

ECON21



Introductory Econometrics

讲义
原

Quiz: Quiz 3 - 4

Quiz 6 - 3

Quiz 10 - 4

21-3 Introduction and Statistics Review

1. Basics of econometrics

① Definition

Analysis of economic data with statistical methods

② Steps of doing econometric research

计量 → 实证分析

▫ Steps in an empirical analysis

- Step 1: Carefully pose a question.
- Step 2: Specify an economic or conceptual model.
- Step 3: Turn the economic model into an econometric model.
- Step 4: Collect data on the variables and use statistical methods to estimate the parameters, construct confidence intervals for the parameters, and test hypotheses.

③ Economic data sets

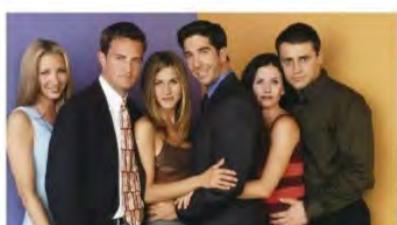
▪ Cross-sectional data 横向数据 represent a random sample

▪ Time series data 时间序列数据

▪ Pooled cross sections 混合横断面数据 cross sections are drawn independently

▪ Panel/Longitudinal data 时序数据

Cross Sectional Data



2003

Time Series Data



1993

2003

2013

Pooled Cross Sectional Data



1993

2003

2013

Panel Data



④ Causality & Ceteris paribus

Relationship between x & y $\xrightarrow{\textcircled{X}}$ Causality
Relationship between x & $y + c$ $\xrightarrow{\textcircled{Y}}$ Causality

Control other relevant factors the same \Rightarrow Make x independent with other relevant factors \Rightarrow Random sampling
 $(\Delta u = 0)$ (Zero conditional mean)

2. Statistics review

① Variance

$$\text{Var}(ax+by) = a^2 \text{Var}(x) + b^2 \text{Var}(y) + 2ab \text{Cov}(x, y)$$

Recall: Risk = $\sqrt{\text{Var}(w_1x_1 + \dots + w_nx_n)} = \text{Std}(w_1x_1 + \dots + w_nx_n)$

② Covariance

a. Sample covariance

$$\text{Sample Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1} = \frac{1}{n-1} [\sum x_i y_i - \frac{\sum x_i \bar{y}}{n}]$$

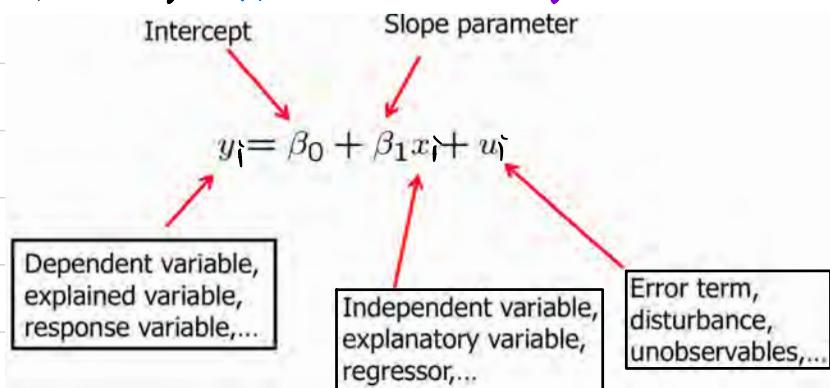
b. Two properties

$$\text{Cov}(ax, by) = ab \text{Cov}(x, y)$$

$$\text{Cov}(ax+bx, cy+dz) = ac \text{Cov}(x, y) + ad \text{Cov}(x, z) + bc \text{Cov}(y, x) + bd \text{Cov}(y, z)$$

4-6 Simple Linear Regression Model (Bivariate Linear Regression Model)

1. Actual case (散点图) ← 同个 x_i 有各个 y_i :



2. Average: Population regression function (PRF)

① Zero conditional mean assumption

$$(u \& y: E(u) = 0 \text{ (by defining } \beta_0\text{)})$$

$$(u \& x: E(u|x) = E(u))$$

通过调整来达到

$$\forall x_i, E(u|x_i) = 0$$

$$\Rightarrow E(u|x) = 0$$

② PRF

$$E(\beta_1 x | x) = \beta_1 x$$

$$E(y|x) = E(\beta_0 + \beta_1 x + u|x) = \beta_0 + \beta_1 x + E(u|x) = \beta_0 + \beta_1 x$$

$$\leftarrow y|x \sim N(\beta_0 + \beta_1 x, \sigma^2)$$

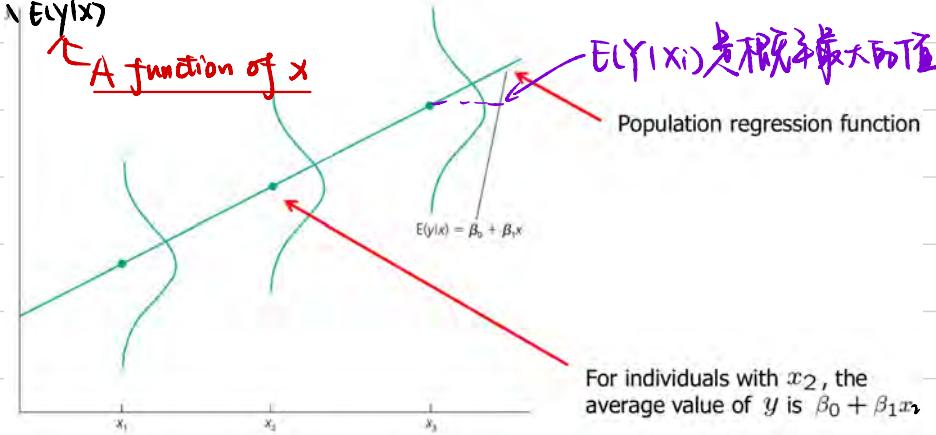


参数=我们想用的 SLR

解释 Causality relationship

Significance: The average value of y given x can be expressed as a linear function of x .

③ $E(y|x)$



3. Prediction: Sample regression function / Ordinary Least Squares regression function (OLS regression function)

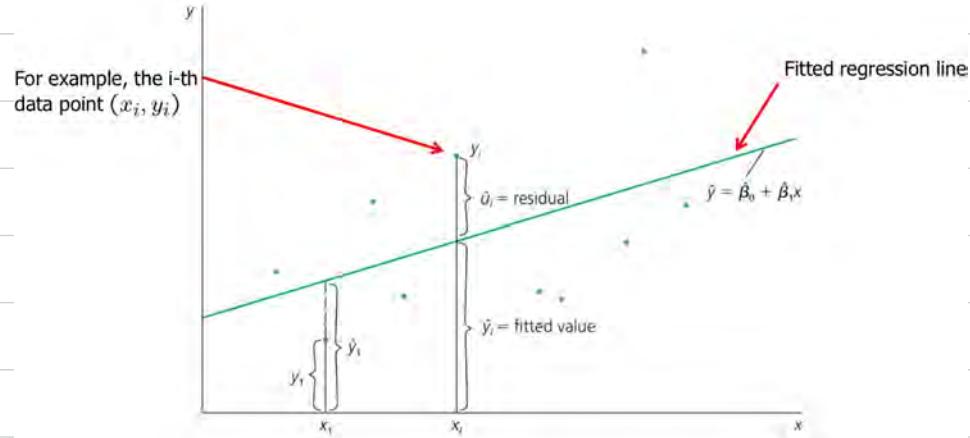
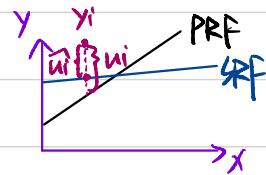
① Overview

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

The mistake is called "residual": $u_i = y_i - \hat{y}_i$

$$u_i = y_i - E(y|x_i)$$

$$\hat{u}_i \text{ (Residual)} = y_i - \hat{y}_i$$



\hat{x}_j : Holding other factors fixed, ...

② OLS (Point estimation)

a. Process

$$\min \sum_{i=1}^n \hat{u}_i^2 \rightarrow \min \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \rightarrow \text{FOC: } \begin{cases} \hat{\beta}_0: \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \\ \hat{\beta}_1: \sum x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \end{cases}$$

$$\rightarrow \begin{cases} \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} = \text{Corr}(x, y) \frac{s_y}{s_x}, \text{ where } \begin{cases} s_x = \sqrt{\sum (x_i - \bar{x})^2} \\ s_y = \sqrt{\sum (y_i - \bar{y})^2} \end{cases} \\ \hat{\sigma}_e^2 = \frac{\sum \hat{u}_i^2}{n-2} \end{cases}$$

$n \rightarrow n-2$: Biased \rightarrow Unbiased

若 $\hat{u}_i \sim N(0, \sigma_e^2)$ (MM, 课本 P49)

$$E(u) = 0$$

$$\text{Cov}(x, u) = E(xu) = 0$$

b. Basic properties

i) $E(\hat{u}_i) = 0 \quad \sum_{i=1}^n \hat{u}_i = 0$

FOC of $\hat{\beta}_0 \rightarrow \sum (y_i - \hat{y}_i) = 0 \rightarrow \sum \hat{u}_i = 0$

ii) $E(\hat{u}_i | x_i) = 0 \quad \sum x_i \hat{u}_i = 0$

FOC of $\hat{\beta}_1 \rightarrow \sum x_i (y_i - \hat{y}_i) = 0 \rightarrow \sum x_i \hat{u}_i = 0$

iii) $(x, y) \text{ 独立} \rightarrow \bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}$

$E(\hat{f}(x) \hat{u}) = 0$

TRIVIAL: $\sum x_i$

Implications

i) $\bar{y} = \bar{\hat{y}}$ $y_i = \hat{y}_i + \hat{u}_i \rightarrow \sum y_i = \sum \hat{y}_i + \sum \hat{u}_i \rightarrow \bar{y} = \bar{\hat{y}}$

ii) $\sum \hat{u}_i = 0 \quad \sum x_i \hat{u}_i = 0 \rightarrow (\hat{\beta}_0 + \hat{\beta}_1 x_i) \hat{u}_i = \sum \hat{y}_i \hat{u}_i = 0$

iii) $\sum \hat{y}_i \hat{u}_i = 0$

iv) $\sum y_i \hat{u}_i \neq 0$

$\sum y_i \hat{u}_i = \sum (\hat{y}_i + \hat{u}_i) \hat{u}_i = \sum \hat{u}_i^2$

c. Statistical properties

5 Assumptions (Gauss-Markov assumptions)

- a. Linear relation: $y_i = \beta_0 + \beta_1 x_i + u_i$
- b. Random sampling
- c. Existence of sample variation: $\sum (x_i - \bar{x})^2 > 0 \Rightarrow \hat{\beta}_1 \neq \bar{y}/\bar{x}; \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$
- d. Zero conditional mean: $\forall i, E(u_i | x_i) = 0$
- e. Constant variance: $\text{Var}(u_i | x_i) = \sigma^2$
(Homoskedasticity) \leftarrow (Heteroskedasticity)

ホンセキ

アヘンスドシテスリ

ホンバフ

$\Delta E(\hat{\beta}_1)$

$$E(\hat{\beta}_1) = \beta_1 \text{ (Unbiased)}$$

Proof:

$$\begin{aligned} E(\hat{\beta}_1) &= E\left(\frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}\right) = E\left(\frac{\sum(x_i - \bar{x})y_i}{\sum(x_i - \bar{x})^2}\right) = E\left(\frac{\sum(x_i - \bar{x})(\beta_0 + \beta_1 x_i + u_i)}{\sum(x_i - \bar{x})^2}\right) = E\left(\frac{\beta_0 \sum(x_i - \bar{x}) + \beta_1 \sum(x_i - \bar{x})x_i + \sum u_i(x_i - \bar{x})}{\sum(x_i - \bar{x})^2}\right) \\ &= E\left(\frac{0 + \beta_1 \sum(x_i - \bar{x}) + \sum u_i(x_i - \bar{x})}{\sum(x_i - \bar{x})^2}\right) = E(\beta_1) \frac{\sum x_i(x_i - \bar{x})}{\sum(x_i - \bar{x})^2} + E\left(\frac{\sum u_i(x_i - \bar{x})}{\sum(x_i - \bar{x})^2}\right) = E(\beta_1) \frac{\sum x_i(x_i - \bar{x})}{\sum(x_i - \bar{x})(x_i - \bar{x})} + E\left(\frac{\sum u_i x_i - \sum u_i \bar{x}}{\sum(x_i - \bar{x})^2}\right) \\ &= E(\beta_1) \frac{\sum x_i(x_i - \bar{x})}{\sum x_i(x_i - \bar{x}) - \sum \bar{x}(x_i - \bar{x})} + E\left(\frac{\sum u_i x_i}{\sum(x_i - \bar{x})^2}\right) = \beta_1 + 0 = \beta_1 \end{aligned}$$

$\Delta E(\hat{\beta}_0)$

$$E(\hat{\beta}_0) = \beta_0 \text{ (Unbiased)}$$

Proof:

$$E(\hat{\beta}_0) = E(\bar{y} - \hat{\beta}_1 \bar{x}) = y - \beta_1 x = \beta_0$$

$\Delta \text{Var}(\hat{\beta}_1)$

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum(x_i - \bar{x})^2}$$

Proof:

$$\text{Var}(\hat{\beta}_1) = \text{Var}\left(\frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}\right) = \text{Var}\left(\frac{\sum(x_i - \bar{x})y_i}{\sum(x_i - \bar{x})^2}\right) = \frac{\sum(x_i - \bar{x})^2 \text{Var}(y_i)}{(\sum(x_i - \bar{x})^2)^2} = \frac{1}{\sum(x_i - \bar{x})^2} \text{Var}(y_i) = \frac{\sigma^2}{\sum(x_i - \bar{x})^2}$$

Logic: 在“Conditioned on x_i ”的前提下算出 $E(\text{Var})$

后发现与 x_1, x_2, \dots, x_n 无关，所以去掉后项也成立
Conditioned on x_i $\Rightarrow x_i$ are constant

$\Delta \text{Var}(\hat{\beta}_0)$

$$\text{Var}(\hat{\beta}_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2} \right)$$

Proof:

$$\begin{aligned} \text{Var}(\hat{\beta}_0) &= \text{Var}(\bar{y} - \hat{\beta}_1 \bar{x}) = \text{Var}(\bar{y}) + \bar{x}^2 \frac{\sigma^2}{\sum(x_i - \bar{x})^2} - 2 \text{Cov}(\bar{y}, \hat{\beta}_1 \bar{x}) \\ &= \text{Var}\left(\frac{\sum y_i}{n}\right) + \bar{x}^2 \frac{\sigma^2}{\sum(x_i - \bar{x})^2} - 2\bar{x} \text{Cov}\left(\frac{\sum y_i}{n}, \frac{\sum(x_i - \bar{x})y_i}{\sum(x_i - \bar{x})^2}\right) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2} \right) + \frac{2\bar{x}}{n \sum(x_i - \bar{x})^2} \text{Cov}(\sum y_i, \sum(x_i - \bar{x})y_i) \\ &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2} \right) + \frac{2\bar{x}}{n \sum(x_i - \bar{x})^2} \sum(x_i - \bar{x}) \text{Cov}(y_i, y_i) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2} \right) + \frac{2\bar{x}}{n \sum(x_i - \bar{x})^2} \sum(x_i - \bar{x}) \sigma_i = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2} \right) \end{aligned}$$

↑
 $\sigma^2 \uparrow$, $\text{Var}(\hat{\beta}_1) \uparrow$ (Mistake 越大, 模型参数越差)
↑
 $n \uparrow$, $\text{Var}(\hat{\beta}_1) \uparrow$ (样本越大, 则不线性)
 $\bar{x} \downarrow$, $\text{Var}(\hat{\beta}_1) \uparrow$ (样本越集中, 则不线性)

$$\Delta E(\hat{\beta}^2) = \sigma^2 \text{ (当 n 趋近于 } \infty \text{ 时, 就是 unbiased)}$$

$$\text{sd}(\hat{\beta}) \Rightarrow \text{se}(\hat{\beta})$$

Standard error of the regression: $\hat{\sigma}_{(-Y)}$ \longrightarrow Standard error of $\hat{\beta}$

$E(\hat{\sigma}^2) = \sigma^2 \checkmark$ Unbiased

(Standard error) $\hat{\sigma}$ Estimator

$$E(\hat{\sigma}) = \sigma \times \text{Biased}$$

③ Goodness of fit

a. Variation

$$y_i = \hat{y}_i + \hat{u}_i$$

$$\rightarrow y_i - \bar{y} = \hat{y}_i - \bar{y} + \hat{u}_i$$

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

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

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

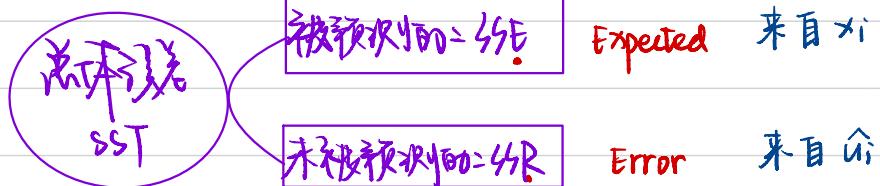
$$= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (\hat{u}_i)^2$$

$$SST = \left(\begin{array}{c} y_i \\ \hat{y}_i \\ \bar{y} \end{array} \right) \text{ SSR}$$

$$SSE = \left(\begin{array}{c} \hat{u}_i \end{array} \right)$$

↑ 描述 y_i 的波动

$$\rightarrow SST = SSE + SSR, \text{ where } \begin{cases} SST = \sum_{i=1}^n (y_i - \bar{y})^2 & \text{Total variation } y \\ SSE = \sum_{i=1}^n (\hat{u}_i)^2 & \text{Explained part } x \\ SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 & \text{Unexpected part } u \end{cases}$$



b. R-square

$$R^2 = \frac{SSE}{SST} = 1 - \frac{SSR}{SST}$$

$$R^2 \in [0, 1], R^2 \uparrow, \text{Goodness of fit} \uparrow$$

(注意, R^2 只适用于相同 model 下)

R^2 小, 说明预测性差

R^2 小, 不能说明线性模型差 (判断线性模型好坏的依据是是否满足 5 个 assumptions)

g. Alternative expression of R-squared

$$\begin{aligned} R^2 &= \frac{SSE}{SST} \\ &= \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} \\ &= \frac{[\sum (y_i - \bar{y})^2]^2}{\sum (y_i - \bar{y})^2 \cdot \sum (\hat{y}_i - \bar{y})^2} \\ &\quad \downarrow \sum (\hat{y}_i - \bar{y})^2 \\ &= \sum (y_i + \hat{u}_i - \bar{y})(\hat{y}_i - \bar{y}) \\ &= \sum (y_i + \hat{u}_i)(\hat{y}_i - \bar{y}) - \sum \bar{y}(\hat{y}_i - \bar{y}) \\ &= \sum (y_i - \bar{y})(\hat{y}_i - \bar{y}) + \sum \hat{u}_i(\hat{y}_i - \bar{y}) \\ \therefore R^2 &= \frac{\sum (y_i - \bar{y})(\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2 \cdot \sum (\hat{y}_i - \bar{y})^2} \\ &= \text{corr}(y_i, \hat{y}_i)^2 \end{aligned}$$

④ Change of units

a. $\Delta x \rightarrow \% \Delta x$ (Change in %dec)

$$\Delta x = \frac{\% \Delta x}{100}$$

Intercept 不變
 Slope 變
 R^2 不變

percentage change

b. $\% \Delta x \rightarrow \Delta \ln(x)$

$$\Delta \ln(x) \approx \frac{\% \Delta x}{100}$$

Intercept 變
 Slope 變
 R^2 變

proof. Since $x \rightarrow 0$, $\ln(x+1) \rightarrow x$, then $\Delta \ln(x) = \ln(x_0 + \Delta x) - \ln(x_0) = \ln(1 + \frac{\Delta x}{x_0}) \approx \frac{\Delta x}{x_0} = \frac{\% \Delta x}{100}$

c. Summary of forms of variables (因为都是 "y對x = $\beta_0 + \beta_1 x^n$ " 的形式, 所以是 LR)

Model	Dep. Var.	Indep. Var.	β_1
Level-Level	y	x	$\Delta y = \beta_1 \Delta x$
Level-Log	y	$\log(x)$	$\Delta y = (\beta_1 / 100) \% \Delta x$
Log-Level	$\log(y)$	x	$\% \Delta y = (100 \beta_1) \Delta x$
Log-Log	$\log(y)$	$\log(x)$	$\% \Delta y = \beta_1 \% \Delta x$

← "Constant elasticity model"

4. Regression through the origin ($\beta_0 = 0$)

① Regression model (ASSUMPTION)

$$y = \beta_0 + \beta_1 x + u$$

② Sample regression function

a. Overview

$$\tilde{y} = \tilde{\beta}_1 x$$

b. Point estimate

$$\hat{\beta}_1 = \frac{\sum x_i y_i}{\sum x_i^2}$$

c. Statistical properties

$$\Delta E(\tilde{\beta}_1)$$

$$E(\tilde{\beta}_1) = \beta_0 \frac{\sum x_i}{\sum x_i^2} + \beta_1 \text{ (biased)}$$

proof:

$$E(\tilde{\beta}_1) = E\left(\frac{\sum x_i y_i}{\sum x_i^2}\right) = E\left(\frac{\sum x_i(\beta_0 + \beta_1 x_i + u_i)}{\sum x_i^2}\right) = E\left(\frac{\sum x_i \beta_0}{\sum x_i^2}\right) + E\left(\frac{\sum x_i \beta_1 x_i}{\sum x_i^2}\right) + E\left(\frac{\sum x_i u_i}{\sum x_i^2}\right) = \beta_0 \left(\frac{\sum x_i}{\sum x_i^2}\right) + \beta_1 + 0$$

$$= \beta_0 \frac{\sum x_i}{\sum x_i^2} + \beta_1 \neq \beta_1$$

$$\Delta \text{Var}(\tilde{\beta}_1)$$

$$\text{Var}(\tilde{\beta}_1) = \frac{\sigma^2}{\sum x_i^2}$$

proof:

$$\text{Var}(\tilde{\beta}_1) = \text{Var}\left(\frac{\sum x_i y_i}{\sum x_i^2}\right) = \frac{\sum x_i^2}{(\sum x_i^2)^2} \sigma^2 = \frac{\sigma^2}{\sum x_i^2}$$

d. Goodness of fit

$$R^2 = 1 - \frac{\sum (y_i - \tilde{\beta}_1 x_i)^2}{\sum (y_i - \bar{y})^2} \quad \bar{y} = 0, \text{ and when } y_i \text{ & } \tilde{y}_i \text{ have different signs, then } R^2 < 0.$$

At this time, it means that regression effectiveness: $\bar{y} > \tilde{y} > \tilde{y}$

指正 y_i 的不 \bar{y}

L7-10 Multiple Linear Regression Model — Estimation & Properties

1. Motivation

① Ceteris paribus

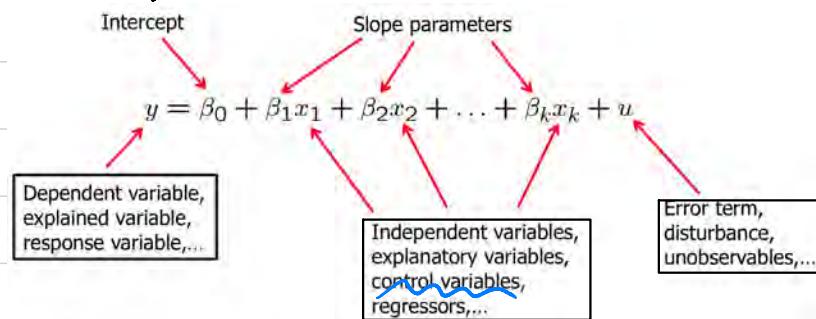
SLR 假设是 x_i 与 u 不相关 (CP 或 II, i.e., $E(u|x_i) = E(u)$)。

因此, 我们引进 x_i 的变量, 从假设是 x_i 与 u 不相关 (CP 或 II), 研究 causality $x_1 \rightarrow x_2 \rightarrow \dots$ 观察中的变化

② Function form

MLR (线性函数) 模型为, e.g. $\begin{cases} x = x \\ y = x^2 \end{cases}$ (单个多元是取决于 p_i 的两个数)
线性的是 p_i

2. Actual Case



3. Prediction = Sample regression function

① Overview

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 \quad \text{For } i \quad \begin{matrix} x_{i1} \\ x_{i2} \end{matrix}, \hat{u}_i = y_i - \hat{y}_i$$

(当然, 不止于此形式)

② OLS (Point estimation)

a. Process

$$\min \sum_{i=1}^n \hat{u}_i^2 \rightarrow \min \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 x_{i2})^2 \rightarrow \text{FOC: } \begin{cases} \hat{\beta}_0: \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 x_{i2}) = 0 \\ \hat{\beta}_1: \sum x_{i1} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 x_{i2}) = 0 \\ \hat{\beta}_2: \sum x_{i2} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 x_{i2}) = 0 \end{cases} \quad \begin{matrix} \hat{\beta}_j \text{ 表示 } \text{Partial effect} \\ \text{Holding other independent variables fixed} \end{matrix}$$

We don't need to compute the exact form of $\hat{\beta}_j$, but we need to master "Partialling out" Interpretation form of $\hat{\beta}_j$

△ Reg x_{i1} on x_{i2} , find residual \hat{r}_{i1}

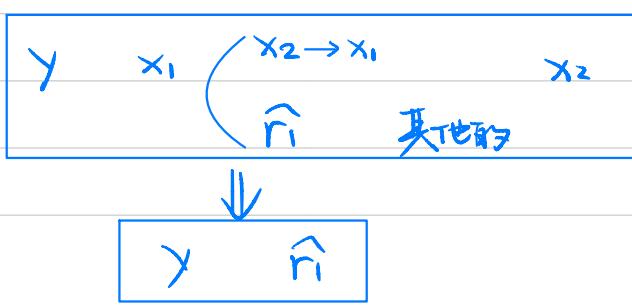
因子部分, 而非整体

△ Reg y_i on \hat{r}_{i1} , find $\hat{\beta}_1$

$$\hat{\beta}_1 = \frac{\sum (\hat{r}_{i1} - \bar{\hat{r}}_{i1})(y_i - \bar{y}_i)}{\sum (\hat{r}_{i1} - \bar{\hat{r}}_{i1})^2} = \frac{\sum (\hat{r}_{i1} - \bar{\hat{r}}_{i1})(y_i - \bar{y}_i)}{\sum (\hat{r}_{i1} - \bar{\hat{r}}_{i1})^2}$$

$$\hat{\beta}_1 = \frac{\sum (\hat{r}_{i1} - \bar{\hat{r}}_{i1})(y_i - \bar{y}_i)}{\sum (\hat{r}_{i1} - \bar{\hat{r}}_{i1})^2} \quad (= \frac{\sum (\hat{r}_{i1} - 0)(y_i - \bar{y}_i)}{\sum (\hat{r}_{i1} - 0)^2})$$

$$\hat{\sigma}^2 = \frac{\sum \hat{u}_i^2}{n-(k+1)} \quad (\text{因为 } \hat{u}_i = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 x_{i2})$$



b. Basic properties

$$(i) E(\hat{u}_i) = 0 \quad \sum_{i=1}^n \hat{u}_i = 0$$

$$(ii) E(u_i | x_1, \dots, x_n) = 0 \quad \sum_{i=1}^n x_{ij} \hat{u}_i = 0$$

$$(iii) (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_k) \text{ WRF 上 } \bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}_1 + \dots + \hat{\beta}_k \bar{x}_k$$

Implications

$$\left. \begin{array}{l} (i) \bar{y} = \bar{\bar{y}} \quad y_i = \bar{y} + \hat{u}_i \rightarrow \frac{1}{n} \sum y_i = \frac{1}{n} \sum \bar{y} + \frac{1}{n} \sum \hat{u}_i \rightarrow \bar{y} = \bar{\bar{y}} \\ (ii) \sum_{i=1}^n \hat{u}_i = 0 \quad \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_k x_{ik}) \hat{u}_i = 0 + 0 + \dots + 0 = 0 \end{array} \right\}$$

c. Statistical properties

5 Assumptions (Gauss-Markov assumptions)

a. Linear relation : $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + u_i$

b. Random sampling

c. No perfect collinearity : Δ Existence of sample variation : For $\forall j$, $\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 > 0$
 Δ No exact linear relationship between independent variables, e.g. $x_1 = 2x_2 + 3x_3$ x , $x_1 = 2x_2 + 3x_3$ v

d. Zero conditional mean : $\forall i$, $E(u_i | x_{i1}, \dots, x_{in}) = 0$ $\xrightarrow{\text{Exogeneity assumption}}$ $\xrightarrow{\text{NO ENDGENOUS}}$

Endogenous : x_j is correlated with u_i

Exogenous : x_j is uncorrelated with u_i

e. Constant variance : $\forall i$, $\text{Var}(u_i | x_{i1}, \dots, x_{in}) = \sigma^2$

(Homoskedasticity)

$\Delta E(\hat{\beta}_j)$ (Normal cases)

$$E(\hat{\beta}_j) = \beta_j$$

proof:

$$\begin{aligned} E(\hat{\beta}_j) &= E\left(\frac{\sum_{i=1}^n (\hat{r}_{ij} y_i)}{\sum_{i=1}^n \hat{r}_{ij}^2}\right) = E\left(\frac{\sum_{i=1}^n \hat{r}_{ij}(\beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + u_i)}{\sum_{i=1}^n \hat{r}_{ij}^2}\right) \\ &\quad \xrightarrow{\text{由 } \hat{r}_{ij} \text{ 是 uncorrelated with } \hat{r}_{ij} \text{ 的线性组合}} E\left(0 + 0 + \dots + \frac{\sum_{i=1}^n \hat{r}_{ij} \beta_j x_{ij}}{\sum_{i=1}^n \hat{r}_{ij}^2} + \dots + 0 + \frac{\sum_{i=1}^n \hat{r}_{ij} u_i}{\sum_{i=1}^n \hat{r}_{ij}^2}\right) \\ &= E\left(\frac{\sum_{i=1}^n \hat{r}_{ij} \beta_j x_{ij} + \sum_{i=1}^n \hat{r}_{ij} u_i}{\sum_{i=1}^n \hat{r}_{ij}^2}\right) = E\left(\frac{\sum_{i=1}^n \hat{r}_{ij} x_{ij} \beta_j}{\sum_{i=1}^n \hat{r}_{ij}^2}\right) + E\left(\frac{\sum_{i=1}^n \hat{r}_{ij} u_i}{\sum_{i=1}^n \hat{r}_{ij}^2}\right) = E\left(\frac{\sum_{i=1}^n \hat{r}_{ij} x_{ij} \beta_j}{\sum_{i=1}^n \hat{r}_{ij}^2}\right) + 0 \\ &= E\left(\frac{\sum_{i=1}^n \hat{r}_{ij}^2 \beta_j}{\sum_{i=1}^n \hat{r}_{ij}^2}\right) + E\left(\frac{\sum_{i=1}^n \hat{r}_{ij} x_{ij} \beta_j}{\sum_{i=1}^n \hat{r}_{ij}^2}\right) = \beta_j + 0 = \beta_j \end{aligned}$$

A Function of x_i

$\Delta E(\hat{\beta}_j)$ (Overspecifying case) R.P.: 包含无关变量 x_j

$$E(\hat{\beta}_j) = 0$$

proof:

$$\begin{aligned} E(\hat{\beta}_j) &= \beta_j = 0 \quad \text{注意: } (\beta_j = 0) \\ &\quad (\hat{\beta}_j \neq 0 \text{ (不一定)}) \end{aligned}$$

$\Delta E(\hat{\beta}_j)$ (Underspecifying/Omitting cases)

Exp1: Relevant variable x_2 is omitted

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + u_i$$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1}$$

$$E(\hat{\beta}_1) = \beta_1 + \beta_2 \tilde{s}_1, \quad \text{实际上, } \tilde{s}_1 = \hat{\beta}_1 + \hat{\beta}_2 \tilde{\delta}_1$$

where \tilde{s}_1 is the estimated slope in SLR where reg x_2 on x_1

$E(\hat{\beta}_1)$ (来自 β_1) \rightarrow Type of Bias (Upward / Positive) \rightarrow Downward / Negative

proof:

$$\begin{aligned} E(\hat{\beta}_1) &= E\left(\frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1)(y_i - \bar{y}_i)}{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2}\right) = E\left(\frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1)(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + u_i)}{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2}\right) \\ &= E\left(\beta_1 + \beta_2 \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1)x_{i2}}{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2}\right) \geq \beta_1 + \beta_2 \cdot \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)}{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2} \end{aligned}$$

$$> \beta_1 + \beta_2 \tilde{s}_1$$

$$\tilde{s}_1$$

Bias 的正负

$$\beta_2 x_2 \text{ 与 } \beta_0 \text{ 的关系}$$

$$\tilde{s}_1 \text{ 与 } x_2 \text{ 与 } x_1 \text{ 的关系}$$

Exp2: Relevant variable x_3 is omitted

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + u_i$$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2}$$

Similarly, we get

$$E(\hat{\beta}_1) = \beta_1 + \beta_3 \tilde{s}_1$$

$$\text{where } \tilde{s}_1 = \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1)(x_{i3} - \bar{x}_3)}{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2}$$

proof:

$$\begin{aligned} E(\hat{\beta}_1) &= E\left(\frac{\sum_{i=1}^n (\hat{r}_{i1} y_i)}{\sum_{i=1}^n (\hat{r}_{i1})^2}\right) = E\left(\frac{\sum_{i=1}^n \hat{r}_{i1}(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + u_i)}{\sum_{i=1}^n (\hat{r}_{i1})^2}\right), \\ &= \beta_1 + \frac{\sum_{i=1}^n \hat{r}_{i1} \hat{r}_{i3} \beta_3}{\sum_{i=1}^n (\hat{r}_{i1})^2} + E\left(\frac{\sum_{i=1}^n \hat{r}_{i1} x_{i1} \beta_1}{\sum_{i=1}^n (\hat{r}_{i1})^2}\right) + E\left(\frac{\sum_{i=1}^n \hat{r}_{i1} \beta_2 x_{i2}}{\sum_{i=1}^n (\hat{r}_{i1})^2}\right) \\ &= \beta_1 + \frac{\sum_{i=1}^n \hat{r}_{i1} \hat{r}_{i3} \beta_3}{\sum_{i=1}^n (\hat{r}_{i1})^2} + E\left(\frac{\sum_{i=1}^n \hat{r}_{i1} (\hat{r}_{i1} + \tilde{s}_1)}{\sum_{i=1}^n (\hat{r}_{i1})^2}\right) + \beta_3 \frac{\sum_{i=1}^n \hat{r}_{i1} x_{i3}}{\sum_{i=1}^n (\hat{r}_{i1})^2} \\ &= \beta_1 + \beta_3 \frac{\sum_{i=1}^n \hat{r}_{i1} x_{i3}}{\sum_{i=1}^n (\hat{r}_{i1})^2} \\ &= \beta_1 + \beta_3 \tilde{s}_1 \end{aligned}$$

A) $\text{Var}(\hat{\beta}_j)$ (Normal cases)

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 (1 - R_j^2)}$$

R_j^2 reg x_j on other independent variables for Goodness of fit

proof:

$$\text{Var}(\hat{\beta}_j) = \text{Var}\left(\frac{\sum_{i=1}^n \hat{r}_{ij} y_i}{\sum_{i=1}^n \hat{r}_{ij}^2}\right) = \text{Var}\left(\frac{\sum_{i=1}^n \hat{r}_{ij} (\beta_0 + \beta_1 x_{i1} + \dots + \beta_n x_{in} + u_i)}{\sum_{i=1}^n \hat{r}_{ij}^2}\right) = \text{Var}\left(\frac{\sum_{i=1}^n \hat{r}_{ij} \hat{\beta}_j x_{ij} + \sum_{i=1}^n \hat{r}_{ij} u_i}{\sum_{i=1}^n \hat{r}_{ij}^2}\right)$$

$$> \text{Var}\left(\frac{\sum_{i=1}^n \hat{r}_{ij} \hat{\beta}_j (\hat{x}_j + \hat{r}_{ij}) + \sum_{i=1}^n \hat{r}_{ij} u_i}{\sum_{i=1}^n \hat{r}_{ij}^2}\right) > \text{Var}(\hat{\beta}_j + \frac{\sum_{i=1}^n \hat{r}_{ij} u_i}{\sum_{i=1}^n \hat{r}_{ij}^2})$$

$$> \text{Var}(\frac{\sum_{i=1}^n \hat{r}_{ij} u_i}{\sum_{i=1}^n \hat{r}_{ij}^2}) = \frac{\sum_{i=1}^n (\hat{r}_{ij})^2}{(\sum_{i=1}^n \hat{r}_{ij}^2)^2} \text{Var}(u_i) = \frac{1}{\sum_{i=1}^n \hat{r}_{ij}^2} \sigma^2 = \frac{\sigma^2}{SSR_j}$$

Since $R_j^2 = 1 - \frac{SSR_j}{SST_j}$, then $SSR_j = (1 - R_j^2) SST_j$

$$\text{then } \frac{\sigma^2}{SSR_j} = \frac{\sigma^2}{SST_j(1 - R_j^2)} > \frac{\sigma^2}{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 (1 - R_j^2)}$$

$$SST_j = \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 \times n \sigma^2$$

是因为 sample 是有限的

(i) Multicollinearity $\Rightarrow R_j^2 \rightarrow 1$ for some j

T^2 大, 只要 $R^2 \neq 1$, 有时就不符合 Ass 3

(ii) $R_j^2 \rightarrow 1 \Rightarrow \text{Var}(\hat{\beta}_j) \rightarrow \infty$ 即: 多重共线性越强, Ass 3 越被违背, MCR 越好模型越差, 参数的变动越大

(iii) $\text{se}(\hat{\beta}_j) = \sqrt{\frac{\sigma^2}{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 (1 - R_j^2)}}$ 注意: 此时是 se , 故用 $\hat{\sigma}^2$ 而不是 $\sigma^2 = \frac{\sum_{i=1}^n \hat{u}_i^2}{n-k+1}$

④ $\text{Var}(\hat{\beta}_1)$ (Special cases)

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$$

$$\tilde{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 \quad 2$$

$$\tilde{y} = \tilde{\beta}_0 + \tilde{\beta}_1 x_1 \quad 1$$

$$\text{Var}(\tilde{\beta}_1) = \frac{\sigma^2}{\sum(x_{1i} - \bar{x}_1)^2}$$

(无论是 overspecifying / underspecifying, 都有: $\text{Var}(\tilde{\beta}_1) < \text{Var}(\hat{\beta}_1)$)

原因:

$$\begin{cases} \text{若 } \beta_2 = 0 \text{ (忽略MLR overspecifying)} \rightarrow \\ E(\hat{\beta}_1) = \beta_1 \\ E(\tilde{\beta}_1) = \beta_1 \\ \text{Var}(\tilde{\beta}_1) < \text{Var}(\hat{\beta}_1) \text{ (当不一致时)} \end{cases}$$

$$\begin{cases} \text{若 } \beta_2 \neq 0 \text{ (忽略MLR underspecifying)} \rightarrow \\ E(\hat{\beta}_1) = \beta_1 \\ E(\tilde{\beta}_1) \neq \beta_1 \\ \text{Var}(\tilde{\beta}_1) < \text{Var}(\hat{\beta}_1) \text{ (当不一致时)} \end{cases}$$

Trade off between bias & variance

Bias 上升
variance 下降

⇒ 由于 (Bias 不可接受)
Variance 可通过提高来减小, 所以当 underspecifying 时, 我们选择 MLR

③ Goodness of fit

R² 可以表示的前 k 个 相同 sample

$$R^2 = \frac{SSE}{SST} = 1 - \frac{SSR}{SST} \quad (\text{definition})$$

$$= \text{Corr}(y_i, \hat{y}_i) \stackrel{?}{=} \left(\frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{y})}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2} \sqrt{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}} \right)^2 \quad (\text{Intuitive explanation})$$

当我们引进更多 independent variables, $R^2 \rightarrow 1 \rightarrow R^2$ 说明模型目前已引入 x_i 来解释 y_i
如果 x_i 是无关变量, 则 $R^2 \rightarrow 0$

④ Variance inflation factor (VIF 表示 Multicollinearity 的一个量度)

$$\text{VIF}_j = \frac{1}{1-R_j^2} \quad (\text{当共线性越强, VIF 越大})$$

$$(\text{Var}(\hat{\beta}_j)) = \frac{\sigma^2}{SST_j(1-R_j^2)} = \frac{\sigma^2}{SST_j} \text{VIF}_j$$

\hookrightarrow 值不要大于 10 (> 10 Strong multicollinearity)
 相比于 VIF_j , 我们更在意 $\text{Var}(\hat{\beta}_j)$ (< 10 Little multicollinearity)

4. Regression through the origin

① Regression model (ASSUMPTION)

$$Y = \beta_1 X_1 + \dots + \beta_k X_k + U$$

② Sample regression function

At this time, similar to SLR, $\tilde{\beta}_1, \dots, \tilde{\beta}_k$ are biased

(丢掉 β_0) \rightarrow β 会溢或 Bias
 丢掉 relevant X_i

5. Efficiency of OLS: The Gauss-Markov Theorem

The Gauss-Markov Theorem (满足5个假设的OLS) = BLUE

B = Best (with smallest variance) 对于每一个 $\hat{\beta}_j$ 都成立

L = Linear (每个 estimator 都是 y_i 的 linear combination)

U = Unbiased

E = Estimator (Data \rightarrow Estimator)

① Under OLS, parameters are unbiased linear estimator

$$SCE: \hat{\beta}_1 = \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2}$$

都是 linear combination of y_i

$$MLE: \hat{\beta}_j = \frac{\sum r_{ij} \hat{y}_i}{\sum r_{ij}^2}$$

② Under OLS, the variances of unbiased estimator are smallest

\Rightarrow Smallest variance \neq Efficiency !!!

Smallest variance

Efficiency: To Unbiasedness By 方差最小, smallest variance

L11-14 Multiple Linear Regression Model — statistical inference

1. Classical linear model (CLM) assumption

Classical linear model = Gauss-Markov assumption + ass b:

$$u_i \sim N(0, \sigma^2)$$

(其实就是 combine ass 4 & ass 5, 并进一步推导 → Distribution)

$$\rightarrow \begin{cases} \hat{\beta}_j \sim N(\beta_j, \frac{\sigma^2}{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 (1 - R_j^2)}) \\ y|x \sim N(\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k, \sigma^2) \end{cases}$$

$$E(y|x) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

2. T-statistic

① Definition (T-statistic)

$$\frac{\hat{\beta}_j - \beta_j}{\text{se}(\hat{\beta}_j)} \sim t(n-k+1)$$

$\hat{\beta}_j$ 的自由度是 $n - k + 1$ (因为 β_0, \dots, β_k 已经被估计了)

$$(这是由 \frac{\hat{\beta}_j - \beta_j}{\text{se}(\hat{\beta}_j)} \sim N(0, 1))$$

② Properties

可以直观理解

△ t distribution is more spread out than standard normal distribution

△ $n \uparrow, df \uparrow$, the difference between t-distribution & standard normal distribution ↓ (especially for $df > 120$)
(t distribution vs df)

3. T-HT

① Step 1 (写背景信息)

$$\begin{cases} H_0 \\ H_1 \quad \begin{cases} \text{One-sided} \\ \text{Two-sided} \end{cases} \end{cases}$$

Significance level (Type I error) $\alpha = P(\text{Type I error}) = P(T \in C | H_0)$

$\beta = P(\text{Type II error}) = P(T \notin C | H_1)$

② Step 2 (计算统计量)

$$t_{\hat{\beta}_j} = \frac{\hat{\beta}_j - \beta_j}{\text{se}(\hat{\beta}_j)} \quad (\text{Based on } H_0)$$

$$t = \frac{(\hat{\beta}_1 - \hat{\beta}_2) - (\beta_1 - \beta_2)}{\text{se}(\hat{\beta}_1 - \hat{\beta}_2)}$$

In STATA: vce

★注意：STATA 中的 t 值都是基于 $H_0: \beta_j = 0$

T 大 → H_0 is rejected → $\beta_j \neq 0$ → x_j is statistically significant

③ Step 3 (判断)

- a. Critical region
 $|t_{\hat{\beta}_j}| \leq |C|$: Fail to reject
 $|t_{\hat{\beta}_j}| > |C|$: Reject
- ↓ ↓
 计算 查表

☆注意: C 的下标及 Significance level

One-sided

Two-sided 无序减半

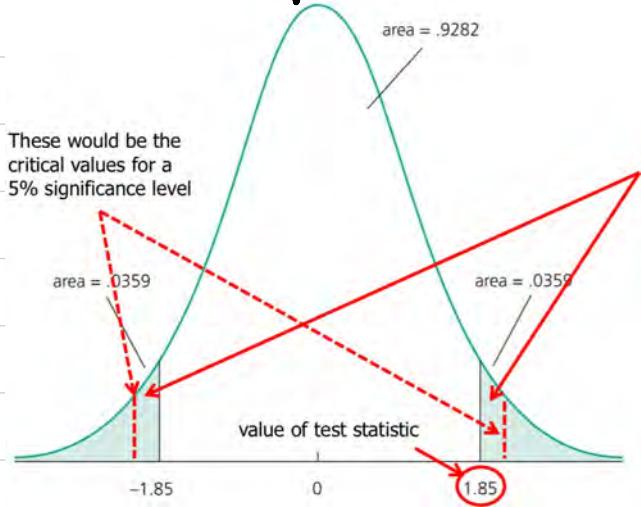
b. p-value

$$p\text{-value} = P(|T| > t)$$

- $p\text{-value} > \alpha$: Fail to reject
 $p\text{-value} < \alpha$: Reject

☆注意: STATA 报出的 p-value 是 two-sided

(因为 STATA 默认 Two-sided HT)



c. Confidence intervals

$$CI: [\hat{\beta}_j - z_{\alpha/2} \text{se}(\hat{\beta}_j), \hat{\beta}_j + z_{\alpha/2} \text{se}(\hat{\beta}_j)]$$

4. F-statistic

① Motivation

(T-statistic: Simple HT)

(F-statistic: Multiple HT / Joint HT) Unrestricted model → Restricted model

当发现 Strong multicollinearity 时, 我们把研究对象从一个简化为一组。

② Definition

$$F > \frac{(SSR_r - SSR_{ur}) / (df_r - df_{ur})}{SSR_{ur} / df_{ur}} \sim F(df_r - df_{ur}, df_{ur})$$

(在STATA中, 用 test 命令)

5. F-HT

① Step 1

H_0 (Multiple exclusion) $\rightarrow H_0: \beta_j = 0$

Overall significant (Overall exclusion)

General linear restriction $\left(\begin{array}{l} \beta_1=1, \beta_2=2, \beta_3=3 \\ \beta_1=\beta_2, \beta_3=1 \end{array} \right)$ (if H0 is true, then linear combination)

H_1 : H_0 is not true (F-HT 只能检测双侧)

Significance level

不能用 R-squared 来算 F-stat. $\because T_r \neq SSt_r$
(Restricted model 差值不为零)

② Step 2

计算 F-statistic

③ Step 3

a. Critical region

$F \leq C$: Fail to reject

$F > C$: Reject

b. P-value

$P\text{-value} = P(F > f)$

$P\text{-value} \geq \alpha$: Fail to reject

$p\text{-value} < \alpha$: Reject

VARIABLES	Dependent Variable: log(salary)			
	(1)	(2)	(3)	(4)
years	0.069*** (5.684)	0.068*** (5.592)	0.070*** (5.776)	0.071*** (5.703)
gamesyr	0.013*** (4.742)	0.016*** (10.079)	0.011*** (5.202)	0.020*** (15.023)
bavg	0.001 (0.887)	0.001 (1.331)	0.001 (0.691)	---
hrunsyr	0.014 (0.899)	0.036*** (4.964)	---	0.014 ---
rbisyr	0.011 (1.500)	---	0.017*** (5.117)	---
Constant	11.192*** (38.752)	11.021*** (41.476)	11.275*** (41.197)	11.224*** (103.625)
Observations	353	353	353	353
R-squared	0.628	0.625	0.627	0.597

▲ 10% → 5% → 1%

6. Relationship between F & T statistics

① F-T

$$F_{1, df_{ur}} = t^2 |df_{ur}|$$

② F-R²

$$F = \frac{(R_{ur}^2 - R_r^2) / (df_r - df_{ur})}{(1 - R_{ur}^2) / df_{ur}}$$

proof:

$$F = \frac{SSR_p - SSR_{ur}}{SSR_{ur} / df_{ur}}$$

Since $R^2 = 1 - \frac{SSR}{SST}$, then $SSR = (1 - R^2) SST$

$$\text{then } F = \frac{[(1 - R_r^2) SST - (1 - R_{ur}^2) SST] / (df_r - df_{ur})}{(1 - R_{ur}^2) SST / df_{ur}} = \frac{(R_{ur}^2 - R_r^2) / (df_r - df_{ur})}{(1 - R_{ur}^2) / df_{ur}}$$

7. Change of units (不變方程變lnx)

Units change \rightarrow Estimator's change \rightarrow R^2 stays unchanged
T stat & F stat stay unchanged

L15-1b Multiple Linear Regression Model — Further issues

1. Large sample vs Small sample

① Motivation

In CLM assumptions,

in practice, some of them cannot be satisfied. Especially, the b^{th} one: Normality assumption.

② Consistency

a. Definition

$$\lim \Pr(\theta_n - \theta) < \varepsilon = 1, \quad \forall \varepsilon > 0 \text{ and } n \rightarrow \infty \triangleq \text{plim } \theta_n = \theta$$

Ass 1 – Ass 9 \longrightarrow Unbiasedness & Consistency

proof

$$\text{In MLR, } \hat{\beta}_1 = \beta_1 + \frac{\sum (x_{1i} - \bar{x}_1) u_i}{\sum (x_{1i} - \bar{x}_1)^2} = \beta_1 + \frac{\frac{1}{n} \sum (x_{1i} - \bar{x}_1) u_i}{\frac{1}{n} \sum (x_{1i} - \bar{x}_1)^2}$$

$$\text{plim } \hat{\beta}_1 = \beta_1 + \frac{\text{Cov}(x_1, u)}{\text{Var}(x_1)}$$

For OLS, (Ass 9 ✓)

By Ass 9 (Zero conditional mean), we get $\text{Cov}(x_1, u) = 0$, Weaker than Ass 9

So, $\text{plim } \hat{\beta}_1 = \beta_1$ (Consistent)

For underspecified OLS, (Ass 9 ✗)

For example, $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + v \rightarrow \tilde{y} = \beta_0 + \beta_1 x_1 + u$

$$\begin{aligned} \text{Cov}(x_1, u) &= \text{Cov}(x_1, \beta_2 x_2 + v) = \beta_2 \text{Cov}(x_1, x_2) + \text{Cov}(x_1, v) \\ &= \beta_2 \text{Cov}(x_1, x_2) \end{aligned}$$

$$\rightarrow \text{plim } \hat{\beta}_1 = \beta_1 + \beta_2 \frac{\text{Cov}(x_1, x_2)}{\text{Var}(x_1)} = \beta_1 + \beta_2 \beta_1 \quad (\text{Inconsistent})$$

b. Interpretation

$n \uparrow, \theta_n \rightarrow \theta$

Make the estimator get closer to the true value

Although not all useful estimators are unbiased, virtually all economists agree that **consistency** is a minimal requirement for an estimator. The Nobel Prize-winning econometrician Clive W. J. Granger once remarked, "If you can't get it right as n goes to infinity, you shouldn't be in this business." The implication is that, if your estimator of a particular population parameter is not consistent, then you are wasting your time.

c. Example

$$\text{plim } (\bar{x}_n) = \mu$$

③ Asymptotic normality

Normality \longrightarrow Asymptotic normality

CLT
Large sample

2. Standardized effect

① Derivation

a. Simple linear regression

$$\begin{aligned} Y_i &= \hat{\beta}_0 + \hat{\beta}_1 X_i + \hat{u}_i \\ \bar{Y} &= \hat{\beta}_0 + \hat{\beta}_1 \bar{X} \end{aligned} \longrightarrow \frac{Y_i - \bar{Y}}{\sigma_Y} = \hat{\beta}_1 \frac{\hat{\sigma}_X}{\hat{\sigma}_Y} \frac{X_i - \bar{X}}{\hat{\sigma}_X} + \frac{\hat{u}_i}{\hat{\sigma}_Y}$$

標準化變量 → 標準化斜率

$$\longrightarrow z_y = \hat{b}_1 z_x + \frac{\hat{u}_i}{\hat{\sigma}_Y}, \quad \hat{b}_1 = \hat{\beta}_1 \frac{\hat{\sigma}_X}{\hat{\sigma}_Y}$$

b. Multiple linear regression

$$z_y = \hat{b}_1 z_1 + \dots + \hat{b}_K z_K + \text{error}, \quad \hat{b}_j = \hat{\beta}_j \frac{\hat{\sigma}_j}{\hat{\sigma}_Y}$$

斜率和截距的標準化系数

Beta coefficient / Standardized coefficient

② Significance

The significance here is to standardize the effect, then to compare the economic significance / effect
从而解决单位不同的无法比较的问题

3. More on Functional form

① Natural log

a. Example

$$\Delta \ln x = \frac{\% \Delta x}{100}$$

When estimated using the data in HPRICE2, we obtain

$$\widehat{\log(price)} = 9.23 - .718 \log(nox) + .306 rooms$$
$$(0.19) \quad (.066) \quad (.019)$$
$$n = 506, R^2 = .514.$$

[6.7]

Thus, when nox increases by 1%, $price$ falls by .718%, holding only $rooms$ fixed. When $rooms$ increases by one, $price$ increases by approximately $100(.306) = 30.6\%$.

When β_i is not too large, then $\log(y_i)$ is a good estimate.

b. Vantages & Drawbacks

△ Vantages: Mitigate heteroskedasticity → Better fit the assumption

Narrow down the range of dependent variable → Alleviate outlier influence

△ Drawbacks: The variables taking natural log must be positive, and the variable had better be away from 0

We can get $\log(y)$, but not the original one y

c. Applications

△ Often used for dollar amount (wage, salary...)

△ Not used for (year)

percentage-point variables

② Quadratic form

$$\widehat{wage} = 3.73 + .298 exper - .0061 exper^2$$
$$(3.5) \quad (.041) \quad (.0009)$$

当同一个 effect 同时以一次型与二次型出现在 regression model 中时，

通过系数判断其关系影响，通过 F 值进行分层。

③ Interaction term

$$\text{Original: } y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + u$$

$$\text{Transformed: } y = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 (x_1 - \bar{x}_1)(x_2 - \bar{x}_2) + u \rightarrow \text{可以用于计算 } x_1 | x_2 = \bar{x}_2 \text{ sample mean } \beta_3 \text{ partial effect}$$

4. More on R^2

① Adjusted / Corrected R^2

$$R^2 = 1 - \frac{SSR}{SST} = 1 - \frac{SSR/n}{SST/n} = 1 - \frac{\hat{\sigma}_u^2}{\hat{\sigma}_y^2}$$

$$\bar{R}^2 = 1 - \frac{\hat{\sigma}_u^2}{\hat{\sigma}_y^2} = 1 - \frac{SSR/(n-k-1)}{SST/(n-1)}$$

② Relationship between R^2 & \bar{R}^2

$$\bar{R}^2 = 1 - (1 - R^2) \frac{n-1}{n-k-1}$$

(\bar{R}^2 可能是负的)

③ Pros & Cons

Pro: 对于3个独立变量，Goodness-of-fit也下降 (用于比较不同model, 同-sample)

Con: The ratio of two unbiased estimators is not an unbiased estimator (统计学不认为 \bar{R}^2 大于 R^2)

注意: R^2 和 \bar{R}^2 都不能作为 Goodness-of-fit 的指标来确定 dependent variable 的形式

决定 ind variable 形式 ✓

决定 dep variable 形式 ✗

L17-18 Dummy variables

1. Motivation

When we want to incorporate some qualitative-information variables (e.g. Gender, region...) into our regression model, then at this time, they are called dummy variables.

(Binary variable | Zero-one variable)

2. Single dummy variable

① Process

a. Introduce dummy variables

$$wage = \beta_0 + \delta_0 female + \beta_1 educ + u$$

= the wage gain/loss if the person is a woman rather than a man (holding other things fixed)



Dummy variable:
 $female = 1$ if the person is a woman
 $female = 0$ if the person is man

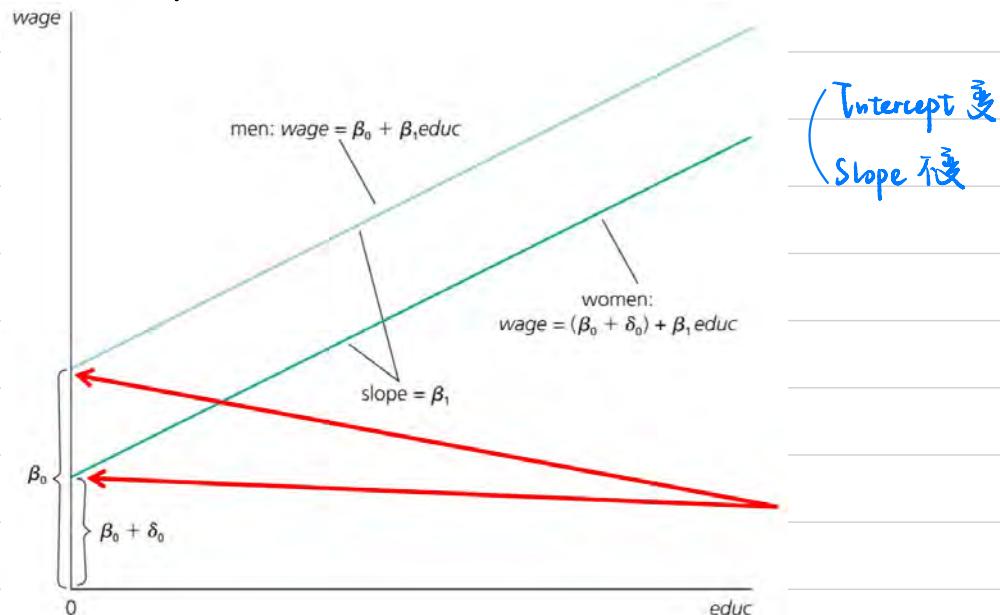
The base category is men.
(benchmark category)

b. Compute the coefficient of dummy variables

$$\delta_0 = E(wage|female = 1, educ)$$

$$-E(wage|female = 0, educ)$$

② Graphic interpretation



③ Dummy variable trap

Dummy variables to all categories 不可以在 model 中同时出现, (i.e. 不要引进 base case)

完美共线性 Perfect collinearity

3. Multiple dummy variables

① Process

a. Introduce & transform the dummy variables

- Two qualitative variables: *married* and *female*

$$\text{To } \hat{\log(wage)} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{married} + \dots ?$$

An important limitation of this model is that the marriage premium is assumed to be the same for men and women; this is relaxed in the following example.

marrmale (married males): *married* = 1, *female* = 0

marrfem (married females): *married* = 1, *female* = 1

singmale (single males): *married* = 0, *female* = 0

singfem (single females): *married* = 0, *female* = 1

$$\hat{\log(wage)} = .321 + .213 \text{marrmale} - .198 \text{marrfem}$$

$$- .110 \text{singfem} + .079 \text{educ} + .027 \text{exper} - .00054 \text{exper}^2$$

b. Compute the coefficients

→ w. base category ($\sigma = 0$) 为基准, 相互比较 effect

② Deal with ordinal variables with multiple dummy variables

Ordinal variables (0, 1, 2, 3, ...) → Multiple dummy variables (0, 1)

a. Example 1

Example: City credit ratings and municipal bond interest rates

Municipal bond rate	Credit rating from 0-4 (0=worst, 4=best)
↓	↓
$MBR = \beta_0 + \beta_1 CR + \text{other factors}$	

This specification would probably not be appropriate as the credit rating only contains ordinal information. A better way to incorporate this information is to define dummies:

$$MBR = \beta_0 + \delta_1 CR_1 + \delta_2 CR_2 + \delta_3 CR_3 + \delta_4 CR_4 + \text{other factors}$$

这样的表现, 相比于之前模型,
保留更多的局部性特征

b. Example 2

Example: effects of law school rankings on starting salaries

(LAWSCH85.dta)

	$\hat{\log(salary)} = 9.17 + .700 \text{top10} + .594 \text{r11to25} + .375 \text{r26to40}$
↓	$+ .263 \text{r41to60} + .132 \text{r61to100} + .0057 \text{LSAT}$
↓	$+ .014 \text{GPA} + .036 \log(\text{libvol}) + .0008 \log(\text{cost})$

Holding other things fixed,
students from top 10 schools
earn 70% more than from
schools below 100

$$n = 136, R^2 = .911, \bar{R}^2 = .905$$

4. Interactions involving dummy variables

① Motivation

Interaction of quantitative variable & quantitative variable \longrightarrow quantitative variable

Interaction of dummy variable & quantitative variable \longrightarrow dummy variable

② Example

□ Restricted model (same regression for both groups)

College grade point average	Standardized aptitude test score	High school rank percentile
-----------------------------	----------------------------------	-----------------------------

$$cumgpa = \beta_0 + \beta_1 sat + \beta_2 hspc + \beta_3 tohrs + u$$

Total hours spent in college courses

□ Unrestricted model (contains full set of interactions)

$$cumgpa = \beta_0 + \delta_0 female + \beta_1 sat + \delta_1 female sat + \beta_2 hspc + \delta_2 female hspc + \beta_3 tohrs + \delta_3 female tohrs + u$$

③ Hypothesis testing (Whether there is group difference)

a. Method 1: General F-test: R \longrightarrow UR

□ Null hypothesis

$$H_0 : \delta_0 = 0, \delta_1 = 0, \delta_2 = 0, \delta_3 = 0$$

All interaction effects are zero, i.e.,
the same regression coefficients
apply to men and women

b. Method 2: Chow-Test

$$SSR_{UR} = SSR_1 + SSR_2$$

$$SSR_p = \text{pooled } SSR$$

$$\rightarrow F = \frac{(SSR_p - SSR_{UR}) / (df_{UR} - df_{UR})}{SSR_{UR} / df_{UR}} = \frac{(SSR_p - (SSR_1 + SSR_2)) / (n - (k+1) - n + 2(k+1))}{(SSR_1 + SSR_2) / (n - 2(k+1))} = \frac{(SSR_p - (SSR_1 + SSR_2)) / (k+1)}{(SSR_1 + SSR_2) / (n - 2(k+1))}$$

这里的 SSR_p , SSR_1 , SSR_2 来自 Restricted model

5. Binary dependent variable

① Model

Linear regression model \rightarrow Linear probability model (LPM)

不是分类

② Derivation

Since PRF:

$$E(y|x) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

Conditional expectation: $E(y|x) = 1 \cdot P(y=1|x) + 0 \cdot P(y=0|x) = P(y=1|x)$

Then we get $P(y=1|x) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$

$$\beta_j = \frac{\partial P(y=1|x)}{\partial x_j}$$

Linear probability model (LPM) 不是分类 $P(y=1|x)$

③ Sample regression function

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_k x_k$$

④ Downside of LPM

a. \hat{y} can be out of range: $\hat{y} > 1$ & $\hat{y} < 0$ can happen and cannot be explained

b. Constant partial effect

c. Heteroskedasticity: $\text{Var}(y|x) = P(y=1|x)[1 - P(y=1|x)]$ (取大于x的那部分)

proof:

$$\begin{aligned} \text{Var}(y|x) &= E(y^2|x) - E(y|x)^2 = E(y|x) - E(y|x)^2 = E(y|x)[1 - E(y|x)] \\ &= (\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k)(1 - \beta_0 - \beta_1 x_1 - \dots - \beta_k x_k) \\ &= P(y=1|x)[1 - P(y=1|x)] \end{aligned}$$

⑤ Goodness of fit

Rather than R^2 , we use "correct percent" to measure goodness-of-fit

First, we define $\hat{y}_i = \begin{cases} 1, & \text{if } \hat{y}_i \geq 0.5 \\ 0, & \text{if } \hat{y}_i < 0.5 \end{cases}$ (y_i 则不用 predict 因为 $y_i = 0 \text{ or } 1$)

Second, we observe the results and compute correct percent

		$inlf_i$		Total
1975		0	1	
	$inlf_i$	0	1	
0	203	122		325
1	78	350		428
Total	281	472		753

goodness of fit = correct percent = $\frac{203 + 350}{753}$

L18 Instrumental variables → Ass 9

1. Motivation

被遗漏的变量: 因为不可观测 (unobservable)

Endogeneity problem (omit some independent variables) → Biass

遗漏的独立变量 \hat{u} 与 x_i 是相关的, \rightarrow 导致回归模型有偏.

Solutions to deal with Endogeneity problem:

Proxy variables method : 用 proxy 变量代替不可观测的独立变量

Instrumental variables method (IV) : 本原理论

Fixed effect method : 固定效应

2. Instrumental variables (IV)

① Definition 定义: 将 endogenous variable 替换为 exogenous variable

x 's instrumental variable is z , which satisfies the following three conditions:

a. z does not appear in the regression

b. $\text{Cov}(z, x) \neq 0$ z 与 endogenous variable x 有关

c. $\text{Cov}(z, u) = 0$

② IV-estimator

If we assume the existence of IV z , then we have $\text{Cov}(z, u) = 0$

$$\rightarrow \text{Cov}(z, y - \beta_0 - \beta_1 x) = \text{Cov}(z, y) - \beta_1 \text{Cov}(z, x) = 0$$

$$\rightarrow \beta_1 = \frac{\text{Cov}(z, y)}{\text{Cov}(z, x)}$$

$$\rightarrow \hat{\beta}_{IV} = \frac{\widehat{\text{Cov}}(z, y)}{\widehat{\text{Cov}}(z, x)} = \frac{\frac{n}{n} \sum (z_i - \bar{z})(y_i - \bar{y})}{\frac{n}{n} \sum (z_i - \bar{z})(x_i - \bar{x})}$$

③ Quality of IV-estimator

$$\text{plim } \hat{\beta}_1, \text{OLS} = \beta_1 + \text{Corr}(x, u) \frac{\sigma_u}{\sigma_x}$$

$$\text{plim } \hat{\beta}_1, \text{IV} = \beta_1 + \frac{\text{Corr}(z, u)}{\text{Corr}(z, x)} \frac{\sigma_u}{\sigma_x}$$

Poor/weak instrumental variables: $\text{Corr}(z, x)$ is small

④ Two stage least square estimation (2SLS estimation)

We use 2SLS to deal with endogenous variable:

First stage: reg this endogenous variables on other exogenous variables (Reduced form regression)

Second stage: reg y on all the exogenous variables

$$\text{eq. } y_1 = \beta_0 + \beta_1 y_2 + \beta_2 z_1 + \dots + \beta_k z_{k-1} + u_1 \quad (\text{y}_2 \text{ endo}, z_1, \dots, z_{k-1} \text{ exo})$$

$$1^{\text{st}} \text{ stage: } \hat{y}_2 = \hat{\pi}_0 + \hat{\pi}_1 z_1 + \dots + \hat{\pi}_{k-1} z_{k-1} + \hat{\pi}_k z_k \quad \boxed{\text{Additional exogenous variable}}$$

$$2^{\text{nd}} \text{ stage: } y_1 = \beta_0 + \beta_1 \hat{y}_2 + \beta_2 z_1 + \dots + \beta_k z_{k-1} + \text{error}$$

⑤ Test after 2SLS estimation

$$H_0: \delta_1 = 0$$

$$H_1: \delta_1 \neq 0$$

$$y_1 = \beta_0 + \beta_1 y_2 + \beta_2 z_1 + \dots + \beta_k z_{k-1} + u_1$$

Reduced form regression: Variable that is suspected to be endogenous

$$y_2 = \pi_0 + \pi_1 z_1 + \dots + \pi_{k-1} z_{k-1} + \pi_k z_k + v_2$$

Variable y_2 is exogenous if and only if v_2 is uncorrelated with u_1 , i.e. if the parameter δ_1 is zero in the regression:

$$u_1 = \delta_1 v_2 + e_1$$

The null hypothesis of exogeneity of y_2 is rejected, if in this regression the parameter δ_1 is significantly different from zero

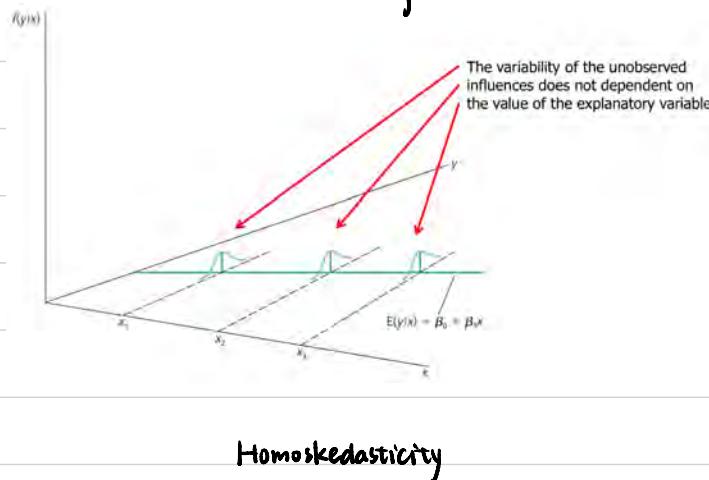
Test equation:

$$y_1 = \beta_0 + \beta_1 y_2 + \beta_2 z_1 + \dots + \beta_k z_{k-1} + \delta_1 \hat{v}_2 + e_1$$

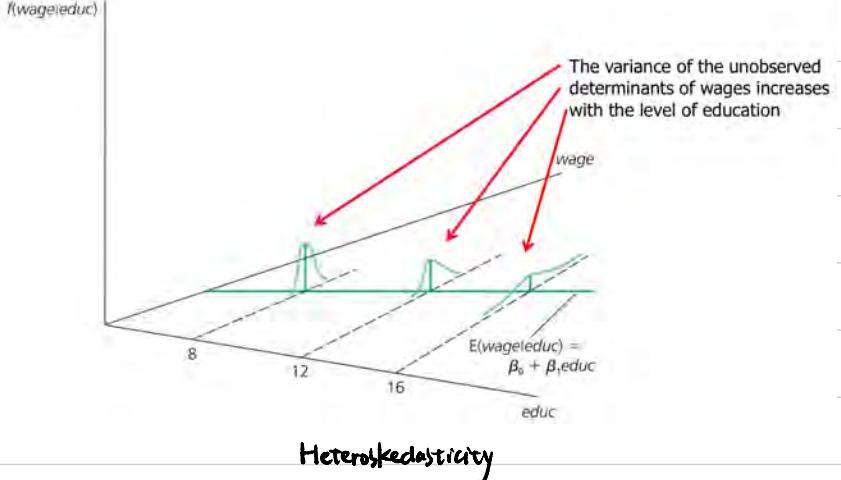
$$y_2 \text{ exogenous} \rightarrow \text{Cov}(y_2, u_1) = 0 \rightarrow \text{Cov}(v_2, u_1) = 0 \rightarrow \delta_1 = 0$$

L19 Heteroskedasticity → Ass 5

1. Hom & Hetero-skedasticity



$\rightarrow \text{Var}(u_i | \mathbf{x}) / \text{Var}(y) \text{ 不变} : \text{Var}(u_i | \mathbf{x}_i) = \sigma^2$



$\rightarrow \text{Var}(u_i | \mathbf{x}) / \text{Var}(y) \text{ 变} : \text{Var}(u_i | \mathbf{x}_i) = \sigma^2 h(x_i)$

2. Consequences of Heteroskedasticity

- ① OLS is still consistent
- ② F-test & T-test are not valid
- ③ Interpretation of R^2 is not changed
- ④ OLS is no longer BLUE, there may be more efficient estimators

3. Testing for heteroskedasticity

① BP test (u^2 与 x_i 线性不相关)

$$\begin{aligned} H_0: \text{Homoskedasticity} &: \text{Var}(u|\mathbf{x}) = \sigma^2 \\ H_1: \text{Heteroskedasticity} & \longrightarrow \begin{cases} H_0: \delta_1 = \dots = \delta_k = 0 \\ H_1: H_0 \text{ is not true} \end{cases} \end{aligned}$$

$$\text{Under } H_0: \text{Var}(u|\mathbf{x}) = E(u^2|\mathbf{x}) - [E(u|\mathbf{x})]^2 = E(u^2|\mathbf{x}) = \sigma^2$$

$$\longrightarrow E(u^2|\mathbf{x}_1, \dots, \mathbf{x}_k) = E(u^2) = \sigma^2$$

$$\longrightarrow \hat{u}^2 = \delta_0 + \delta_1 x_1 + \dots + \delta_k x_k + \text{error}, \text{ where } \delta_1 = \dots = \delta_k = 0 \quad (*) \quad \text{此时用 } \hat{u}^2 \text{ 替代 } u^2$$

→ F test

② White test (u^2 与 x_i 线性、二次型不相关)

a. Original White test

Start from (*), we have $\hat{u}^2 = \delta_0 + \delta_1 x_1 + \dots + \delta_k x_k + \delta_{k+1} x_1^2 + \dots + \delta_{2k} x_k^2 + \delta_{2k+1} x_1 x_2 + \dots + \delta_n x_{k-1} x_k + \text{error}$, where $\delta_1 = \dots = \delta_n = 0$

BP: 从一阶导至二阶，这个推导更正确

b. Alternative White test

$$\hat{u}^2 = \delta_0 + \delta_1 \bar{y} + \delta_2 \bar{y}^2 + \text{error}, \text{ where } \delta_1 = \delta_2 = 0$$

4. Heteroskedasticity-robust estimation → $\hat{\sigma}_{\text{Var}}$ (没有消除) 没有消除 Heteroskedasticity, R 使用了一种方法
 ① Implementation

In STATA:

reg ... , robust 加一个 "robust"

② Results

Large sample → Closer to homogeneity (Estimator 不变)

Variance of estimator 变大变小

Robust inference → OLS isn't BLUE
 WLS → OLS is BLUE

5. Weighted least squares estimation (WLS) estimation

① Procedure of WLS estimation

If heteroskedasticity, then we def $\text{Var}(u_i|x_i) = \sigma^2 h(x_i)$, $h(x_i) = h_i > 0$

$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + u_i$ → Var会减小

$$\rightarrow \left(\frac{y_i}{\sqrt{h_i}} \right) = \beta_0 \frac{1}{\sqrt{h_i}} + \beta_1 \frac{x_{i1}}{\sqrt{h_i}} + \dots + \beta_k \frac{x_{ik}}{\sqrt{h_i}} + \frac{u_i}{\sqrt{h_i}}$$

$$\rightarrow y_i^* = \beta_0 x_{i0}^* + \beta_1 x_{i1}^* + \dots + \beta_k x_{ik}^* + u_i^*$$

(这里还有一个)

② Examples

a. Example 1: $E(u_i|x_i)$

Example: Savings and income

$$sav_i = \beta_0 + \beta_1 inc_i + u_i, \text{Var}(u_i|inc_i) = \sigma^2 inc_i$$

$$\left[\frac{sav_i}{\sqrt{inc_i}} \right] = \beta_0 \left[\frac{1}{\sqrt{inc_i}} \right] + \beta_1 \left[\frac{inc_i}{\sqrt{inc_i}} \right] + u_i^* \quad \text{Note that model has}$$

The transformed model is homoskedastic

$$E(u_i^{*2}|x_i) = E \left[\left(\frac{u_i}{\sqrt{h_i}} \right)^2 |x_i \right] = \frac{E(u_i^2|x)}{h_i} = \frac{\sigma^2 h_i}{h_i} = \sigma^2$$

b. Example 2: $E(u_i|x_i)$: Homo → Hetero

Example: factors affect pension plan contribution (individual-level)

Firm-level equation

Contribution to pension plan for employee e in firm i	Earnings and age for employee e in firm i	Percentage firm contributes to plan	Heteroskedasticity error term
---	---	-------------------------------------	-------------------------------

$$contribution_{i,e} = \beta_0 + \beta_1 earnings_{i,e} + \beta_2 age_{i,e} + \beta_3 mrate_{i,e} + u_{i,e}$$

Average contribution to pension plan in firm i

Average earnings and age in firm i

Percentage firm contributes to plan

Average error across all employees in firm i

$$contribution_i = \beta_0 + \beta_1 earnings_i + \beta_2 age_i + \beta_3 mrate_i + \bar{u}_i$$

Let m_i denote the number of employees at each firm

$$\Rightarrow \text{Var}(\bar{u}_i) = \text{Var} \left(\frac{1}{m_i} \sum_{e=1}^{m_i} u_{i,e} \right) = \sigma^2 / m_i \quad \text{Heteroskedasticity error at the firm-level}$$

- Suppose, however, that we only have firm-level data, i.e., average values of contributions, earnings, and age by firm

If errors are homoskedastic at the employee level, WLS with weights equal to firm size m_i should be used.

→ 对变量取对数之后，要使用 WLS 估计为 Homogeneity 模型

③ WLS with 2PM

$$\text{In LPM, } \text{Var}(y|x) = E(y^2|x) - [E(y|x)]^2 = E(y|x)[1 - E(y|x)] = P(y=1|x)(1 - P(y=1|x)) = y_i(1-y_i)$$

$$\rightarrow \sigma^2 h_i = y_i(1-y_i) \quad \xrightarrow{\text{Defn}} \sigma^2 R_{yy}$$

→ Taking the above \hat{h}_i , we do WLS

$$\begin{cases} y_i < 0 \rightarrow \hat{y}_i = 0.01 \\ y_i > 0 \rightarrow \hat{y}_i = 0.99 \end{cases} \xrightarrow{\text{To guarantee } \hat{h}_i > 0}$$

b. Feasible generalized least squares estimation (FGLS)

① Setup

If heteroskedasticity, then we def $\text{Var}(u|x) = \sigma^2 h(x) = \sigma^2 \exp(\delta_0 + \delta_1 x_1 + \dots + \delta_K x_K)$

General form of heteroskedasticity, exp-function is used to ensure positivity

② Procedure of FGLS estimation

$$\text{Var}(u|x) = \sigma^2 \exp(\delta_0 + \delta_1 x_1 + \dots + \delta_K x_K)$$

$$\rightarrow u^2 = \sigma^2 \exp(\delta_0 + \delta_1 x_1 + \dots + \delta_K x_K) \quad \text{(Error term)}$$

$$\rightarrow \log(u^2) = \underbrace{\delta_0 + \delta_1 x_1 + \dots + \delta_K x_K}_\text{e} + e$$

$$\rightarrow \log(\hat{u}^2) = \hat{\delta}_0 + \hat{\delta}_1 x_1 + \dots + \hat{\delta}_K x_K$$

$$\rightarrow h(x_i) = \exp[\log(\hat{u}^2)]$$

$$\rightarrow WLS$$

□ A Feasible GLS Procedure to Correct for Heteroskedasticity:

- 1) Estimate the model by OLS and obtain the residuals, \hat{u}_i .
- 2) Create $\log(\hat{u}_i^2)$ by first squaring the OLS residuals and then taking the natural log.
- 3) Run the regression of $\log(\hat{u}_i^2)$ on all explanatory variables and obtain the fitted values, \hat{g}_i
- 4) Exponentiate the fitted values: $\hat{h}_i = \exp(\hat{g}_i)$
- 5) Find WLS estimation using weights $1/\hat{h}_i$

L20-21 Time Series

1. Static model

① Definition

We focus on the variables at the same time t

② Examples

$$inf_t = \beta_0 + \beta_1 unem_t + u_t$$

$$mrdrtet = \beta_0 + \beta_1 convrte_t + \beta_2 unem_t + \beta_3 yngmle_t + u_t$$

2. Finite distributed lag model (FDL Model)

① Definition

We focus on the variables with a time lag

② Population model

$$Y_t = \beta_0 + \beta_1 x_{t,1} + \dots + \beta_k x_{t,k} + u_t$$

(Transitory (Short-run) effect: 某一处未尽的系数

Permanent (Long-run) effect: 全部系数和

② Assumptions

a. Ass 1: Linear model

b. Ass 2: No perfect collinearity
 (No independent variable is constant
 No perfect linear combination)

c. Ass 3: Zero conditional mean: $E(u_t | X) = 0$ For any period, u_t is uncorrelated with the independent variables in all periods

"strict exogeneity" where $X = \begin{bmatrix} x_{11} & \dots & x_{1k} \\ \vdots & & \vdots \\ x_{t1} & \dots & x_{tk} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nk} \end{bmatrix}$

Past
Now
Future

△ Contemporaneous exogeneity: $E(u_t | x_t) = 0$ for the same period

Strict exogeneity: $E(u_t | X) = 0$ for all the periods Stronger

△ Ass 3 implies that in $y_t = \beta_0 + \beta_1 x_t + u_t$: β_1 has no lagged effect on y_t

proof: Contemporaneous exogeneity \Rightarrow

If there is lagged effect, then $x_{t-1} \neq 0$, x_{t-1} is included in x_t , then u_t is correlated with x_t , then Ass 3 is violated

△ Example of violating Ass 3

□ Example: City's murder rate and number of police officers

$$mrdrte_t = \beta_0 + \beta_1 polpct_t + u_t$$

we assume it is uncorrected with $polpct_t$

- City adjusts the size of its police force based on past values of murder rate. That means $polpct_{t+1}$ might be correlated with $mrdrte_t$ and u_t , which violates TS.3

$u_t \rightarrow mrdrte_t \rightarrow polpct_t$
 violate

d. Ass 4: Constant variance: $\text{Var}(u_t | X) = \text{Var}(u_t) = \sigma^2$

e. Ass 5: No serial correlation: $\text{Corr}(u_t, u_s | X) = 0, \forall t \neq s$

$\text{Corr}(x_t, x_s) \neq 0$

$\text{Corr}(u_t, x_s) = 0$

$\text{Corr}(u_t, u_s) = 0$

f. Ass 6: Normality: $u_t \sim N(0, \sigma^2)$

\rightarrow Ass 1 - Ass 3: Unbiasedness

\rightarrow Ass 1 - Ass 5: BLUE (滿意的OLS才具有Gauss-Markov Assumption in T)

\rightarrow Ass 1 - Ass 6: The usual T & F tests are valid (滿意的Normality是T/F Testing的依據)

④ Dummy variables

We use dummy variables to do event study.

Examine whether a particular event influences some outcome.

Example: Factors effect fertility rates (FERTIL3.dta)

Children born per 1,000 women in year t 1913-1984	Tax exemption in year t	Dummy for World War II years (1941-45)	Dummy for availability of contraceptive pill (1963-present)
--	-------------------------	--	---

$$\widehat{gfr}_t = 98.68 + .083 pe_t - 24.24 \boxed{ww2_t} - 31.59 \boxed{pill_t}$$

$$n = 72, R^2 = .473, \bar{R}^2 = .450$$

⑤ Trend

a. Method 1: Add a trend variable

A trend variable should be included when
 Some independent variables display trending behaviour
 Dependent variable displays trending behaviour

Example: Housing investment and prices (HSEINV.dta)

Per capita housing investment	Housing price index
-------------------------------	---------------------

$$\widehat{\log(invpc)} = -.550 + 1.241 \log(price)$$

$$n = 42, R^2 = .208, \bar{R}^2 = .189$$

It looks as if investment and prices are positively related

Example: Housing investment and prices (cont.)

$$\widehat{\log(invpc)} = -.913 + .381 \log(price) + .0098 t$$

$$n = 42, R^2 = .341, \bar{R}^2 = .307$$

There is no significant relationship between price and investment anymore

b. Method 2: Detrending

$$\begin{aligned} y_t &= \alpha_0 + \alpha_1 t + u_t \rightarrow \text{let } \hat{y}_t = \hat{u}_t \\ x_t &= \beta_0 + \beta_1 t + v_t \rightarrow \text{let } \hat{x}_t = \hat{v}_t \end{aligned} \quad) \text{除去} \begin{cases} \text{trend} \\ \text{constant} \end{cases}$$

→ Reg \hat{y}_t on \hat{x}_t (No Intercept)

c. Relation between two methods

$$\begin{aligned} (\text{Adding a trend}) &\rightarrow R^2 \uparrow \\ (\text{Detrending}) &\rightarrow R^2 \end{aligned} \quad \begin{cases} R^2 \text{ of reg detrended } y \text{ on } x \& t \\ = R^2 \text{ of reg detrended } y \text{ on detrended } x \\ \neq R^2 \text{ of reg } y \text{ on } x \& t \end{cases}$$

⑥ Seasonality

a. Method 1: Adding a set of seasonal dummies

b. Method 2: Desasonalizing 自变量 by seasons by dummy variables

L22-2) Panel data

1. Pooled OLS regression (POLS regression)

① Definition

Pooling independent cross sections across time:

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m + \delta_2 d_2 + \dots + \delta_n d_n$$

Mean values of variables. Include dummy variables for every period (时间不 \rightarrow base group)
have changed over time

② How to test whether the β 's are constants over time?

Chow Test \Downarrow (Variables & Time variables's interaction term 焦点都为0 \rightarrow 时间不 \rightarrow variable系数 \rightarrow Slope coefficients constant)

In Chapter 7, we discussed how the Chow test—which is simply an F test—can be used to determine whether a multiple regression function differs across two groups. We can apply that test to two

a. Compute SSR_p (无交互项 $\sum \delta_j d_j = 0$)

b. Compute SSR_1, \dots, SSR_T (所有Time variables $\sum \delta_j d_j = 0$ 且只保留对应时间的交互项)

c. Compute F : $F = \frac{[SSR_p - (SSR_1 + \dots + SSR_T)] / [T-1, k]}{(SSR_1 + \dots + SSR_T) / [n-T, k+1]}$

$$df(SSR_p) = n - k - (T-1) - 1 = (T-1)k$$

$$df(\sum SSR_i) = n - T(k+1)$$

2. Policy analysis with pooled cross section

① Motivation

Two or more independent cross sections can be used to evaluate the impact of a certain event or policy change

② Difference-in-difference estimator (DiD)

a. Definition

If we found that the slope coefficients change after the policy/event's implementation,

we introduce an interaction term. and at this time, the slope of this interaction term is called DiD.

b. Example

After incinerator was built:

$$\widehat{rprice} = 101,307.5 - 30,688.27 \text{nearinc}$$

(3,093.0) (5,827.71)

Dummy equals to one if the house is near the incinerator

The difference in the two coefficients on *nearinc* is

$$\hat{\delta}_1 = -30,688.27 - (-18,824.37) = -11,863.9$$

One has to compare with the situation before the incinerator was built:

$$\widehat{rprice} = 82,517.23 - 18,824.37 \text{nearinc}$$

(2,653.79) (4,744.59)

$$\widehat{rprice} = \beta_0 + \delta_0 \text{after} + \beta_1 \text{nearinc} + \delta_1 \text{after} \cdot \text{nearinc} + u$$

Differential effect of being in the location and after the incinerator was built

3. Panel data analysis

① Panel data regression function

a. Two-period

$$y_{it} = \beta_0 + \delta_0 d_{2t} + \beta_1 x_{it} + (a_i + u_{it})$$

d_{2t} : Dummy variable, $\begin{cases} d_{2t} = 1 & \text{when } t=2 \\ d_{2t} = 0 & \text{when } t=1 \end{cases}$

a_i : Unobserved fixed effect (不变量)

u_{it} : Unobserved idiosyncratic error (随机误差)

b. T-period

$$y_{it} = \beta_0 + \delta_0 d_{2t} + \dots + \delta_T d_{Tt} + \beta_1 x_{1it} + \dots + \beta_K x_{Kit} + (a_i + u_{it})$$

先i再t后k

② First-difference model (FD model)

a. Definition

We subtract Two-period / T-period regression functions, then we are able to eliminate fixed effect, and then get the true relationship between x & y

b. Example

△FD for Two-period panel data

$$\text{crmrte}_{i1987} = \beta_0 + \delta_0 \cdot 1 + \beta_1 \text{unem}_{i1987} + a_i + u_{i1987}$$

$$\text{crmrte}_{i1982} = \beta_0 + \delta_0 \cdot 0 + \beta_1 \text{unem}_{i1982} + a_i + u_{i1982}$$

Subtract: $\Rightarrow \Delta \text{crmrte}_i = \delta_0 + \beta_1 \Delta \text{unem}_i + \Delta u_i$

Fixed effect
drops out!
 Δu_i is uncorrelated

The variance of FD estimator may be large if the independent variable's change over the time is small.

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\text{SST}_x}$$

△FD for T-period panel data

Example: Effect of Michigan job training program on worker productivity of manufacturing firms (JTRAIN.dta)

For years 1987, 1988, 1989:

$$\log(\text{scrap}_{it}) = \beta_0 + \beta_1 d_{88t} + \beta_2 d_{89t} + \beta_3 \text{grant}_{it} + \beta_4 \text{grant}_{it-1} + a_i + u_{it}$$

注意把“Michigan job training program”分开，写成两个新的变量

(因为每年不一样)

Additional job training in 1988 made workers more productive in 1989.

Time-invariant reasons why one firm is more productive than another are controlled for. The important point is that these may be correlated with the other independent variables.

First differenced equation using the years 1987, 1988, 1989:

$$\Delta \widehat{\log}(\text{scrap}_{it}) = - .091 - .096 d_{89t} - .223 \Delta \text{grant}_{it} - .351 \Delta \text{grant}_{it-1}$$

(.091) (.125) (.131) (.258)

③ Fixed effects estimation

a. Definition

Since $\alpha_i = \bar{\alpha}_i$, then we can use time-demeaning method to eliminate fixed effect

$$y_{it} = \beta_0 + \beta_1 x_{it1} + \dots + \beta_K x_{itK} + \alpha_i + u_{it}$$

$$\bar{y}_i = \beta_0 + \beta_1 \bar{x}_{i1} + \dots + \beta_K \bar{x}_{iK} + \bar{\alpha}_i + \bar{u}_i$$

Time demeaning $\rightarrow (y_{it} - \bar{y}_i) = \beta_1 (x_{it1} - \bar{x}_{i1}) + \dots + \beta_K (x_{itK} - \bar{x}_{iK}) + (u_{it} - \bar{u}_i)$ But β_1, \dots, β_K is uncorrelated

At this time, β_1, \dots, β_K are called FE estimator

FE estimator is equivalent to introducing a dummy variable in the original regression and using pooled OLS

b. Example

Fixed-effects estimation using the years 1987, 1988, 1989:

$$\widehat{\log}(scrap_{it}) = -.080 d88_t^* - .247 d89_t^* - .252 grant_{it}^* + .422 grant_{it-1}^*$$

$$(.109) \quad (.133) \quad (.151) \quad (.210)$$

$$n = 162, R^2 = .201$$

Training grants significantly improve productivity (with a time lag)

c. Further discussion

△ The R^2 of the demeaned equation is inappropriate

△ The effect of time-invariant variables cannot be estimated

△ The effect of interactions with time-invariant variables can be estimated

△ Degree of freedom have to be adjusted because the N time average are estimated in addition

d. FD & FE

△ $T=2: FD = FE$

$T>2: FE > FD$

△ FD may be better if there is severe serial correlation in the errors

△ Generally, it's a good idea to compute both

④ Random effects model

a. Assumption

$$\text{Cov}(x_{it}, \alpha_i) = 0 \text{ for } y_{it} = \beta_0 + \beta_1 x_{i1t} + \dots + \beta_K x_{iKt} + \alpha_i + u_{it}$$

b. Downside

We set up an assumption to satisfy zero conditional mean assumption,

but still the Ass 5 "No serial correlation" is not satisfied.

$$\text{Cov}(\alpha_i + u_{it}, \alpha_j + u_{jt}) = \text{Cov}(\alpha_i, \alpha_j) = \sigma^2_\alpha$$

c. Transformation to correct Ass 5.

We introduce $\theta = 1 - \sqrt{\frac{\sigma^2_u}{\sigma^2_u + \sigma^2_\alpha}}$ ($\hat{\theta} = 1 - \sqrt{\frac{\widehat{\sigma}^2_u}{\widehat{\sigma}^2_u + \widehat{\sigma}^2_\alpha}}$), and then we have:

$$(y_{it} - \theta \bar{y}_i) = \beta_0(1-\theta) + \beta_1(x_{i1t} - \theta \bar{x}_{i1}) + \dots + \beta_K(x_{iKt} - \theta \bar{x}_{iK}) + (\alpha_i - \theta \bar{\alpha}_i) + u_{it} - \theta \bar{u}_i$$

This equation is called Quasi-demeaned equation, the coefficients are called RE estimator

$$\theta \in [0, 1] \quad \begin{cases} \rightarrow 0 & OLS \\ \rightarrow 1 & FE \end{cases}$$