

Applied Econometrics

Derek Li

Contents

1	Introduction	3
1.1	Econometric Methodology	3
1.2	Data Types	3
1.2.1	Cross-Section	3
1.2.2	Time Series	3
1.2.3	Repeated Cross-Section	3
1.2.4	Panel Data	3
2	Review of Statistics	4
2.1	Estimator	4
2.2	Sampling Distribution	4
2.3	Confidence Interval	5
2.4	Hypothesis Testing	5
2.4.1	Two Types of Errors	5
2.4.2	Testing Steps	5
2.4.3	p -Value	5
2.5	Proves of Some Theorems and Results (optional)	6
3	Simple Regression	7
3.1	Basic Property	7
3.2	Econometric Model	7
3.3	Estimator: OLS	8
3.4	Properties of OLS	8
3.5	Interpretation	11
4	Multiple Regression: Estimation	12
4.1	Econometric Model	12
4.2	Estimator: OLS	12
4.3	Properties of OLS	12
4.3.1	Omitted Variable Bias	13
4.3.2	Including New Variables	13
5	Multiple Regression: Inference	15
5.1	Confidence Interval	15
5.2	t -Test	15
5.3	F -Test	15
5.4	Example in STATA	16

6	Multiple Regression: OLS Asymptotics	17
6.1	Consistency	17
6.2	Asymptotic Distribution	18
7	Multiple Regression: Further Issues	19
7.1	Beta Coefficients	19
7.2	Qualitative Information: Binary/Dummy Variables	19
8	Heteroskedasticity	20
8.1	Heteroskedasticity-Robust SE	20
8.2	Generalized Least Squares (GLS)	20
8.3	Testing for Heteroskedasticity	21
9	Instrumental Variables	22
9.1	Comparison of OLS and IV	23
9.2	Weak IV	23
9.2.1	Detection of Weak IV	23
9.3	IV Estimation of the Multiple Regression Model	24
9.4	Two Stage Least Squares (2SLS)	24
9.5	Endogeneity Test	24
9.6	Over-Identification Test	25

1 Introduction

Econometrics are quantitative methods of analyzing and interpreting economic data, which need to combine economic theory, math and statistics, data, and statistical or econometrics software.

We focus on estimating economic relationships, testing hypothesis involving economic behavior, and forecasting the behavior of economic variables.

1.1 Econometric Methodology

- Ask a question - statement of theory or hypothesis
- Specification of economic model
- Specification of econometric model
- Collection of Data
- Estimation of the econometric model
- Hypothesis testing
- Prediction or forecasting

1.2 Data Types

There are different data structures: cross-section, time series, repeated cross-section, and panel data.

1.2.1 Cross-Section

Cross-section consists of a sample of individuals, households, firms, countries, etc, taken at a given point in time. Observations are generally independent draws from the population. It is commonly indexed by i as x_i .

1.2.2 Time Series

Time series consists of observations on a variable or several variables over time. Observations are almost never independent of each other. It is commonly indexed by t as x_t .

1.2.3 Repeated Cross-Section

Repeated cross-section consists of two or more cross-sectional data in different points in time, and is different units in different periods. It is commonly index by it as x_{it} .

Example 1.1. Suppose two cross-sectional household surveys are taking in 1985 and 1990. In 1985, a random sample of households is surveyed for variables such as income. In 1990, a new random sample of households is taken using the same survey questions. This is a repeated cross-section data set.

1.2.4 Panel Data

Panel data consists of a time series for each cross-sectional unit in the data set. Observations are independent among units and dependent over time for each unit. It is commonly index by it as x_{it} .

2 Review of Statistics

Suppose $X \sim \mathcal{N}(\mu, \sigma^2)$. Assume a random sample $\{x_1, \dots, x_n\}$, i.e., identically and independently distributed (i.i.d.).

2.1 Estimator

Definition 2.1 (Statistic). A statistic is a function of the data.

Definition 2.2 (Estimator). An estimator is a statistic that is used to estimate the parameter of interest.

Take μ as an example, the proposed estimator is sample average

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n x_i.$$

2.2 Sampling Distribution

We have

$$\mathbb{E}[\bar{X}_n] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[x_i] = \frac{1}{n} \sum_{i=1}^n \mu = \mu,$$

i.e., \bar{X}_n is an unbiased estimator of μ . Besides,

$$\text{Var}[\bar{X}_n] = \text{Var}\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n^2} \text{Var}\left[\sum_{i=1}^n x_i\right] \stackrel{\text{independent}}{=} \frac{1}{n^2} \sum_{i=1}^n \text{Var}[x_i] = \frac{\sigma^2}{n}.$$

Definition 2.3 (Consistency). Let W_n be an estimator of θ based on a sample Y_1, \dots, Y_n . Then W_n is a consistent estimator of θ if for every $\varepsilon > 0$,

$$P(|W_n - \theta| > \varepsilon) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

As commonly stated, an estimator is called consistent when its sampling distribution becomes more and more concentrated around the parameters of interest as the sample size increases. Note that \bar{X}_n is a consistent estimator of μ .

Theorem 2.1. If $Y_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$, $Y_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$, then

$$Y_1 + Y_2 \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2 + 2\text{Cov}(Y_1, Y_2)),$$

By theorem, we have the sampling distribution

$$\bar{X}_n \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right),$$

provided $X \sim \mathcal{N}(\mu, \sigma^2)$.

2.3 Confidence Interval

Theorem 2.2. If Y has $\mathbb{E}[Y] = \mu$, $\text{Var}[Y] = \sigma^2$, then $Z = \frac{Y - \mu}{\sigma}$ is such that $\mathbb{E}[Z] = 0$, $\text{Var}[Z] = 1$.

By theorem, we have

$$Z = \frac{\bar{X}_n - \mu}{\sqrt{\sigma^2/n}} \sim \mathcal{N}(0, 1).$$

Therefore,

$$\begin{aligned} 1 - \alpha &= P\left(-z_{\frac{\alpha}{2}} \leq Z \leq z_{\frac{\alpha}{2}}\right) = P\left(-z_{\frac{\alpha}{2}} \leq \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \leq z_{\frac{\alpha}{2}}\right) \\ &= P\left(\bar{X}_n - z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}\right), \end{aligned}$$

i.e., $(1 - \alpha)\%$ confidence interval is

$$\left[\bar{X}_n - z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}, \bar{X}_n + z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right].$$

2.4 Hypothesis Testing

2.4.1 Two Types of Errors

We want to test H_0 against H_1 and there are two types of errors:

- (1) Reject H_0 when H_0 is true;
- (2) Reject H_1 when H_1 is true.

There is a trade-off between two types of errors: if $P(\text{Reject } H_0 | H_0)$ decreases, then $P(\text{Reject } H_1 | H_1)$ increases. But we fix the $P(\text{Reject } H_0 | H_0)$ at a very small level, i.e., we reject H_0 when there is strong evidence against it.

2.4.2 Testing Steps

- (1) Fix the $P(\text{Reject } H_0 | H_0)$ at some level α .
- (2) Define a test statistic and we have to know the distribution of test statistic under H_0 .
- (3) Define rejection region, $\{|t| > c\}$, where c is the critical value and

$$\alpha = P(|t| > c | H_0) = 1 - P(-c < t < c | H_0).$$

- (4) Check the result and decide whether reject H_0 or not.

2.4.3 p -Value

Given the observed value of the test statistic, p -value is the smallest significance level α at which the null would be rejected. The smaller the p -value, the greatest the evidence against H_0 .

2.5 Proves of Some Theorems and Results (optional)

Theorem 2.3 (Markov's Inequality). If X is a nonnegative random variable and $a > 0$, then

$$P(X \geq a) \leq \frac{\mathbb{E}[X]}{a}.$$

Proof. Since X is a nonnegative random variable, we have

$$\begin{aligned} \mathbb{E}[X] &= \int_0^\infty xf(x)dx = \mathbb{E}[X] = \int_0^a xf(x)dx + \int_a^\infty xf(x)dx \\ &\geq \int_a^\infty xf(x)dx \geq \int_a^\infty af(x)dx = a \int_a^\infty f(x)dx = aP(X \geq a). \end{aligned}$$

Hence

$$P(X \geq a) \leq \frac{\mathbb{E}[X]}{a}.$$

□

Theorem 2.4 (Chebyshev Inequality). For any $b > 0$,

$$P(|X - \mathbb{E}[X]| \geq b) \leq \frac{\text{Var}[X]}{b^2}.$$

Proof. By Markov's Inequality, we have

$$P((X - \mathbb{E}[X])^2 \geq b^2) \leq \frac{\mathbb{E}[(X - \mathbb{E}[X])^2]}{b^2} = \frac{\text{Var}[X]}{b^2}.$$

Therefore,

$$P(|X - \mathbb{E}[X]| \geq b) \leq \frac{\text{Var}[X]}{b^2}.$$

□

Theorem 2.5 (Weak Law of Large Numbers). Let X_1, \dots, X_n be a sequence of independent random variables with $\mathbb{E}[X_i] = \mu$. Let $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$, then

$$\forall \varepsilon > 0, \lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| > \varepsilon) = 0.$$

Proof. By Chebyshev Inequality, for all $\varepsilon > 0$,

$$P(|\bar{X}_n - \mathbb{E}[\bar{X}_n]| > \varepsilon) \leq \frac{\text{Var}[\bar{X}_n]}{\varepsilon^2} \Leftrightarrow 0 \leq P(|\bar{X}_n - \mu| > \varepsilon) \leq \frac{\sigma^2}{n\varepsilon}.$$

Since

$$\lim_{n \rightarrow \infty} \frac{\sigma^2}{n\varepsilon} = 0,$$

then

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| > \varepsilon) = 0.$$

□

It follows that \bar{X}_n is a consistent estimator of μ .

3 Simple Regression

3.1 Basic Property

Property 3.1. $\sum_{i=1}^n (x_i - \bar{x}) = 0$.

Proof. $\sum_{i=1}^n (x_i - \bar{x}) = \sum_{i=1}^n x_i - n\bar{x} = n\bar{x} - n\bar{x} = 0$. □

Property 3.2. $\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i(y_i - \bar{y}) = \sum_{i=1}^n (x_i - \bar{x})y_i$.

Proof. On the one hand,

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

On the other hand,

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

□

Corollary 3.1. $\sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i(x_i - \bar{x})$.

Property 3.3. $\mathbb{E}[Y] = \mathbb{E}[\mathbb{E}[Y|X]]$.

3.2 Econometric Model

Let (Y, X, U) be random variables with joint distribution s.t. $Y = g(X, U)$, where Y is dependent variable, X is explanatory (regressor/covariate) variable, and U is unobservable variable.

Assumption 1 (Linear in Parameters). $Y = \beta_0 + \beta_1 X + U$.

Example 3.1. $y = \beta_0 + \beta_1 x^2 + u$, then $\frac{\partial y}{\partial x} = 2\beta_1 x$.

Example 3.2. $\ln y = \beta_0 + \beta_1 \ln x + u$, then $\frac{\partial \ln y}{\partial \ln x} = \beta_1 \approx \frac{\Delta y\%}{\Delta x\%}$.

Assumption 2 (Zero Conditional Mean). $\mathbb{E}[U|X] = \mathbb{E}[U]$, and $\mathbb{E}[U] = 0$.

Example 3.3. Let Y be wage, X be training program, and U be ability. If training is assigned randomly, then X and U are fully independent. Nevertheless, if let X be education, then the education may influence the ability so that $\mathbb{E}[U|X=0] \neq \mathbb{E}[U|X=1]$, then A2 is violated.

From A1 and A2, we have

$$\begin{aligned} \mathbb{E}[Y|X] &\stackrel{A1}{=} \mathbb{E}[\beta_0 + \beta_1 X + U|X] = \beta_0 + \beta_1 \mathbb{E}[X|X] + \mathbb{E}[U|X] \\ &\stackrel{A2}{=} \beta_0 + \beta_1 X + \mathbb{E}[U] = \beta_0 + \beta_1 X. \end{aligned}$$

Assumption 3 (Random Sample). $\{(x_i, y_i), i = 1, \dots, n\}$ is i.i.d.

Assumption 4 (Sample Variation). $\{x_1, \dots, x_n\}$ are not all the same.

Assumption 5 (Homoscedasticity). $\text{Var}[U|X] = \sigma_U^2$.

From A1 and A5, we have

$$\text{Var}[Y|X] \stackrel{\text{A1}}{=} \text{Var}[\beta_0 + \beta_1 X + U|X] = \text{Var}[U|X] \stackrel{\text{A5}}{=} \sigma_U^2.$$

3.3 Estimator: OLS

We want to solve

$$\min_{\hat{\beta}_0, \hat{\beta}_1} \sum_{i=1}^n Q(\hat{\beta}_0, \hat{\beta}_1) = \sum_{i=1}^n \left(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right)^2.$$

Let

$$\begin{aligned} \frac{\partial Q}{\partial \hat{\beta}_0} &= - \sum_{i=1}^n 2 \left(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right) = 0, \\ \frac{\partial Q}{\partial \hat{\beta}_1} &= - \sum_{i=1}^n 2 \left(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right) x_i = 0. \end{aligned}$$

Therefore,

$$\hat{\beta}_0 = \bar{y} - \beta_1 \bar{x}.$$

Besides,

$$\begin{aligned} \sum_{i=1}^n \left(x_i y_i - \hat{\beta}_0 x_i - \hat{\beta}_1 x_i^2 \right) &= \sum_{i=1}^n (x_i y_i - x_i \bar{y} + \hat{\beta}_1 x_i \bar{x} - \hat{\beta}_1 x_i^2) \\ &= \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \bar{y} + \sum_{i=1}^n \hat{\beta}_1 x_i \bar{x} - \sum_{i=1}^n \hat{\beta}_1 x_i^2 \\ &= \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} + n \hat{\beta}_1 \bar{x}^2 - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = 0, \end{aligned}$$

and thus

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2} = \frac{\sum_{i=1}^n (x_i y_i - x_i \bar{y})}{\sum_{i=1}^n (x_i^2 - x_i \bar{x})} = \frac{\sum_{i=1}^n x_i (y_i - \bar{y})}{\sum_{i=1}^n x_i (x_i - \bar{x})} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

3.4 Properties of OLS

Property 3.4. $\mathbb{E}[\hat{\beta}_1] = \beta_1$.

Proof. We have

$$\begin{aligned}
\hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(\beta_0 + \beta_1 x_i + u_i)}{\sum_{i=1}^n (x_i - \bar{x})^2} \\
&= \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \left[\sum_{i=1}^n (x_i - \bar{x})\beta_0 + \sum_{i=1}^n (x_i - \bar{x})\beta_1 x_i + \sum_{i=1}^n (x_i - \bar{x})u_i \right] \\
&= \beta_1 + \frac{\sum_{i=1}^n (x_i - \bar{x})u_i}{\sum_{i=1}^n (x_i - \bar{x})^2}.
\end{aligned}$$

By A4, $\sum_{i=1}^n (x_i - \bar{x})^2 \neq 0$, and thus,

$$\begin{aligned}
\mathbb{E} \left[\hat{\beta}_1 | x_1, \dots, x_n \right] &= \mathbb{E} \left[\beta_1 + \frac{\sum_{i=1}^n (x_i - \bar{x})u_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \middle| x_1, \dots, x_n \right] \\
&= \beta_1 + \mathbb{E} \left[\frac{\sum_{i=1}^n (x_i - \bar{x})u_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \middle| x_1, \dots, x_n \right] \\
&= \beta_1 + \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) \mathbb{E}[u_i | x_1, \dots, x_n] \\
&\stackrel{\text{A3}}{=} \beta_1 + \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) \mathbb{E}[u_i | x_i] \\
&\stackrel{\text{A2}}{=} \beta_1 + \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) \cdot 0 = \beta_1.
\end{aligned}$$

Therefore, $\mathbb{E} [\hat{\beta}_1] = \mathbb{E} \left[\mathbb{E} [\hat{\beta}_1 | x_1, \dots, x_n] \right] = \mathbb{E} [\beta_1] = \beta_1$. As a consequence, OLS is unbiased. \square

Property 3.5. $\text{Var} \left[\hat{\beta}_1 | x_1, \dots, x_n \right] = \frac{\sigma_u^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$

Proof. We have

$$\begin{aligned}
\text{Var} [\hat{\beta}_1] &= \text{Var} \left[\beta_1 + \frac{\sum_{i=1}^n (x_i - \bar{x}) u_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \middle| x_1, \dots, x_n \right] \\
&= \frac{1}{\left(\sum_{i=1}^n (x_i - \bar{x})^2 \right)^2} \text{Var} \left[\sum_{i=1}^n (x_i - \bar{x}) u_i \middle| x_1, \dots, x_n \right] \\
&\stackrel{\text{A3}}{=} \frac{1}{\left(\sum_{i=1}^n (x_i - \bar{x})^2 \right)^2} \sum_{i=1}^n (x_i - \bar{x})^2 \text{Var}[u_i | x_1, \dots, x_n] \\
&\stackrel{\text{A3}}{=} \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \text{Var}[u_i | x_i] \stackrel{\text{A5}}{=} \frac{\sigma_u^2}{\sum_{i=1}^n (x_i - \bar{x})^2}.
\end{aligned}$$

□

Notice that larger $\sum_{i=1}^n (x_i - \bar{x})^2$ implies smaller $\text{Var} [\hat{\beta}_1]$.

Theorem 3.1 (Gauss-Markov Theorem). Under A1 to A5, OLS is the best linear unbiased estimator.

Definition 3.1 (Standard Error). Let $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$, $\hat{u}_i = y_i - \hat{y}_i$, $\hat{\sigma}_U^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{u}_i^2$, define

$$\text{SE}(\hat{\beta}_1) = \sqrt{\widehat{\text{Var}}[\hat{\beta}_1]} = \sqrt{\frac{\hat{\sigma}_U^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}.$$

Property 3.6. $\sum_{i=1}^n \hat{u}_i = 0$.

Proof. We have

$$\begin{aligned}
\sum_{i=1}^n \hat{u}_i &= \sum_{i=1}^n (y_i - \hat{y}_i) = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = \sum_{i=1}^n (y_i - \bar{y} + \hat{\beta}_1 \bar{x} - \hat{\beta}_1 x_i) \\
&= n\bar{y} - n\bar{y} + n\hat{\beta}_1 \bar{x} - n\hat{\beta}_1 \bar{x} = 0.
\end{aligned}$$

□

Property 3.7. $\sum_{i=1}^n \hat{u}_i x_i = 0$.

Proof. We have

$$\begin{aligned}
\sum_{i=1}^n \hat{u}_i x_i &= \sum_{i=1}^n (y_i - \hat{y}_i) x_i = \sum_{i=1}^n (y_i - \bar{y} + \hat{\beta}_1 \bar{x} - \hat{\beta}_1 x_i) x_i \\
&= \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} + \hat{\beta}_1 \left(n\bar{x}^2 - \sum_{i=1}^n x_i^2 \right) = 0.
\end{aligned}$$

□

Definition 3.2. $R^2 = 1 - \frac{\text{SSR}}{\text{SST}}$, where $\text{SSR} = \sum_{i=1}^n \hat{u}_i^2$, $\text{SST} = \sum_{i=1}^n (y_i - \bar{y})^2$.

3.5 Interpretation

Table 3.1: Summary of Functional Forms Involving Logarithms

Model	Dependent Variable	Independent Variable	Interpretation of β_1
Level-Level	y	x	$\Delta y = \beta_1 \Delta x$
Level-Log	y	$\ln x$	$\Delta y = \frac{\beta_1}{100} \% \Delta x$
Log-Level	$\ln y$	x	$\% \Delta y = 100 \beta_1 \Delta x$
Log-Log	$\ln y$	$\ln x$	$\% \Delta y = \beta_1 \% \Delta x$

Example 3.4 (Level-Log). Suppose

$$hours = 33 + 45.1 \ln(wage),$$

then a 1% increase in $wage$ increases the weekly hours worked by about 0.451.

Example 3.5 (Log-Level). Suppose

$$\ln(wage) = 2.78 + 0.094educ,$$

then one more year of education increases hourly wage by about 9.4%.

4 Multiple Regression: Estimation

4.1 Econometric Model

Assumption 1 (Linear in Parameters). $Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k + U$. We can write Y as

$$Y = \begin{pmatrix} 1 & X_1 & \cdots & X_k \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} = \mathbf{X}\beta + U,$$

where \mathbf{X} is $1 \times (k+1)$, β is $(k+1) \times 1$.

Assumption 2 (Zero Conditional Mean). $\mathbb{E}[U|\mathbf{X}] = 0$.

Assumption 3 (Random Sample). $\{(y_i, x_{1i}, \cdots, x_{ki}), i = 1, \cdots, n\}$ is i.i.d.

Assumption 4 (No Perfect Collinearity). There is no exact linear relationship among the explanatory variables.

Example 4.1. Let $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + U$, $X_3 = X_1 + X_2$. Then

$$Y = \beta_0 + (\beta_1 + \beta_3)X_1 + (\beta_2 + \beta_3)X_2 + U,$$

which violates No Perfect Collinearity.

Example 4.2. Let $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \beta_3 (X_1 \cdot X_2) + U$, then

$$\frac{\partial Y}{\partial X_1} = \beta_1 + 2\beta_2 X_1 + \beta_3 X_2,$$

and thus the model does not violate No Perfect Collinearity.

Assumption 5 (Homoscedasticity). $\text{Var}[U|X_1, \cdots, X_k] = \sigma_U^2$.

From A1 and A5, we have

$$\text{Var}[Y|\mathbf{X}] = \sigma_U^2.$$

4.2 Estimator: OLS

We want to solve

$$\min_{\hat{\beta}_0, \cdots, \hat{\beta}_k} \sum_{i=1}^n \left(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \cdots - \hat{\beta}_k x_{ki} \right)^2.$$

With F.O.C., we have $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$.

4.3 Properties of OLS

Property 4.1. R^2 increases when the model includes more explanatory variables.

Property 4.2 (Partialling Out). Regress X_1 on X_2, \dots, X_k :

$$X_1 = \alpha_0 + \alpha_1 X_2 + \dots + \alpha_{k-1} X_k + \Omega.$$

Thus

$$\hat{\Omega}_i = X_{i1} - \hat{X}_{i1} = X_{i1} - (\hat{\alpha}_0 + \hat{\alpha}_1 X_{i2} + \dots + \hat{\alpha}_{k-1} X_{ik}).$$

Regress Y_i on $\hat{\Omega}_i$:

$$Y_i = \gamma_0 + \gamma_1 \hat{\Omega}_i + V_i,$$

then

$$\hat{\gamma}_1 = \frac{\sum_{i=1}^n \hat{\Omega}_i Y_i}{\sum_{i=1}^n \hat{\Omega}_i^2} = \hat{\beta}_1.$$

Therefore, $\hat{\beta}_1$ measures the effect of X_1 on Y after X_2, \dots, X_k have been partialled out.

Property 4.3. Under A1 to A4, OLS is unbiased, i.e., $\mathbb{E}[\hat{\beta}_j] = \beta_j, j = 0, \dots, k$.

Property 4.4. Under A1 to A5,

$$\text{Var}[\hat{\beta}_j] = \frac{\sigma_U^2}{\left[\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 \right] (1 - R_j^2)},$$

for $j = 1, \dots, k$, where R_j^2 is the R^2 of the regression of X_j on all other X 's.

4.3.1 Omitted Variable Bias

Suppose the true model is

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + U,$$

with $\mathbb{E}[U|X_1, X_2] = 0$.

Now ignore X_2 and suppose $X_2 = \alpha_0 + \alpha_1 X_1 + V, \mathbb{E}[V|X_1] = 0$, and thus

$$\begin{aligned} Y &= \beta_0 + \beta_1 X_1 + \beta_2(\alpha_0 + \alpha_1 X_1 + V) + U = (\beta_0 + \beta_2 \alpha_0) + (\beta_1 + \beta_2 \alpha_1) X_1 + (U + \beta_2 V) \\ &:= \delta_0 + \delta_1 X_1 + \varepsilon. \end{aligned}$$

If $\mathbb{E}[\varepsilon|X_1] = 0$, then $\tilde{\delta}_1$ for δ_1 is unbiased, i.e.,

$$\mathbb{E}[\hat{\delta}_1] = \delta_1 = \beta_1 + \beta_2 \alpha_1,$$

then $\text{Bias}(\hat{\delta}_1) = \beta_2 \alpha_1$. Because the bias in this case arises from omitting x_2 , we call $\beta_2 \alpha_1$ the omitted variable bias.

4.3.2 Including New Variables

Suppose we have $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + U$, and we consider whether X_{k+1} should be included in the model.

We consider these cases:

- If $\beta_{k+1} = 0$, then we should not include X_{k+1} .
- If $\beta_{k+1} \neq 0$ and X_{k+1} is uncorrelated with all other X 's:
 - There is no omitted variable bias problem.
 - $\Delta R_j^2 = 0$ and σ_U^2 decreases, and thus $\text{Var} \left[\hat{\beta}_j \right]$ decreases.

Hence we should include X_{k+1} .

- If $\beta_{k+1} \neq 0$ and X_{k+1} is correlated with other X 's:
 - Excluding X_{k+1} leads to omitted variable bias.
 - R_j^2 increases and σ_U^2 decreases, and thus the change in $\text{Var} \left[\hat{\beta}_j \right]$ is unclear.

We should still include X_{k+1} to avoid omitted variable bias.

5 Multiple Regression: Inference

Assumption 6 (Normality). $U|\mathbf{X} \sim \mathcal{N}(0, \sigma_U^2)$

With A6 we have

$$Y|\mathbf{X} \sim \mathcal{N}(\mathbf{X}\beta, \sigma_U^2).$$

Property 5.1. Under A1 to A6,

$$\hat{\beta}_j \sim \mathcal{N}\left(\beta_j, \frac{\sigma_U^2}{\left[\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2\right] (1 - R_j^2)}\right),$$

for $j = 1, \dots, k$.

Corollary 5.1. We have

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{\text{Var}[\hat{\beta}_j]}} \sim \mathcal{N}(0, 1).$$

Corollary 5.2. We have

$$T = \frac{\hat{\beta}_j - \beta_j}{\text{SE}(\hat{\beta}_j)} \sim t_{(n-k-1)},$$

$$\text{where } \text{SE}(\hat{\beta}_j) = \sqrt{\frac{\hat{\sigma}_U^2}{\left[\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2\right] (1 - R_j^2)}}, \hat{\sigma}_U^2 = \frac{1}{n-k-1} \sum_{i=1}^n (\hat{u}_i)^2 = \frac{1}{n-k-1} \sum_{i=1}^n (y_i - \mathbf{x}_i \hat{\beta})^2.$$

5.1 Confidence Interval

The $(1 - \alpha)\%$ confidence interval for β_j is

$$\left[\hat{\beta}_j \pm c \cdot \text{SE}(\hat{\beta}_j) \right].$$

5.2 t -Test

We want to test

$$H_0 : \beta_j = \beta_j^0 \text{ against } H_1 : \beta_j \neq \beta_j^0,$$

and we use the test statistic

$$T = \frac{\hat{\beta}_j - \beta_j^0}{\text{SE}(\hat{\beta}_j)}.$$

5.3 F -Test

Suppose the unrestricted model $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + U$. F -test is a joint test to test last q coefficients are zero, i.e., we want to test

$$H_0 : \beta_{k-q+1} = \dots = \beta_k = 0 \text{ against } H_1 : \text{At least one of them is not 0.}$$

Under H_0 , we have the restricted model $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_{k-q} X_{k-q} + U$, and the test statistic is

$$F = \frac{(\text{SSR}_r - \text{SSR}_{ur})/q}{\text{SSR}_{ur}/(n - k - 1)} \sim F_{(q, n-k-1)}$$

under H_0 .

5.4 Example in STATA

Source	SS	df	MS	Number of obs = n $F(\quad , \quad) = F$ Prob > F = p -value of the F -test R-squared = $1 - \frac{SSR}{SST}$ Adj R-squared Root MSE = $\sqrt{\widehat{\sigma}_U^2}$
Model	SSE	k	$\frac{SSE}{k}$	
Residual	SSR	$n - k - 1$	$\frac{SSR}{n-k-1}$	
Total	SST	$n - 1$	$\frac{SST}{n-1}$	

Y	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
X1	$\hat{\beta}_1$	$SE \left(\hat{\beta}_1 \right)$	$\frac{\hat{\beta}_j}{SE(\hat{\beta}_j)}$	p -value of t -test	$\left[\hat{\beta}_j \mp c \cdot SE \left(\hat{\beta}_j \right) \right]$
X2	$\hat{\beta}_2$	$SE \left(\hat{\beta}_2 \right)$			
_cons	$\hat{\beta}_0$	$SE \left(\hat{\beta}_0 \right)$			

6 Multiple Regression: OLS Asymptotics

6.1 Consistency

Definition 6.1. $\hat{\theta}_n$ is a consistent estimator of θ if

$$\forall \varepsilon > 0, P\left(\left|\hat{\theta}_n - \theta\right| < \varepsilon\right) \rightarrow 1,$$

as $n \rightarrow \infty$, and we write $\hat{\theta}_n \xrightarrow{P} \theta$.

Property 6.1. Let $\hat{\theta}_n \xrightarrow{P} \theta, \hat{\alpha}_n \xrightarrow{P} \alpha$, then

- (1) $\hat{\theta}_n + \hat{\alpha}_n \xrightarrow{P} \theta + \alpha$;
- (2) $\hat{\theta}_n \cdot \hat{\alpha}_n \xrightarrow{P} \theta \cdot \alpha$;
- (3) $\frac{\hat{\theta}_n}{\hat{\alpha}_n} \xrightarrow{P} \frac{\theta}{\alpha}$, provided $\alpha \neq 0$.
- (4) $g(\hat{\theta}_n) \xrightarrow{P} g(\theta)$, provided that g is continuous.

Theorem 6.1 (Law of Large Numbers). Let X_1, \dots, X_n be i.i.d. with mean $\mathbb{E}[X] = \mu$, then

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{P} \mathbb{E}[X] = \mu.$$

Corollary 6.1. Let X_1, \dots, X_n be i.i.d. with mean $\mathbb{E}[X] = \mu$, then

$$\text{Var}[\bar{X}] = \frac{\sigma^2}{n} \rightarrow 0,$$

as $n \rightarrow \infty$.

Theorem 6.2. Under A1 to A4, OLS is consistent, i.e.,

$$\hat{\beta}_j \xrightarrow{P} \beta_j, j = 1, \dots, k.$$

Proof of the simple case. We have

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} = \beta_1 + \frac{\sum_{i=1}^n (x_i - \bar{x})u_i/n}{\sum_{i=1}^n (x_i - \bar{x})^2/n}.$$

We have

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})u_i = \frac{1}{n} \sum_{i=1}^n (x_i - \mathbb{E}[x_i] + \mathbb{E}[x_i] - \bar{x})u_i = \frac{1}{n} \sum_{i=1}^n (x_i - \mathbb{E}[x_i])u_i + \frac{1}{n} \sum_{i=1}^n (\mathbb{E}[x_i] - \bar{x})u_i.$$

By A3, $\frac{1}{n} \sum_{i=1}^n (x_i - \mathbb{E}[x_i])u_i$ is i.i.d., and by LLN

$$\frac{1}{n} \sum_{i=1}^n (x_i - \mathbb{E}[x_i])u_i \xrightarrow{P} \mathbb{E}[(X - \mathbb{E}[X])U] = \mathbb{E}[(X - \mathbb{E}[X])(U - \mathbb{E}[U])] = \text{Cov}(X, U).$$

In fact,

$$\begin{aligned}\text{Cov}(X, U) &= \mathbb{E}[(X - \mathbb{E}[X])U] = \mathbb{E}[XU] - \mathbb{E}[U]\mathbb{E}[X] \\ &= \mathbb{E}[\mathbb{E}[XU|X]] = \mathbb{E}[\mathbb{E}[X]\mathbb{E}[U|X]] = 0.\end{aligned}$$

Also,

$$\frac{1}{n} \sum_{i=1}^n (\mathbb{E}[x_i] - \bar{x})u_i = \frac{1}{n} \sum_{i=1}^n (\mathbb{E}[x] - \bar{x})u_i \xrightarrow{P} 0.$$

Besides,

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 \xrightarrow{P} \mathbb{E}[X^2] - (\mathbb{E}[X])^2 = \text{Var}[X].$$

Therefore,

$$\hat{\beta}_1 \xrightarrow{P} \beta_1 + \frac{0}{\text{Var}[X]} = \beta_1.$$

□

6.2 Asymptotic Distribution

Definition 6.2 (Asymptotic Distribution). Let $\{Z_1, \dots, Z_n, \dots\}$ be a sequence of random variables s.t.

$$\forall z, P(Z_n \leq z) = F_{Z_n}(z) \rightarrow F_Z(z) = P(Z \leq z),$$

as $n \rightarrow \infty$. We say Z is an asymptotic distribution of Z_n , and write $Z_n \overset{a}{\sim} Z$.

Theorem 6.3 (Central Limit Theorem). Let $\{X_1, \dots, X_n\}$ be a random sample with mean μ and variance σ^2 , then

$$\frac{\bar{X}_n - \mu}{\sqrt{\sigma^2/n}} \overset{a}{\sim} \mathcal{N}(0, 1).$$

Theorem 6.4. Under A1 to A5,

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{\text{Var}[\hat{\beta}_j]}} \overset{a}{\sim} \mathcal{N}(0, 1),$$

and

$$\frac{\hat{\beta}_j - \beta_j}{\text{SE}(\hat{\beta}_j)} \overset{a}{\sim} \mathcal{N}(0, 1).$$

With the theorem above, the usual CI, t -test and F -test are asymptotically valid.

7 Multiple Regression: Further Issues

7.1 Beta Coefficients

Suppose $y = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \cdots + \hat{\beta}_k x_k + \hat{u}$, then

$$y_i - \bar{y} = \hat{\beta}_1(x_1 - \bar{x}_1) + \cdots + \hat{\beta}_k(x_k - \bar{x}_k) + \hat{u}_i,$$

and thus

$$\frac{y - \bar{y}}{\hat{\sigma}_y} = \frac{\hat{\sigma}_1}{\hat{\sigma}_y} \hat{\beta}_1 \frac{x_1 - \bar{x}_1}{\hat{\sigma}_1} + \cdots + \frac{\hat{\sigma}_k}{\hat{\sigma}_y} \hat{\beta}_k \frac{x_k - \bar{x}_k}{\hat{\sigma}_k} + \frac{\hat{u}}{\hat{\sigma}_y}.$$

Let $\hat{b}_j = \frac{\hat{\sigma}_j}{\hat{\sigma}_y} \hat{\beta}_j$ and $z_j = \frac{x_j - \bar{x}_j}{\hat{\sigma}_j}$, then

$$z_y = \hat{b}_1 z_1 + \cdots + \hat{b}_k z_k + error,$$

where \hat{b}_j is the standardized coefficient or beta coefficient.

Beta coefficients have different interpretation: if x_j increases by one standard deviation, then \hat{y} changes by \hat{b}_j standard deviations.

7.2 Qualitative Information: Binary/Dummy Variables

Example 7.1. Consider the following simple model of hourly wage determination:

$$wage = \beta_0 + \beta_1 female + \beta_2 educ + u.$$

Because $female = 1$ when the person is female and $female = 0$ when the person is male, the parameter β_1 has the following interpretation: β_1 is the difference in hourly wage between females and males, given the same amount of education and the same error term u . In this example, we choose males to be the base group or benchmark group

8 Heteroskedasticity

If homoscedasticity assumption is dropped, the formula for $\text{SE}(\hat{\beta}_j)$ is wrong, OLS is not guaranteed to be BLUE (but still unbiased and consistent), and the t -test and CI are invalid.

8.1 Heteroskedasticity-Robust SE

First, consider the simple regression, and we have

$$\hat{\beta}_1 = \beta_1 + \frac{\sum_{i=1}^n (x_i - \bar{x})u_i}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

Under A1 to A4,

$$\text{Var}[\hat{\beta}_1|x_i] = \frac{1}{SXX^2} \text{Var}\left[\sum_{i=1}^n (x_i - \bar{x})u_i\right] = \frac{1}{SXX^2} \sum_{i=1}^n (x_i - \bar{x})^2 \text{Var}[u_i|x_i] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma_i^2}{SXX^2}.$$

Since

$$\text{Var}[u|x] = \mathbb{E}[u^2|x] - (\mathbb{E}[u|x])^2 \stackrel{A2}{=} \mathbb{E}[u^2|x],$$

then

$$\widehat{\text{Var}}[\hat{\beta}_1|x_i] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \hat{u}_i^2}{SXX^2}.$$

In multiple regression, $\widehat{\text{Var}}[\hat{\beta}_j]$ is a consistent estimator of $\text{Var}[\hat{\beta}_j]$ and thus

$$T = \frac{\hat{\beta}_j - \beta_j}{\text{SE}(\hat{\beta}_j)} \stackrel{a}{\sim} \mathcal{N}(0, 1).$$

In STATA, using `reg y x, robust`.

8.2 Generalized Least Squares (GLS)

Assume that

$$\text{Var}[u|\mathbf{x}] = \sigma_u^2 h(\mathbf{x}),$$

where $h(\mathbf{x})$ is some function of the explanatory variables and must be positive. For example, $h(\mathbf{x}) = e^{\gamma_0 + \gamma_1 x_1 + \dots + \gamma_k x_k}$. Therefore,

$$\sigma_u^2 h(\mathbf{x}) = \text{Var}[u|\mathbf{x}] \stackrel{A2}{=} \mathbb{E}[u^2|\mathbf{x}] \Rightarrow \sigma_u^2 = \frac{\mathbb{E}[u^2|\mathbf{x}]}{h(\mathbf{x})} = \mathbb{E}\left[\left(\frac{u}{\sqrt{h(\mathbf{x})}}\right)^2 \middle| \mathbf{x}\right] := \mathbb{E}[(u^*)^2|\mathbf{x}] = \text{Var}[u^*|\mathbf{x}].$$

For original equation $y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + u$, we divide both side by $\sqrt{h(\mathbf{x})}$ and have

$$\frac{y}{\sqrt{h(\mathbf{x})}} = \beta_0 \frac{1}{\sqrt{h(\mathbf{x})}} + \beta_1 \frac{x_1}{\sqrt{h(\mathbf{x})}} + \dots + \beta_k \frac{x_k}{\sqrt{h(\mathbf{x})}} + \frac{u}{\sqrt{h(\mathbf{x})}},$$

or

$$y^* = \beta_0 x_0^* + \beta_1 x_1^* + \dots + \beta_k x_k^* + u^*$$

with $\text{Var}[u^*|\mathbf{x}] = \sigma_u^2$. Hence, generalized least squares (GLS) estimators are BLUE. Note that in most case we do not know $h(\mathbf{x})$ and thus we use heteroskedasticity-robust.

8.3 Testing for Heteroskedasticity

We want to test

$$H_0 : \text{Var}[u|\mathbf{x}] = \sigma_u^2 \text{ against } H_1 : \text{Var}[u|\mathbf{x}] \neq \sigma_u^2.$$

We regress y on \mathbf{x} and have

$$\hat{u}_i = y_i - \mathbf{x}_i \hat{\beta},$$

and then estimate the regression

$$\hat{u}^2 = \delta_0 + \delta_1 x_1 + \cdots + \delta_k x_k + \text{Interaction terms} + v.$$

Then, we can do the F -test on $H_0 : \delta_1 = \delta_2 = \cdots = 0$.

9 Instrumental Variables

If $\text{Cov}(X, U) \neq 0$, OLS is inconsistent. Some reasons cause $\text{Cov}(X, U) \neq 0$: omitted variable bias, measurement error, simultaneity, and/or sample selection. We can solve this problem by instrumental variables.

Definition 9.1 (Instrument Variable). If Z is observed and is s.t. $\text{Cov}(Z, U) = 0$ (valid instrumental variable) and $\text{Cov}(Z, X) \neq 0$ (relevant instrumental variable), then we call Z an instrument variable.

Example 9.1. Let $Y = \text{wage}$, $X = \text{education}$, $U = \text{unobserved ability}$. $Z = \text{SIN}$ is not IV since $\text{Cov}(Z, X) = 0$. $Z = \text{IQ}$ is not IV since $\text{Cov}(Z, U) \neq 0$. Note that good proxies are bad IV. $Z = \text{tuition subsidies}$ could be IV, since $\text{Cov}(Z, X) \neq 0$, and $\text{Cov}(Z, U)$ may be 0 depends on how subsidies are allocated.

Consider simple regression with $\text{Cov}(X, U) \neq 0$, i.e., dropping A2 (zero conditional mean).

Assumption 2'. $\text{Cov}(Z, U) = 0, \text{Cov}(Z, X) \neq 0$.

Assumption 3'. $\{(y_i, x_i, z_i), i = 1, \dots, n\}$ is i.i.d.

Assumption 5'. $\text{Var}[U|Z] = \sigma_U^2$,

Property 9.1. $\beta_1 = \frac{\text{Cov}(Y, Z)}{\text{Cov}(X, Z)}$ under A1, A2' and A3'.

Proof. We have $\mathbb{E}[Y] = \beta_0 + \beta_1 \mathbb{E}[X] + \mathbb{E}[U]$ and thus

$$\mathbb{E}[Y]\mathbb{E}[Z] = \beta_0 \mathbb{E}[Z] + \beta_1 \mathbb{E}[X]\mathbb{E}[Z] + \mathbb{E}[U]\mathbb{E}[Z].$$

Besides, since $YZ = \beta_0 Z + \beta_1 XZ + UZ$, then

$$\mathbb{E}[YZ] = \beta_0 \mathbb{E}[Z] + \beta_1 \mathbb{E}[XZ] + \mathbb{E}[UZ].$$

Hence,

$$\mathbb{E}[YZ] - \mathbb{E}[Y]\mathbb{E}[Z] = \beta_1 (\mathbb{E}[XZ] - \mathbb{E}[X]\mathbb{E}[Z]) + \mathbb{E}[UZ] - \mathbb{E}[U]\mathbb{E}[Z],$$

i.e.,

$$\text{Cov}(Y, Z) = \beta_1 \text{Cov}(X, Z) + \text{Cov}(U, Z) \Rightarrow \beta_1 = \frac{\text{Cov}(Y, Z)}{\text{Cov}(X, Z)}.$$

□

The estimator is

$$\hat{\beta}_1^{\text{IV}} = \frac{\frac{1}{n} \sum_{i=1}^n (z_i - \bar{z})(y_i - \bar{y})}{\frac{1}{n} \sum_{i=1}^n (z_i - \bar{z})(x_i - \bar{x})}.$$

Property 9.2. Under A1, A2' and A3',

$$\hat{\beta}_1^{\text{IV}} \xrightarrow{P} \frac{\text{Cov}(Z, Y)}{\text{Cov}(Z, X)} = \beta_1,$$

i.e., $\hat{\beta}_1^{\text{IV}}$ is consistent.

Property 9.3. $\hat{\beta}_1^{\text{IV}}$ is biased.

Property 9.4. Under A1, A2', A3' and A5',

$$\widehat{\text{Var}} \left[\hat{\beta}_1^{\text{IV}} \right] = \frac{\hat{\sigma}_U^2}{\sum_{i=1}^n (x_i - \bar{x})^2 R_{XZ}^2},$$

where R_{XZ}^2 is the R^2 obtained by regression X on Z .

Property 9.5. Under A1, A2', A3' and A5',

$$T = \frac{\hat{\beta}_1^{\text{IV}} - \beta_1}{\text{SE}(\hat{\beta}_1^{\text{IV}})} \stackrel{a}{\sim} \mathcal{N}(0, 1).$$

In STATA, using `ivreg y (x = z), (robust)`.

9.1 Comparison of OLS and IV

- If $\text{Cov}(X, U) = 0$ and $\text{Cov}(Z, U) = 0$, then $\hat{\beta}^{\text{OLS}} \xrightarrow{P} \beta$ and $\hat{\beta}^{\text{IV}} \xrightarrow{P} \beta$. Also under homoscedasticity,

$$\widehat{\text{Var}} \left[\hat{\beta}^{\text{OLS}} \right] = \frac{\hat{\sigma}_u^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \leq \frac{\hat{\sigma}_u^2}{\sum_{i=1}^n (x_i - \bar{x})^2 R_{XZ}^2} = \widehat{\text{Var}} \left[\hat{\beta}^{\text{IV}} \right].$$

Therefore, we use OLS.

- If $\text{Cov}(X, U) \neq 0$ and $\text{Cov}(Z, U) = 0$, then $\hat{\beta}^{\text{OLS}}$ is inconsistent to β and $\hat{\beta}^{\text{IV}} \xrightarrow{P} \beta$. Though $\widehat{\text{Var}} \left[\hat{\beta}^{\text{OLS}} \right] \leq \widehat{\text{Var}} \left[\hat{\beta}^{\text{IV}} \right]$, we use IV.
- If $\text{Cov}(X, U) \neq 0$ and $\text{Cov}(Z, U) \neq 0$, then

$$\hat{\beta}^{\text{OLS}} \xrightarrow{P} \beta + \frac{\text{Cov}(X, U)}{\text{Var}[X]}, \hat{\beta}^{\text{IV}} \xrightarrow{P} \beta + \frac{\text{Cov}(Z, U)}{\text{Cov}(Z, X)}.$$

In general, it is unclear which one is better. But if $\text{Cov}(Z, X)$ is very small, even if $|\text{Cov}(Z, U)| < |\text{Cov}(X, U)|$, we still have

$$\left| \frac{\text{Cov}(X, U)}{\text{Var}[X]} \right| < \left| \frac{\text{Cov}(Z, U)}{\text{Cov}(Z, X)} \right|.$$

9.2 Weak IV

When $\text{Cov}(Z, X)$ is small, we say Z is weak. There are two problems from weak IV:

- $\widehat{\text{Var}} \left[\hat{\beta}^{\text{IV}} \right] = \frac{\hat{\sigma}_u^2}{\sum_{i=1}^n (x_i - \bar{x})^2 R_{XZ}^2}$ becomes larger since R_{XZ}^2 is smaller;
- $\frac{\hat{\beta}^{\text{IV}} - \beta}{\text{SE}(\hat{\beta}^{\text{IV}})} \stackrel{a}{\not\sim} \mathcal{N}(0, 1)$ and thus usual CI, t and F test are invalid.

9.2.1 Detection of Weak IV

We regress X on Z and get the F statistics. If $F < 10$, then Z is a weak IV; otherwise, Z is a strong IV.

9.3 IV Estimation of the Multiple Regression Model

Consider the structural equation

$$Y = \beta_0 + \beta_1 X + \beta_2 Z_1 + \beta_3 Z_2 + \cdots + \beta_k Z_{k-1} + U,$$

where X is endogenous regressor ($\text{Cov}(X, U) \neq 0$), Z_1, \dots, Z_{k-1} are exogenous regressor ($\text{Cov}(Z_j, U) = 0$), and Z_k, \dots, Z_q are instrumental variables ($\text{Cov}(Z_j, U) = 0$). We say $Z = (Z_1, \dots, Z_{k-1}, Z_k, \dots, Z_q)$ is exogenous.

We call

$$X = \pi_0 + \pi_1 Z_1 + \cdots + \pi_q Z_q + V = Z\pi + V$$

as reduced form equation, where $\text{Cov}(Z_j, V) = 0$.

In this case we have

$$\begin{aligned} 0 \neq \text{Cov}(X, U) &= \text{Cov}(Z\pi + V, U) = \text{Cov}(\pi_0 + \pi_1 Z_1 + \cdots + \pi_q Z_q + V, U) \\ &= \pi_1 \text{Cov}(Z_1, U) + \cdots + \pi_q \text{Cov}(Z_q, U) + \text{Cov}(V, U) \\ &= \text{Cov}(V, U). \end{aligned}$$

In STATA, using `y (x = zk zk+1 ... zq) z1 ... zk-1, robust`.

9.4 Two Stage Least Squares (2SLS)

In first stage, we regress X on all Z 's (Z_1, \dots, Z_q) and get

$$\hat{X}_i = Z_i \hat{\pi}.$$

In second stage, we regress Y on \hat{X}_i and all exogenous regressors Z_1, \dots, Z_{k-1} (no IV included).

Mathematically,

$$\begin{aligned} Y &= \beta_0 + \beta_1 X + \beta_2 Z_1 + \cdots + \beta_k Z_{k-1} + U \\ &= \beta_0 + \beta_1 (\pi_0 + \pi_1 Z_1 + \cdots + \pi_{k-1} Z_{k-1} + \pi_k Z_k + \cdots + \pi_q Z_q + V) + \beta_2 Z_1 + \cdots + \beta_k Z_{k-1} + U \\ &= (\beta_0 + \beta_1 \pi_0) + (\beta_1 \pi_1 + \beta_2) Z_1 + \cdots + (\beta_1 \pi_{k-1} + \beta_k) Z_{k-1} + \beta_1 \pi_k Z_k + \cdots + \beta_1 \pi_q Z_q + (U + \beta_1 V) \\ &:= \alpha_0 + \alpha_1 Z_1 + \cdots + \alpha_q Z_q + \varepsilon. \end{aligned}$$

Since $\alpha_k = \beta_1 \pi_k, \dots, \alpha_q = \beta_1 \pi_q$ and $\alpha_j, \pi_j, j = k, \dots, q$ can be estimated, then we can find $\hat{\beta}_1$. For $\alpha_1 = \beta_1 \pi_1 + \beta_2, \dots, \alpha_{k-1} = \beta_1 \pi_{k-1} + \beta_k$, we can also find $\hat{\beta}_i, i = 2, \dots, k$.

In STATA, using `ivreg 2sls y (x = ...)`.

9.5 Endogeneity Test

We know that OLS is BLUE, and unless there is a compelling reason, we should use OLS. If we have a suspicious OLS estimator due to potential endogeneity and a plausible 2SLS estimator, we can test if $\text{Cov}(X, U) = 0$.

We cannot use OLS to test this since the sample covariance between x_i and $\hat{u}_i = y_i - x_i \hat{\beta}^{\text{OLS}}$ is zero by construction $\sum_{i=1}^n \hat{u}_i x_i = 0$. But we can use a valid instrument to test.

Suppose $Y = \beta_0 + \beta_1 X + U$ and $X = \pi_0 + \pi_1 Z + V$ with $\text{Cov}(Z, U) = \text{Cov}(Z, V) = 0$, then $\text{Cov}(X, U) = 0$ iff $\text{Cov}(V, U) = 0$, which is what we want to test. Let

$$U = \delta V + e,$$

where $\text{Cov}(V, e) = 0, \mathbb{E}[e] = 0$. U and V are uncorrelated iff $\delta = 0$.

Therefore, we have

$$Y = \beta_0 + \beta_1 X + \delta V + e,$$

which is called control function and is numerically equivalent to 2SLS.

We take endogeneity test following three steps:

- Step 1: Regress X on Z (OLS) and get $\hat{V} = X - Z\hat{\pi}$.
- Step 2: Regress Y on X and \hat{V} (OLS).
- Step 3: Test $H_0 : \delta = 0$ against $H_1 : \delta \neq 0$ (t test).

9.6 Over-Identification Test

Note that control function depends on $\text{Cov}(Z, U) = 0$ but we cannot use IV estimator to test it. In the data, the covariance between z_i and $\hat{u}_i = y_i - x_i \hat{\beta}^{\text{IV}}$ is zero by construction $\sum_{i=1}^n \hat{u}_i z_i = 0$.

Suppose $Y = \beta_0 + \beta_1 X + U$ and $X = \pi_0 + \pi_1 Z_1 + \pi_2 Z_2 + V$ with

$$\text{Cov}(Z_1, U) = \text{Cov}(Z_2, U) = \text{Cov}(Z_1, V) = \text{Cov}(Z_2, V) = 0.$$

If we run 2SLS using Z_1 or Z_2 only, we have $\tilde{\beta}_1^{\text{IV}}$ and $\check{\beta}_1^{\text{IV}}$, respectively. Under the assumption that both Z_1 and Z_2 are valid IVs, we should have $\tilde{\beta}_1^{\text{IV}} \approx \check{\beta}_1^{\text{IV}}$. If they are different, then either Z_1 or Z_2 or both are invalid.

We take over-identification test following three steps:

- Step 1: Estimate structural equation using all IVs and get \hat{U} (2SLS).
- Step 2: Regress \hat{U} on all exogenous variables (OLS) and get R^2 .
- Step 3: Test $H_0 : \text{Cov}(Z_1, U) = \text{Cov}(Z_2, U) = 0$ against H_1 : At least one of covariance is not zero. Under $H_0, nR^2 \stackrel{a}{\sim} \chi_{(q)}^2$ where $q = \text{Number of IVs} - \text{Number of endogenous regressors}$ and $n = \text{Number of observations}$.

We cannot test if we have one IV. If we have more IVs than endogenous regressors, we can run the over-identification test. If we reject H_0 , we only know either Z_1 or Z_2 or both are invalid but we do not know which one is invalid. If we do not reject H_0 , it does not guarantee that both IVs are valid because they could both be invalid and deliver similar estimates, $\tilde{\beta}_1^{\text{IV}} \approx \check{\beta}_1^{\text{IV}}$.