{I}_{mn}$$

Therefore, GLS technique is required for this case.

$$\begin{aligned} E(\mathbf{u}\mathbf{u}'|\mathbf{X}) &= \boldsymbol{\Omega} \\ &= \begin{pmatrix} E(\mathbf{u}^1 \mathbf{u}^{1'}|\mathbf{X}) & E(\mathbf{u}^1 \mathbf{u}^{2'}|\mathbf{X}) & \cdots & E(\mathbf{u}^1 \mathbf{u}^{m'}|\mathbf{X}) \\ E(\mathbf{u}^2 \mathbf{u}^{1'}|\mathbf{X}) & E(\mathbf{u}^2 \mathbf{u}^{2'}|\mathbf{X}) & \cdots & E(\mathbf{u}^2 \mathbf{u}^{m'}|\mathbf{X}) \\ \vdots & \vdots & \ddots & \vdots \\ E(\mathbf{u}^m \mathbf{u}^{1'}|\mathbf{X}) & E(\mathbf{u}^m \mathbf{u}^{2'}|\mathbf{X}) & \cdots & E(\mathbf{u}^m \mathbf{u}^{m'}|\mathbf{X}) \end{pmatrix} \\ &= \begin{pmatrix} \sigma^{11} \mathbf{I}_n & \sigma^{12} \mathbf{I}_n & \cdots & \sigma^{1m} \mathbf{I}_n \\ \sigma^{21} \mathbf{I}_n & \sigma^{22} \mathbf{I}_n & \cdots & \sigma^{2m} \mathbf{I}_n \\ \vdots & \vdots & \ddots & \vdots \\ \sigma^{m1} \mathbf{I}_n & \sigma^{m2} \mathbf{I}_n & \cdots & \sigma^{mm} \mathbf{I}_n \end{pmatrix} \\ &= \begin{pmatrix} \sigma^{11} & \sigma^{12} & \cdots & \sigma^{1m} \\ \sigma^{21} & \sigma^{22} & \cdots & \sigma^{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma^{m1} & \sigma^{m2} & \cdots & \sigma^{mm} \end{pmatrix} \otimes \mathbf{I}_n \\ &= \boldsymbol{\Sigma} \otimes \mathbf{I}_n \quad \boldsymbol{\Sigma} \in \mathbb{R}^{m \times m} \end{aligned}$$

Remark *Kronecker* product \otimes can be expressed as below; for two matrices $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{p \times q}$,

$$\underset{(m \times n)}{\mathbf{A}} \otimes \underset{(p \times q)}{\mathbf{B}} = \underbrace{\begin{pmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \cdots & a_{2n}\mathbf{B} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & a_{m2}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{pmatrix}}_{(mp \times nq)}$$

Property

The inverse of a Kronecker product has below characteristic.

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{AC} \otimes \mathbf{BD}$$

$$\begin{aligned} \Rightarrow (\mathbf{A} \otimes \mathbf{B})(\mathbf{A}^{-1} \otimes \mathbf{B}^{-1}) &= \mathbf{AA}^{-1} \otimes \mathbf{BB}^{-1} \\ &= \mathbf{I}_{mn} \end{aligned}$$

$$\therefore (\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$$

Thus,

$$\mathbf{\Omega}^{-1} = \mathbf{\Sigma}^{-1} \otimes \mathbf{I}_n$$

$$\Rightarrow \hat{\boldsymbol{\beta}}_{\text{GLS}} = \left[\begin{pmatrix} \mathbf{X}^1 & \mathbf{0}' & \cdots & \mathbf{0}' \\ \mathbf{0}' & \mathbf{X}^2 & \cdots & \mathbf{0}' \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}' & \mathbf{0}' & \cdots & \mathbf{X}^m \end{pmatrix} (\mathbf{\Sigma}^{-1} \otimes \mathbf{I}_n) \begin{pmatrix} \mathbf{X}^1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{X}^2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{X}^m \end{pmatrix} \right]^{-1} \begin{pmatrix} \mathbf{X}^1 & \mathbf{0}' & \cdots & \mathbf{0}' \\ \mathbf{0}' & \mathbf{X}^2 & \cdots & \mathbf{0}' \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}' & \mathbf{0}' & \cdots & \mathbf{X}^m \end{pmatrix} (\mathbf{\Sigma}^{-1} \otimes \mathbf{I}_n) \begin{pmatrix} \mathbf{y}^1 \\ \mathbf{y}^2 \\ \vdots \\ \mathbf{y} \end{pmatrix}$$

Let

$$\mathbf{\Sigma}^{-1} = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1m} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{m1} & \sigma_{m2} & \cdots & \sigma_{mm} \end{pmatrix}$$

$$\Rightarrow \hat{\boldsymbol{\beta}}_{\text{GLS}} = \begin{pmatrix} \sigma_{11}\mathbf{X}^1\mathbf{X}^1 & \sigma_{12}\mathbf{X}^1\mathbf{X}^2 & \cdots & \sigma_{1m}\mathbf{X}^1\mathbf{X}^m \\ \sigma_{21}\mathbf{X}^2\mathbf{X}^1 & \sigma_{22}\mathbf{X}^2\mathbf{X}^2 & \cdots & \sigma_{2m}\mathbf{X}^2\mathbf{X}^m \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{m1}\mathbf{X}^m\mathbf{X}^1 & \sigma_{m2}\mathbf{X}^m\mathbf{X}^2 & \cdots & \sigma_{mm}\mathbf{X}^m\mathbf{X}^m \end{pmatrix}^{-1} \begin{pmatrix} \sum_{j=1}^m \sigma_{1j}\mathbf{X}^1\mathbf{y}^j \\ \sum_{j=1}^m \sigma_{2j}\mathbf{X}^2\mathbf{y}^j \\ \vdots \\ \sum_{j=1}^m \sigma_{mj}\mathbf{X}^m\mathbf{y}^j \end{pmatrix}$$

2.3.2. Relationship between GLS and OLS Estimators

When can we obtain same $\hat{\boldsymbol{\beta}}_{\text{GLS}}$ and $\hat{\boldsymbol{\beta}}_{\text{OLS}}$?

Case 1 $\sigma^{ij} = 0 \quad \forall i \neq j \Rightarrow \sigma_{ij} = 0 \quad \forall i \neq j$; then, $\sigma^{ij} = 0 \Leftrightarrow \sigma_{ij} = 0$; i.e. no correlation between observations.

$$\begin{aligned}
\Rightarrow \hat{\beta}_{\text{GLS}} &= \begin{pmatrix} \sigma_{11} \mathbf{X}^1{}' \mathbf{X}^1 & \sigma_{12} \mathbf{X}^1{}' \mathbf{X}^2 & \cdots & \sigma_{1m} \mathbf{X}^1{}' \mathbf{X}^m \\ \sigma_{21} \mathbf{X}^2{}' \mathbf{X}^1 & \sigma_{22} \mathbf{X}^2{}' \mathbf{X}^2 & \cdots & \sigma_{2m} \mathbf{X}^2{}' \mathbf{X}^m \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{m1} \mathbf{X}^m{}' \mathbf{X}^1 & \sigma_{m2} \mathbf{X}^m{}' \mathbf{X}^2 & \cdots & \sigma_{mm} \mathbf{X}^m{}' \mathbf{X}^m \end{pmatrix}^{-1} \begin{pmatrix} \sum_{j=1}^m \sigma_{1j} \mathbf{X}^1{}' \mathbf{y}^j \\ \sum_{j=1}^m \sigma_{2j} \mathbf{X}^2{}' \mathbf{y}^j \\ \vdots \\ \sum_{j=1}^m \sigma_{mj} \mathbf{X}^m{}' \mathbf{y}^j \end{pmatrix} \\
&= \begin{pmatrix} \sigma_{11} \mathbf{X}^1{}' \mathbf{X}^1 & 0 & \cdots & 0 \\ 0 & \sigma_{22} \mathbf{X}^2{}' \mathbf{X}^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_{mm} \mathbf{X}^m{}' \mathbf{X}^m \end{pmatrix}^{-1} \begin{pmatrix} \sigma_{11} \mathbf{X}^1{}' \mathbf{y}^1 \\ \sigma_{22} \mathbf{X}^2{}' \mathbf{y}^2 \\ \vdots \\ \sigma_{mm} \mathbf{X}^m{}' \mathbf{y}^m \end{pmatrix} \\
&= \hat{\beta}_{\text{OLS}}
\end{aligned}$$

In this case, GLS estimator is directly equivalent to OLS estimator.

Case 2 $\mathbf{X}^1 = \mathbf{X}^2 = \cdots = \mathbf{X}^m \quad (\equiv \mathbf{X}^C)$

Above means that every regressions exploit exactly same independent variables. In this case, also, it is natural that GLS estimator and OLS estimator are equivalent.

Lecture Note 3

2014-04-02

3.1. Binary Choice Model

Definition I_i is a dependent variable. ($\forall i \in \{1, 2, \dots, n\}$) Where

$$I_i := \begin{cases} 0 & \text{with probability } 1-p \\ 1 & \text{with probability } p \end{cases}$$

Definition x_i is an independent variable.

With these given dependent and independent variables, what we want to do can be denoted as below.

$$I_i = x_i' \beta + \varepsilon_i$$

Figure 3.1 Ordinary Dependent Variable

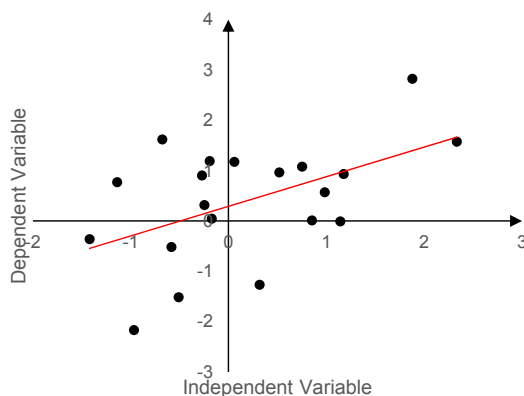


Figure 3.2 Binomial Dependent Variable

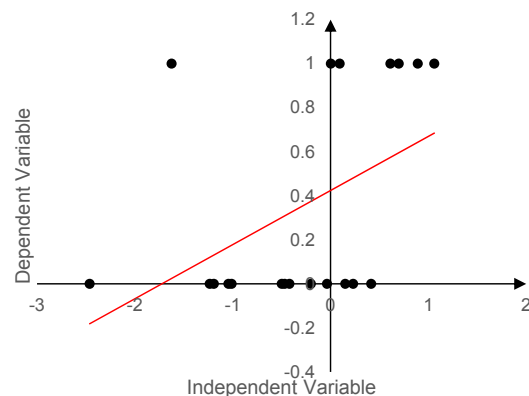


Figure 3.3 Logit Model (or Probit Model)

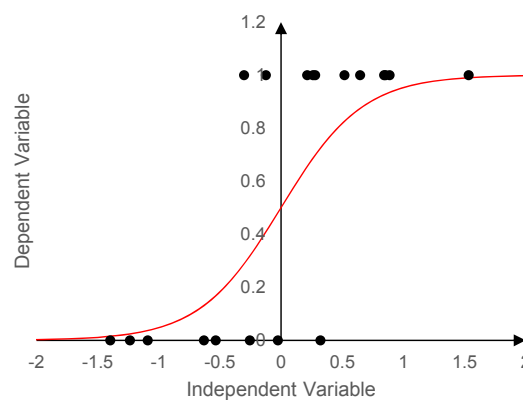


Figure 3.1 shows classical linear regression line. However, if the dependent variable is

binomially distributed, classical linear regression will not be good. Figure 3.2 graphically shows that linear regression can be wrong for this case. Fortunately, we can solve this problem based on Logit model (or Probit model if possible) that is usually adopted in binomial dependent variable case.

$$I_i = F(\mathbf{x}_i' \boldsymbol{\beta}) + \varepsilon_i$$

Where F is a cumulative distribution function, which has below properties.

$$\begin{aligned} F &\in [0, 1] \\ \frac{dF(z)}{dz} &\geq 0 \end{aligned}$$

Note If I_i follows Bernoulli distribution as below, then

$$I_i := \begin{cases} 0 & f_1 = p \\ 1 & f_2 = 1 - p \end{cases}$$

$$\begin{aligned} \Rightarrow E(I_i) &= p \\ \text{Var}(I_i) &= p(1 - p) \end{aligned}$$

$$\begin{aligned} \therefore E(I_i | \mathbf{x}_i) &= F(\mathbf{x}_i' \boldsymbol{\beta}) \\ &= P(I_i = 1 | \mathbf{x}_i) \end{aligned}$$

If F is Standard Normal, then Probit model can be used.

$$F = \Phi(\mathbf{x}_i' \boldsymbol{\beta})$$

If F is a logistic distribution function, then Logit model can be used.

$$F(\mathbf{x}_i' \boldsymbol{\beta}) = \frac{e^{\mathbf{x}_i' \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i' \boldsymbol{\beta}}}$$

Note F is bounded with 0 and 1.

If $\mathbf{x}_i > \mathbf{0}$, then

$$\begin{aligned} \lim_{\boldsymbol{\beta} \rightarrow \infty} F(\mathbf{x}_i' \boldsymbol{\beta}) &= \lim_{\boldsymbol{\beta} \rightarrow \infty} \frac{e^{\mathbf{x}_i' \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i' \boldsymbol{\beta}}} \\ &= \lim_{\boldsymbol{\beta} \rightarrow \infty} \frac{e^{\mathbf{x}_i' \boldsymbol{\beta}}}{e^{\mathbf{x}_i' \boldsymbol{\beta}}} \quad \text{by l'Hôpital's rule} \\ &= 1 \end{aligned}$$

Similarly,

$$\begin{aligned}
 \lim_{\beta \rightarrow -\infty} F(\mathbf{x}_i' \beta) &= \lim_{\beta \rightarrow -\infty} \frac{e^{\mathbf{x}_i' \beta}}{1 + e^{\mathbf{x}_i' \beta}} \\
 &= \frac{0}{1 + 0} \\
 &= 0
 \end{aligned}$$

Figure 3.4 The Use of Logit and Probit Model

