

Econometrics 1. Introduction & Review of Probability and Statistics

Types of Data Sets (P5-6)

• Cross-Sectional Data

- Data collected at a particular point in time across different economic units.
- The ordering of the data doesn't matter in these samples and contains no information

• Time Series Data

- Data collected across time for the same economic unit.
- e.g. $\begin{cases} \text{macroeconomic data: inflation, unemployment, GDP} \\ \text{financial data: interest rates, exchange rates, stock market indices} \end{cases}$

• Panel Data

- Observations of a collection of units over several time periods.
- e.g. $\begin{cases} \text{Household expenditure and income surveys over several years} \\ \text{Market shares of firms in different years} \\ \text{Trade deficits of countries over several years} \end{cases}$

Review of Probability (Chapter 2)

$$\text{Skewness} = \frac{E[(Y - \mu_Y)^3]}{\sigma_Y^3} \quad \text{how much a distribution deviates from symmetry}$$

- $=0$: symmetric
- <0 : left tail is longer
- >0 : right tail is longer

$$\text{Kurtosis} = \frac{E[(Y - \mu_Y)^4]}{\sigma_Y^4} \quad \text{how much mass is in its tails / how much of the variance arises}$$

- $=3$ normal distribution
- >3 : heavy tails ("leptokurtic")
more values in the dist. tails and more values close to the mean (i.e. sharply peaked with heavy tails)
- <3 : ("platykurtic") fewer values in the tails & fewer values close to the mean (i.e. flat peak and more dispersed scores with lighter tails)



$$X \perp \!\!\! \perp Z \Rightarrow \text{Cov}(X, Z) = 0. \quad \text{corr}(X, Z) = \frac{\text{Cov}(X, Z)}{\sqrt{\text{Var}(X) \text{Var}(Z)}} = \frac{0}{\sigma_X \sigma_Z} = 0$$

$$W = X_1^2 + \dots + X_n^2 \sim \chi_n^2. \quad EW = n, \quad \text{Var}W = 2n.$$

$T = \frac{Z}{\sqrt{n}/\sqrt{n}} \sim t_n.$ t_n is the same as the standard normal dist.

$F = \frac{W_1/n}{W_2/m} \sim F_{n, m}$ $F_{n, m}$ is the same as χ_n^2/n dist.

Review of Statistics (Chapter 3)

1. Estimation (unbiasedness, consistency, efficiency) Chebyshev: $P(|X - EX| \geq a) \leq \frac{\text{Var}(X)}{a^2}$

consistency: $\lim_{n \rightarrow \infty} P(\hat{\mu} - \mu \leq \epsilon) = 1 : \hat{\mu} \xrightarrow{P} \mu.$

efficiency: $\text{Var}(\hat{\mu})$ is smaller than the variance of all other estimators

$E(\bar{Y}) = \mu, \text{Var}(\bar{Y}) = \frac{\sigma^2}{n}$. inversely proportional to n .

the spread of the sampling distribution is proportional to $1/\sqrt{n}$

the sampling uncertainty associated with \bar{Y} is proportional to $1/\sqrt{n}$.

LLN: $P(|\bar{Y} - \mu| < \epsilon) \rightarrow 1$ as $n \rightarrow \infty$ ($\text{Var}(\bar{Y}) = \frac{\sigma^2}{n} \rightarrow 0$)

\bar{Y} is BLUE (proof: SW, Ch. 7)

2. Hypothesis Testing

$$S_Y^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2 \xrightarrow{P} \sigma_Y^2$$

$$\begin{aligned} \text{p-value} &= P_{H_0} [|\bar{Y} - \mu_{H_0}| > |\bar{Y}^{\text{act}} - \mu_{H_0}|] \approx P_{H_0} \left[\left| \frac{\bar{Y} - \mu_{H_0}}{S_Y/\sqrt{n}} \right| > \left| \frac{\bar{Y}^{\text{act}} - \mu_{H_0}}{S_Y/\sqrt{n}} \right| \right] \\ &= P_{H_0} [|t| > |t^{\text{act}}|] \approx 2 \phi(-|t^{\text{act}}|) \quad \begin{cases} \text{p-value} < 5\% \\ |t^{\text{act}}| > 1.96. \end{cases} \end{aligned}$$

right tail: p-value = $1 - \phi(t^{\text{act}}) \Leftrightarrow t^{\text{act}} > 1.645$ (5%)

left tail: p-value = $\phi(t^{\text{act}}) \Leftrightarrow t^{\text{act}} < -1.645$ (5%)

3. Confidence Intervals.

$$\begin{aligned} \textcircled{1} \quad (5\%) : \quad \bar{Y} &\pm 1.96 \cdot \frac{S_Y}{\sqrt{n}}. \quad \left\{ \begin{array}{l} \bar{Y}_m - \bar{Y}_w \sim N(\mu_m - \mu_w, \frac{\sigma_m^2}{n_m} + \frac{\sigma_w^2}{n_w}) \\ t = \frac{\bar{Y}_m - \bar{Y}_w - d_0}{SE(\bar{Y}_m - \bar{Y}_w)}, \text{ where } SE(\bar{Y}_m - \bar{Y}_w) = \sqrt{\frac{S_m^2}{n_m} + \frac{S_w^2}{n_w}} \end{array} \right. \\ \textcircled{2} \quad \text{Compare Means: } H_0: \mu_m - \mu_w = d_0. \quad & \end{aligned}$$

Econometrics 2.1 Introduction to Linear Regression - OLS Estimation.

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

↓ intercept ↓ slope ↓ independent var.

dependent var. not itself economically meaningful

error (omitted factors or measurement error)

$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 \hat{X}_i = \text{OLS prediction} + \text{OLS residual}$

$\Rightarrow \text{Var}(Y_i) = \text{Var}(\hat{Y}_i) + \text{Var}(u_i)$

$\text{explained SS} \quad \text{residual SS}$

$\text{SER: standard error of the regression: } \text{SER} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (\hat{u}_i - \bar{u})^2} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n \hat{u}_i^2}$

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

$R^2 = \frac{\text{ESS}}{\text{TSS}} = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \quad R^2 = 0 : \text{ESS} = 0$

$R^2 = 1 : \text{ESS} = \text{TSS}$

$0 \leq R^2 \leq 1.$

$\text{RMSE} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n \hat{u}_i^2}$

Hint: $n-2$ is "degrees of freedom"

The Least Squares Assumptions (Section 4.4)

① $E(u_i | X_i = x) = 0 \Rightarrow \hat{\beta}_1$ is unbiased

"other factors" contained in u_i should be uncorrelated with X_i .

② (X_i, Y_i) are i.i.d. \Rightarrow delivers the sampling distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$.

③ Large outliers in X and/or Y are rare

Outliers can result in meaningless values of $\hat{\beta}_1$.

Technically $E(X^4) < \infty, E(Y^4) < \infty, X, Y$ are bounded.

OLS sensitive to outlier

The Sampling Distribution of the OLS Estimator (Section 4.5).

$$Y_i = \beta_0 + \beta_1 X_i + u_i \Rightarrow Y_i - \bar{Y} = \beta_1 (X_i - \bar{X}) + (u_i - \bar{u})$$

$$\hat{\beta}_1 = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2} = \frac{\sum (X_i - \bar{X}) [\beta_1 (X_i - \bar{X}) + (u_i - \bar{u})]}{\sum (X_i - \bar{X})^2} = \beta_1 + \frac{\sum (X_i - \bar{X})(u_i - \bar{u})}{\sum (X_i - \bar{X})^2} = \beta_1 + \frac{\sum (X_i - \bar{X})u_i}{\sum (X_i - \bar{X})^2}$$

$$E(\hat{\beta}_1) = \beta_1 + E\left[E\left[\frac{\sum (X_i - \bar{X})u_i}{\sum (X_i - \bar{X})^2} \mid X_1, \dots, X_n\right]\right] = \beta_1. \quad (\text{LSA \#1})$$

$$\hat{\beta}_1 - \beta_1 = \frac{\sum (X_i - \bar{X})u_i}{\sum (X_i - \bar{X})^2} = \frac{\frac{1}{n} \sum v_i u_i}{\left(\frac{1}{n} \sum (X_i - \bar{X})^2\right)}, \text{ where } v_i = (X_i - \bar{X})u_i.$$

$$\text{If } n \text{ is large, } \hat{\beta}_1 - \beta_1 \approx \frac{\frac{1}{n} \sum v_i u_i}{\sigma_x^2} \Rightarrow \text{Var}(\hat{\beta}_1 - \beta_1) = \text{Var}(\hat{\beta}_1) = \frac{\text{Var}(v_i) / n}{(\sigma_x^2)^2} = \frac{1}{n} \cdot \frac{\text{Var}(X_i - \bar{X})u_i^2}{\sigma_x^4}$$

$$\text{Var}(\hat{\beta}_1) \propto \frac{1}{n}, E(\hat{\beta}_1) = \beta_1 \Rightarrow \hat{\beta}_1 \xrightarrow{P} \beta_1.$$

$$\left. \begin{aligned} & \text{if } v_i \text{ i.i.d., } E(v) = 0, \text{Var}(v) = \sigma^2 \Rightarrow \text{CLT} \Rightarrow \frac{1}{n} \sum v_i \sim N(0, \sigma^2 / n) \end{aligned} \right\} \hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{n \sigma_x^4}\right)$$

$$v_i = (X_i - \bar{X})u_i, E(v_i) = 0, \text{Var}(v_i) < \infty$$

Econometrics 2.2 Introduction to Linear Regression - Hypothesis Testing

Test β_1 : $t = \frac{\hat{\beta}_1 - \beta_{1,0}}{SE(\hat{\beta}_1)}$, reject if $|t| > 1.96$ (5%).

Requirement:
Large- n : $n=50$

$$\text{Var}(\hat{\beta}_1) = \frac{\text{Var}[(X_i - \bar{X})u_i]}{n(\sigma_x^2)^2} = \frac{\sigma_u^2}{n\sigma_x^4} \Rightarrow \sigma_{\hat{\beta}_1}^2 = \frac{1}{n} \cdot \frac{\frac{1}{n-2} \sum_{i=1}^n \hat{u}_i^2}{\left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \right]^2} \Rightarrow SE(\hat{\beta}_1) = \sqrt{\sigma_{\hat{\beta}_1}^2}$$

$$\sigma_{\hat{\beta}_1}^2 = \frac{1}{n} \cdot \frac{\frac{1}{n-2} \sum_{i=1}^n (X_i - \bar{X})\hat{u}_i^2}{\left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \right]^2}$$

Confidence interval (5%): $\hat{\beta}_1 \pm 1.96 \cdot SE(\hat{\beta}_1)$

Regression when X is Binary (Section 5.3)

$$Y_i = \beta_0 + \beta_1 X_i + u_i \quad \begin{cases} X_i = 0, \quad Y_i = \beta_0 + u_i, \quad E(Y_i | X_i=0) = \beta_0 \\ X_i = 1, \quad Y_i = \beta_0 + \beta_1 + u_i, \quad E(Y_i | X_i=1) = \beta_0 + \beta_1 \end{cases}$$

β_1 = population difference in group means

Heteroskedasticity (異方差) & Homoskedasticity (同方差) (Section 5.4)

$\text{Var}(u | X=x)$ is constant (does not depend on X)

If Homoskedasticity $\Rightarrow \text{Var}(u_i | X_i=x) = \sigma_u^2$

$$\Rightarrow \text{Var}(\hat{\beta}_1) = \frac{\text{Var}[(X_i - \bar{X})u_i]}{n(\sigma_x^2)^2} = \frac{E[(X_i - \bar{X})^2 u_i^2]}{n(\sigma_x^2)^2} = \frac{\sigma_u^2}{n \sigma_x^2}$$

$$\Rightarrow SE(\hat{\beta}_1) = \sqrt{\frac{1}{n} \cdot \frac{\frac{1}{n-2} \sum \hat{u}_i^2}{\frac{1}{n} \sum (X_i - \bar{X})^2}}$$

Advantage: formula is simpler

Worng: too small

Additional Theoretical Foundations of OLS (Section 5.5)

The extended Least Square Assumptions:

1. $E(u | X=x) = 0$ 2. (X_i, Y_i) i.i.d. 3. Large outliers are rare ($E(Y^4) < \infty, E(X^4) < \infty$)

4. u is homoskedastic 5. u is distributed $N(0, \sigma_u^2)$

① Under Assumption 1-4: $\hat{\beta}_1$ is the BLVE (Gauss-Markov theorem) $\xrightarrow{\text{PROVE}}$ Appendix 5.2

$$\text{smallest var}(\hat{\beta}_1). \quad \hat{\beta}_1 - \beta_1 = \frac{\sum (X_i - \bar{X})u_i}{\sum (X_i - \bar{X})^2} = \frac{1}{n} \sum w_i u_i, \text{ where } w_i = \frac{X_i - \bar{X}}{\sum (X_i - \bar{X})^2}$$

② Under Assumption 1-5: $\hat{\beta}_1$ has smallest var of all consistent estimator (linear or nonlinear).

③ not-so-good thing about OLS. — $\begin{cases} \text{GM condition not hold in practice} \rightarrow \text{homoskedasticity} \\ \text{② require normal errors.} \end{cases}$

$\xrightarrow{\text{異方差}} \text{Weighted least squares more efficient}$

Econometrics 3.1 Introduction to Multiple Regression (Chapters 6)

Omitted Variable Bias (Section 6.1)

It arises because of factors influence Y are not included \rightarrow omitted variables.

The bias in the OLS estimator \rightarrow omitted variable bias.

omitted variable "Z" $\begin{cases} \text{A determinant of } Y \text{ (i.e. part of } u) \\ \text{Correlated with } X \text{ (i.e. } \text{cov}(Z, X) \neq 0) \end{cases}$ $\hat{\beta}_1 - \beta_1$

$$\hat{\beta}_1 - \beta_1 = \frac{\sum (X_i - \bar{X}) u_i}{\sum (X_i - \bar{X})^2} = \frac{\frac{1}{n} \sum u_i}{\frac{1}{n} \sum (X_i - \bar{X})^2} \quad (u_i = (X_i - \bar{X}) u_i \approx (X_i - \bar{X}) u_i).$$

Under Assumption 1, $E[(X_i - \bar{X}) u_i] = \text{cov}(X_i, u_i) = 0$.

What if $E[(X_i - \bar{X}) u_i] = \text{cov}(X_i, u_i) = \sigma_{Xu} \neq 0$?

$$\text{In general, } \hat{\beta}_1 - \beta_1 = \frac{\frac{1}{n} \sum (X_i - \bar{X}) u_i}{\frac{1}{n} \sum (X_i - \bar{X})^2} \xrightarrow{P} \frac{\sigma_{Xu}}{\sigma_X^2} = \left(\frac{\sigma_{Xu}}{\sigma_X} \right) \left(\frac{\beta_1}{\sigma_X} \right) \xrightarrow{P} \frac{\sigma_{Xu}}{\sigma_X} \beta_1$$

$$\hat{\beta}_1 \xrightarrow{P} \beta_1 + \frac{\sigma_{Xu}}{\sigma_X} \beta_1.$$

To overcome: ① Run a randomized controlled experiment
② Add omitted variable \rightarrow multiple regression

The population multiple regression model (Section 6.2)

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i.$$

$$\beta_1 = \frac{\partial Y}{\partial X_1} \text{ holding } X_2 \text{ const.} \quad \beta_2 = \frac{\partial Y}{\partial X_2} \text{ holding } X_1 \text{ const.}$$

β_0 = predicted value of Y when $X_1 = X_2 = 0$

$$\text{OLS: } \hat{\beta} = (X'X)^{-1} X' Y$$

Measures of Fit for Multiple Regression (Section 6.4)

$$\text{Actual} = \text{predicted} + \text{residual: } Y_i = \hat{Y}_i + \hat{u}_i. \quad \begin{cases} \text{SER} = \text{std. deviation of } \hat{u}_i \text{ (d.f. correction)} \\ \text{RMSE} = \text{std. deviation of } \hat{u}_i \text{ (d.f. NOT correction)} \\ R^2, \text{ adj-}R^2 = \text{adj-}R^2 = \hat{R}^2 < R^2. \end{cases}$$

$$SER = \sqrt{\frac{1}{n-k-1} \sum \hat{u}_i^2} \quad RMSE = \sqrt{\frac{1}{n} \sum \hat{u}_i^2} \quad R^2 = \frac{ESS}{TSS} = 1 - \frac{SSR}{TSS} \quad ESS = \sum (\hat{Y}_i - \bar{Y})^2, \quad SSR = \sum \hat{u}_i^2, \quad TSS = \sum (Y_i - \bar{Y})^2.$$

$$\hat{R}^2 = 1 - \left(\frac{n-1}{n-k-1} \right) \frac{SSR}{TSS}$$

The least squares Assumption (Section 6.5)

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + u_i, \quad i=1, \dots, n.$$

$$\textcircled{1} \quad E(u | X_1 = x_1, \dots, X_k = x_k) = 0 \quad \textcircled{2} \quad (X_{1i}, \dots, X_{ki}, Y_i) \text{ i.i.d.} \quad \textcircled{3} \quad E(X_{1i}^4) < \infty, \dots, E(X_{ki}^4) < \infty, E(Y_i^4) < \infty$$

$$\textcircled{4} \quad \text{No Multicollinearity.}$$

- # 1. $E(u|X_1=x_1, \dots, X_k=x_k)=0$
 If omitted variable (z) belongs in u (z) correlated with $X \Rightarrow$ condition fails
- # 2. $(X_{1i}, \dots, X_{ki}, Y_i)$ \sim N .
- # 3. Large outliers are rare
- # 4. No perfect multicollinearity.

The sampling Distribution of the OLS Estimator (see section 6.6).

$$\hat{\beta} = (X'X)^{-1}X'Y = (X'X)^{-1}X'(X\beta + u) = \beta + (X'X)^{-1}X'u$$

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, \sigma^2 \Sigma_{\beta\beta}^{-1})$$

Multicollinearity (Section 6.7)

The Dummy Variable Trap

Solution ① Omit one of the groups ② Omit the intercept

Perfect multicollinearity: "drop" one
 Imperfect multicollinearity \rightarrow highly correlated \rightarrow reg. coefficients imprecisely estimated

$$\text{Var}(\hat{\beta}_i) = \frac{\sigma^2}{n} \cdot \left[\frac{1}{1 - R_{X_i, X_i}^2} \right] \cdot \frac{1}{\sigma_{X_i}^2} \quad \text{譯本(英) P215}$$

Econometrics 3.2 Multiple Regression - Hypothesis Testing (Chapter 7)

Hypothesis Tests and CI for a single coefficient (Section 7.1)

$$t = \frac{\hat{\beta}_1 - \beta_{1,0}}{SE(\hat{\beta}_1)}, \quad CI: \hat{\beta}_1 \pm 1.96 \cdot SE(\hat{\beta}_1). \quad (5\%)$$

$$\hat{Var}(\hat{\beta}) = \frac{1}{n} \cdot \left(\frac{1}{n} \sum X_i X_i' \right)^{-1} \left(\frac{1}{n-k-1} \sum X_i X_i' u_i u_i' \right) \left(\frac{1}{n} \sum X_i X_i' \right)^{-1}$$

$$SE(\hat{\beta}) = \sqrt{\hat{Var}(\hat{\beta})}, \quad SE(\hat{\beta}_1) = \sqrt{\hat{Var}(\hat{\beta})}_{11}.$$

Tests of Joint Hypotheses (Section 7.2)

$$H_0: \beta_1 = \beta_2 = 0 \quad \text{vs. } H_1: \text{either } \beta_1 \neq 0 \text{ or } \beta_2 \neq 0.$$

$$\text{P}(|t_1| > 1.96 \text{ or } |t_2| > 1.96) = P(|t_1| > 1.96, |t_2| > 1.96) + P(|t_1| > 1.96, |t_2| \leq 1.96) + P(|t_2| \leq 1.96, |t_1| > 1.96) \\ = 0.05 \times 0.05 + 0.95 \times 0.95 \times 0.05 = 0.0975 \rightarrow \text{not } 0.05$$

Solution ① Bonferroni (Appendix 7.1) \Rightarrow F-statistic.

$$F = \frac{1}{2} \left(\frac{t_1^2 + t_2^2 - 2 \hat{P}_{t_1, t_2} t_1 t_2}{1 - \hat{P}_{t_1, t_2}^2} \right), \quad \hat{P}_{t_1, t_2} = \text{correlation between } t_1 \text{ & } t_2$$

$$\text{if } \hat{P}_{t_1, t_2} \rightarrow 0 \quad (P_{t_1 \perp \perp t_2}) \quad F \approx \frac{1}{2} (t_1^2 + t_2^2) \rightarrow \frac{\chi_{q_0}^2}{q_0} \text{ or } F_{q_0, \infty}$$

critical value (PPT 3.2 P14)

If homoskedasticity only F-statistic.

$$F = \frac{(R^2_{\text{unrestricted}} - R^2_{\text{restricted}}) / q_0}{(1 - R^2_{\text{unrestricted}}) / (n - k_{\text{unrestricted}} - 1)} \quad \begin{matrix} \text{number of restrictions under } H_0 \\ \text{number of regressors.} \end{matrix}$$

$$= \frac{(SSR_{\text{restricted}} - SSR_{\text{unrestricted}}) / q_0}{SSR_{\text{unrestricted}} / (n - k_{\text{unrestricted}} - 1)}$$

Testing Single Restrictions on multiple Coefficients (Section 7.3).

$$H_0: \beta_1 = \beta_2$$

Rearrange ("transform") the regression.

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i = \beta_0 + (\beta_1 - \beta_2) X_{1i} + \beta_2 (X_{1i} + X_{2i}) + u_i \\ = \beta_0 + \beta_1' X_{1i} + \beta_2 W_i + u_i$$

↓ Test

— 3.2 - 1 —

More in Matrix Setup (Section 18.3)

Joint Hypothesis: $R\hat{\beta} = \gamma$

Heteroskedasticity-robust F : $(R\hat{\beta} - \gamma)' [R \text{var}(\hat{\beta}) R']^{-1} (R\hat{\beta} - \gamma) / q$
(If 4 LS assumptions hold) $\xrightarrow{d} F_{q, \infty}$.

Confidence Sets for Multiple Coefficients (Section 7.4)

Joint confidence set of (β_1, β_2)

Let $F(\beta_{1,0}, \beta_{2,0})$ be the (heteroskedasticity-robust) F -stat.

95% confidence set = $\{(\beta_{1,0}, \beta_{2,0}) : F(\beta_{1,0}, \beta_{2,0}) < 3.00\}$
an ellipse: $\{(\beta_1, \beta_2) : \frac{1}{2} \left(\frac{t_1^2 + t_2^2 - 2\hat{P}_{t_1, t_2} t_1 t_2}{1 - \hat{P}_{t_1, t_2}^2} \right) < 3.00\}$
 $F = \frac{1}{2(1 - \hat{P}_{t_1, t_2}^2)} \left[\left(\frac{\hat{\beta}_1 - \beta_{1,0}}{SE(\hat{\beta}_1)} \right)^2 + \left(\frac{\hat{\beta}_2 - \beta_{2,0}}{SE(\hat{\beta}_2)} \right)^2 - 2\hat{P}_{t_1, t_2} \left(\frac{\hat{\beta}_1 - \beta_{1,0}}{SE(\hat{\beta}_1)} \right) \left(\frac{\hat{\beta}_2 - \beta_{2,0}}{SE(\hat{\beta}_2)} \right) \right]$

Digression about measures of fit.

High $R^2(\bar{R}^2)$ $\left\{ \begin{array}{l} \text{regression explains variation in } Y \\ \text{NOT mean eliminated omitted variable bias} \\ \text{NOT mean have unbiased estimator of causal effect } (\beta_i) \\ \text{NOT mean included variables are statistically significant} \\ \text{(must use hypothesis tests)} \end{array} \right.$

Econometrics 4 Non-linear Regression

The general nonlinear population regression function

$$Y_i = f_i(X_{1i}, X_{2i}, \dots, X_{ki}) + \nu_i, \quad i=1, \dots, n.$$

Assumptions.

1. $E(\nu_i | X_{1i}, \dots, X_{ki}) = 0 \Rightarrow f$ is the conditional expectation of Y given X 's.

2. $(X_{1i}, \dots, X_{ki}, Y_i)$ i.i.d.

3. Big outliers are rare

4. No perfect multicollinearity

① Polynomials ② Logarithmic transformation ("percentages" interpretation)

$$1. \text{ Polynomials } Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{1i}^2 + \dots + \beta_r X_{1i}^r + \nu_i$$

Choice of r : ① Plot, t-/F-test / Check sensitivity / Judgment ② Model selection

2. Logarithmic functions of Y / X .

Reason: $\ln(X + \Delta X) - \ln(X) = \ln\left(1 + \frac{\Delta X}{X}\right) \approx \frac{\Delta X}{X}$

$$I. \quad Y_i = \beta_0 + \beta_1 \ln(X_i) + \nu_i \quad II. \quad \ln(Y_i) = \beta_0 + \beta_1 \ln(X_i) + \nu_i \quad IV. \quad \ln(Y_i) = \beta_0 + \beta_1 \ln(X_i) + \nu_i$$

linear = log

log - linear

log - log

$$I. \quad Y_i = \beta_0 + \beta_1 \ln(X_i). \quad \left\{ \begin{array}{l} Y + \Delta Y = \beta_0 + \beta_1 \ln(X + \Delta X) \\ Y = \beta_0 + \beta_1 \ln(X) \end{array} \right. \Rightarrow \Delta Y = \beta_1 \ln\left(1 + \frac{\Delta X}{X}\right) \approx \beta_1 \frac{\Delta X}{X} \quad \left(\beta_1 \approx \frac{\Delta Y}{\Delta X/X} \right)$$

$$II. \quad \ln(Y_i) = \beta_0 + \beta_1 X_i. \quad \left\{ \begin{array}{l} \ln(Y + \Delta Y) = \beta_0 + \beta_1 \ln(X + \Delta X) \\ \ln Y = \beta_0 + \beta_1 \ln(X) \end{array} \right. \Rightarrow \ln\left(1 + \frac{\Delta Y}{Y}\right) \approx \frac{\Delta Y}{Y} \quad \begin{array}{l} \text{percentage change in } X \\ \text{percentage change in } Y \end{array}$$

$$III. \quad \ln(Y_i) = \beta_0 + \beta_1 \ln(X_i) + \nu_i. \quad \beta_1 \approx \frac{\Delta Y/Y}{\Delta X/X}$$

Interactions Between Independent Variables (Section 8.3)

$$(a) 2 \text{ binary: } Y_i = \beta_0 + \beta_1 D_{1i} + \beta_2 D_{2i} + \beta_3 (D_{1i} \times D_{2i}) + \nu_i$$

$$E(Y_i | D_{1i}=0, D_{2i}=d_2) = \underbrace{\beta_0 + \beta_2 d_2}_{\beta_1 + \beta_3 d_2} \quad E(Y_i | D_{1i}=1, D_{2i}=d_2) = \underbrace{\beta_0 + \beta_1 + \beta_2 d_2 + \beta_3 d_2}_{\beta_1 + \beta_3 d_2}$$

$$(b) \text{ binary & continuous: } Y_i = \beta_0 + \beta_1 D_{1i} + \beta_2 X_{1i} + \beta_3 (D_{1i} \times X_{1i}) + \nu_i$$

increment to the effect
of D_1 when $D_2 = 1$

$$\text{if } D_1=0: \quad Y_i = \beta_0 + \beta_2 X_{1i} + \nu_i$$

$$\text{if } D_1=1: \quad Y_i = (\beta_0 + \beta_1) + (\beta_2 + \beta_3) X_{1i} + \nu_i$$

$$\text{if } D_1=0: \quad Y_i = \beta_0 + \beta_1 D_{1i} + \beta_2 X_{1i} + \beta_3 (D_{1i} \times X_{1i})$$

$$\text{if } D_1=1: \quad Y_i = (\beta_0 + \beta_1) + (\beta_2 + \beta_3) X_{1i} + \nu_i$$

$$(c) 2 \text{ continuous}$$

Econometrics 5 Assessing Studies based on multiple Regression

✓ Internal validity : its statistical inferences about causal effects are valid for population being studied

External validity : its statistical inferences can be generalized to other populations and settings

{ estimator of causal effect : unbiased and consistent

hypothesis tests should have the desired significance level

OLS : interval validity of unbiased and consistent
standard errors are computed correctly

- Threats.
- ① omitted variable bias
 - ② misspecification of the functional form
 - ③ measurement error
 - ④ sample selection bias
 - ⑤ simultaneous causality bias

$$\rightarrow E(u_i | x_{i1}, \dots, x_{in}) \neq 0.$$

Omitted variable bias

(i) a determinant of Y (i.e.) correlated with at least one included regressor

solution: ① include ② instrumental variable ③ randomized controlled experiment

Misspecification of the functional form

① continuous : log., interaction, etc. ② discrete : probit/logit

Errors-in-variable bias

$$Y_i = \beta_0 + \beta_1 X_i + u_i = \beta_0 + \beta_1 \tilde{X}_i + [\beta_1 (X_i - \tilde{X}_i) + u_i] = \beta_0 + \beta_1 \tilde{X}_i + \tilde{u}_i$$

unmeasured true value of X imprecise measured version of X

$$\begin{aligned} \text{cov}(\tilde{X}_i, \tilde{u}_i) &= \text{cov}(\tilde{X}_i, \beta_1 (X_i - \tilde{X}_i) + u_i) = \beta_1 \text{cov}(\tilde{X}_i, X_i - \tilde{X}_i) + \text{cov}(\tilde{X}_i, u_i) \\ &= \beta_1 [\text{cov}(\tilde{X}_i, X_i) - \text{Var}(\tilde{X}_i)] + 0 \neq 0 \end{aligned}$$

Typically not equal

Sample selection Bias

Simultaneous Causality Bias

$$Y_i = \beta_0 + \beta_1 X_i + u_i \quad \xrightarrow{\text{large } u_i \rightarrow \text{large } Y_i (Y_i > 0)} \text{cov}(X_i, u_i) \neq 0.$$

$$X_i = \gamma_0 + \gamma_1 Y_i + v_i$$