

Lecture 4

Multiple Linear Regression: Asymptotics

CHUNG-MING KUAN

Department of Finance & CRETA

National Taiwan University

February 10, 2022

Lecture Outline

1 Multiple Linear Regression: Asymptotics

- Limitation of Classical Assumptions
- Review of LLN and CLT
- Consistency
- Asymptotic Normality
- Estimation of Asymptotic Covariance Matrix
- Large Sample Tests
- Tests of Conditional Homoskedasticity

2 Instrumental Variable Estimation

- Endogeneity
- Instrumental Variable Estimation
- Two-Stage Least Squares Estimation

Are Classical Assumptions Reasonable?

- Non-random x_j : If y and x_j are all economic variables, it is **not** reasonable to assume that, while y is random, x_j are all non-random. When x_j are random variables, it would be difficult to infer the statistical properties of the OLS estimators.
 - $\mathbb{E}(\hat{\beta}) = \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}]$ cannot be expressed as $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbb{E}(\mathbf{y})$.
 - $\text{var}(\hat{\beta}) = \text{var}((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y})$ cannot be expressed as $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\text{var}(\mathbf{y})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$.
- When x_j are random, we need to consider the conditional variance $\text{var}(y_i|x_{i1}, \dots, x_{ik})$ and the effect of **conditional heteroskedasticity**, i.e., $\text{var}(y_i|x_{i1}, \dots, x_{ik})$ change with some variables x_{ij} , for deriving $\text{var}(\hat{\beta})$.

- Normality: There is no guarantee that the data we consider are normally distributed; for example, binary variables (taking values zero and one) and count data (taking finitely many values: 1, 2, ...) are definitely non-normal. Then, the statistics derived earlier no longer have the t or F distribution.
- If \mathbf{X} is random, the OLS estimator $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ need not be normally distributed even when \mathbf{y} is. In fact, $\hat{\beta}$ has an unknown complex distribution when \mathbf{X} is random. This renders hypothesis testing difficult.
- **Question:** Is it possible to draw statistical inference **without** Classical Assumptions?
Ans: Yes, we can derive **asymptotic properties** of the OLS estimator and related statistics under weaker and reasonable conditions.

Law of Large Numbers

Khinchine's Weak Law of Large Numbers (WLLN)

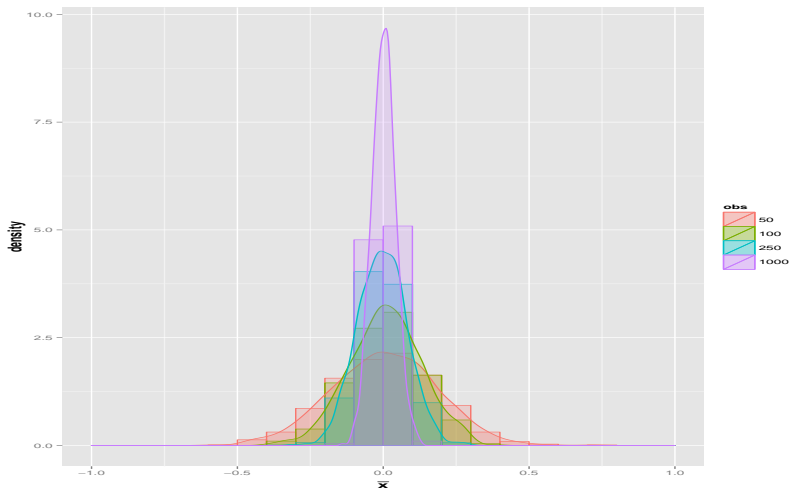
Let z_1, \dots, z_n be i.i.d. with mean μ_o . Then, $\bar{z}_n = n^{-1} \sum_{i=1}^n z_i$ **converges in probability** to μ_o , in the sense that, for any $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}\{|\bar{z}_n - \mu_o| > \varepsilon\} = 0.$$

This is denoted as $\bar{z}_n \xrightarrow{\mathbb{P}} \mu_o$.

- A WLLN ensures that \bar{z}_n , the sample average of z_i , is essentially close to μ_o ; the probability that \bar{z}_n deviates from the true mean by a certain amount approaches zero when the sample size becomes large.
- Note that i.i.d. random variables without a finite mean (e.g. Cauchy variables) do **not** obey a WLLN.

Simulations of LLN: $t(5)$ Variable



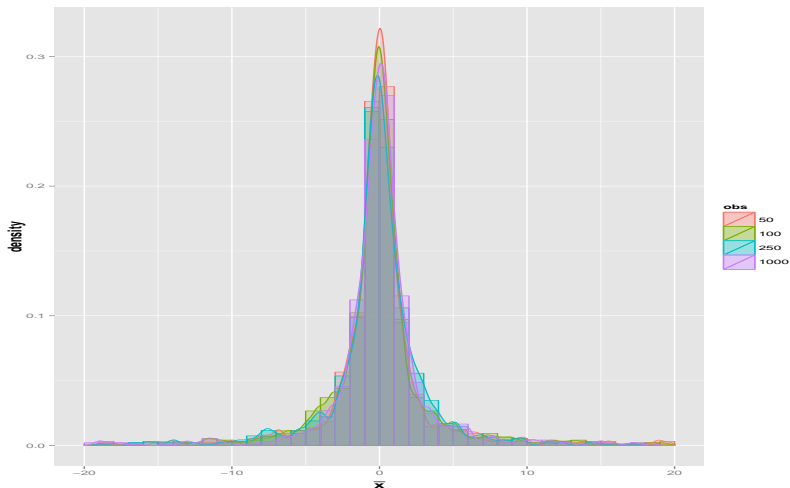
Note: i.i.d. $t(5)$ variables with 1000 replications.

Simulation Procedure

- 1 Generate a sample of n ($n = 100$ for example) realizations z_i from a given distribution F ($t(5)$ in our example) and compute the sample average: $\bar{z}_n = n^{-1} \sum_{i=1}^n z_i$.
- 2 Replicate the step above R ($R = 1000$ for example) times and obtain 1000 sample averages \bar{z}_n .
- 3 The frequency plot of these \bar{z}_n is the simulated (finite-sample) density function of \bar{z}_n under F .

This frequency plot would be more concentrated around the true mean 0 when n becomes larger; that is, it is less likely that \bar{z}_n would be far away from 0 when n is large.

Simulations of LLN: $t(1)$ Variable



Note: i.i.d. $t(1)$ variables with 1000 replications.

Central Limit Theorem

Lindeberg-Lévy's Central Limit Theorem (CLT)

Let z_1, \dots, z_n be i.i.d. with mean μ_o and variance $\sigma_o^2 > 0$. Then,

$$\frac{\sqrt{n}(\bar{z}_n - \mu_o)}{\sigma_o} = \frac{(\bar{z}_n - \mu_o)}{\sigma_o/\sqrt{n}} \xrightarrow{D} \mathcal{N}(0, 1),$$

where \xrightarrow{D} stands for **convergence in distribution**.

A CLT ensures that the distributions of **suitably normalized** sample averages are essentially close to the standard normal distribution in the limit, regardless of the original distributions of z_i . Any random irregularities that are not governed by the standard normal distribution will be wiped out in the limit.

Given that z_1, \dots, z_n are i.i.d. with mean μ_o and variance σ_o^2 , $\mathbb{E}(\bar{z}_n) = \mu_o$, and $\text{var}(\bar{z}_n) = \sigma_o^2/n$. As $\sigma_o^2/n \rightarrow 0$ as $n \rightarrow \infty$, \bar{z}_n ought to have a **degenerate** distribution at μ_o in the limit. Dividing \bar{z}_n by its standard deviation σ_o/\sqrt{n} , we have

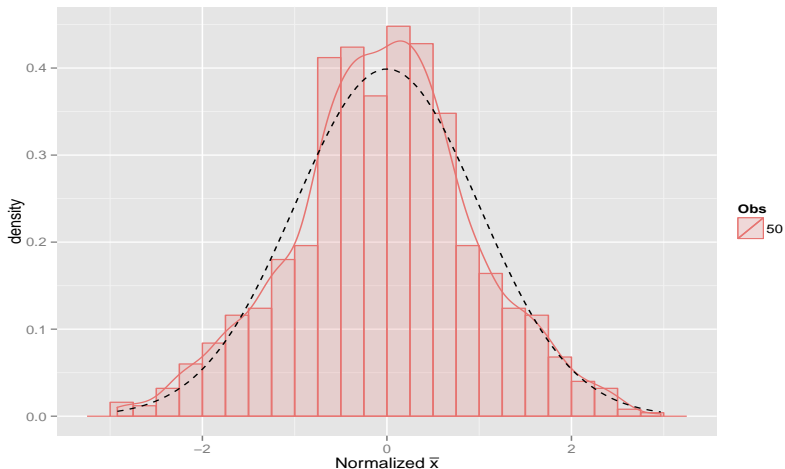
$$\frac{(\bar{z}_n - \mu_o)}{\sigma_o/\sqrt{n}} = \frac{\sqrt{n}(\bar{z}_n - \mu_o)}{\sigma_o},$$

which has mean zero and constant variance one for all n . This prevents the probability mass from shrinking towards a single point in the limit. Note that the factor \sqrt{n} characterizes the **rate of convergence** of \bar{z}_n .

Remarks:

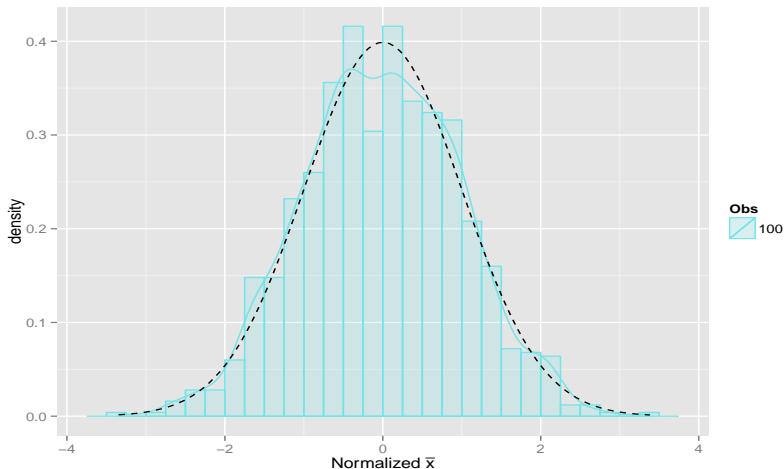
- i.i.d. random variables need **not** obey this CLT if they do not have a finite variance, e.g., $t(2)$ variables.
- Both LLN and CLT may hold for **non-i.i.d.** random variables under stronger conditions (e.g., higher order moments exist).

Simulations of CLT: $t(5)$ Variable; Sample 50



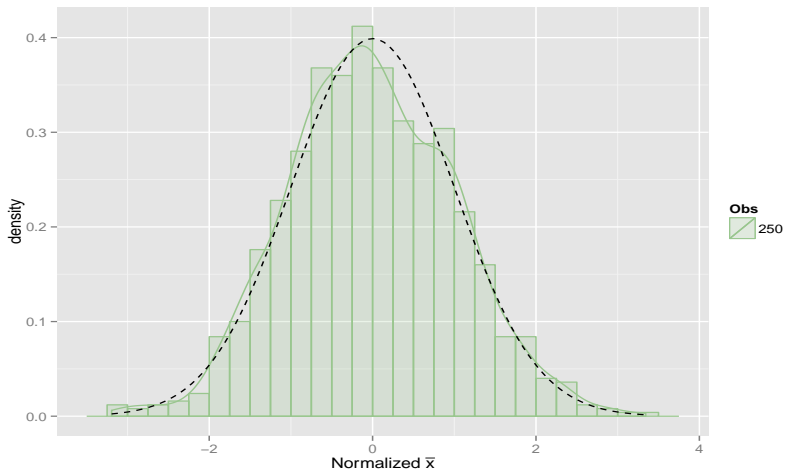
Note: i.i.d. $t(5)$ variables with 1000 replications.

Simulations of CLT: $t(5)$ Variable; Sample 100



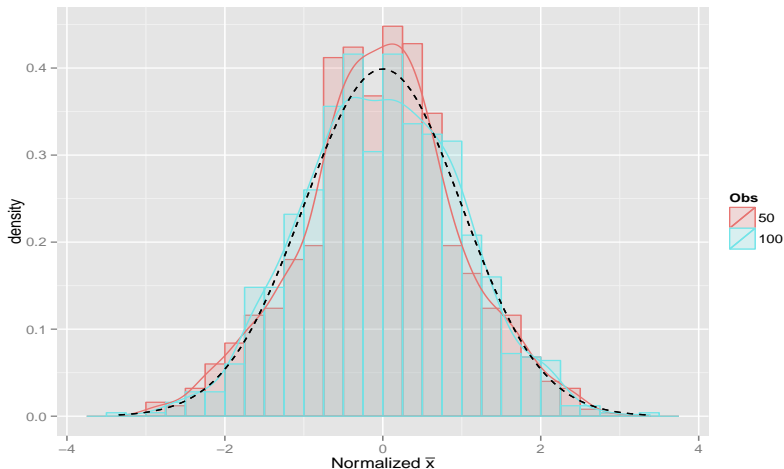
Note: i.i.d. $t(5)$ variables with 1000 replications.

Simulations of CLT: $t(5)$ Variable; Sample 250



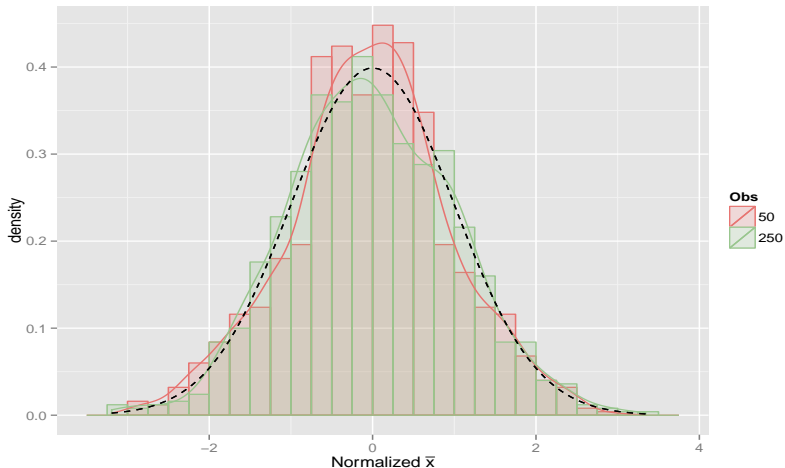
Note: i.i.d. $t(5)$ variables with 1000 replications.

Comparison of CLT: $t(5)$ Variable; Samples 50 & 100



Note: i.i.d. $t(5)$ variables with 1000 replications.

Comparison of CLT: $t(5)$ Variable; Samples 50 & 250



Note: i.i.d. $t(5)$ variables with 1000 replications.

Multivariate Versions of LLN and CLT

- Let $\mathbf{z}_1, \dots, \mathbf{z}_n$ be i.i.d. random vectors with mean $\boldsymbol{\mu}_o$. We say $\bar{\mathbf{z}}_n \xrightarrow{\mathbb{P}} \boldsymbol{\mu}_o$ if each element of $\bar{\mathbf{z}}_n$ converges in probability to the corresponding element in $\boldsymbol{\mu}_o$.
- Let $\mathbf{z}_1, \dots, \mathbf{z}_n$ be i.i.d. with mean $\boldsymbol{\mu}_o$ and positive definite covariance matrix $\boldsymbol{\Sigma}_o$. The multivariate CLT reads:

$$\boldsymbol{\Sigma}_o^{-1/2} \sqrt{n}(\bar{\mathbf{z}}_n - \boldsymbol{\mu}_o) = \boldsymbol{\Sigma}_o^{-1/2} \frac{1}{\sqrt{n}} \sum_{i=1}^n (\mathbf{z}_i - \boldsymbol{\mu}_o) \xrightarrow{D} \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

- **Cramér-Wold Device:** Let $\boldsymbol{\zeta}_n = \boldsymbol{\Sigma}_o^{-1/2} \sqrt{n}(\bar{\mathbf{z}}_n - \boldsymbol{\mu}_o)$ and $\boldsymbol{\nu} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Then, $\boldsymbol{\zeta}_n \xrightarrow{D} \boldsymbol{\nu}$ if, and only if, $\mathbf{a}'\boldsymbol{\zeta}_n \xrightarrow{D} \mathbf{a}'\boldsymbol{\nu}$ for all \mathbf{a} with $\|\mathbf{a}\| = 1$. That is, $\boldsymbol{\zeta}_n$ obeys a multivariate CLT if, and only if, any linear combination of $\boldsymbol{\zeta}_n$ obeys a univariate CLT.

OLS Consistency

Let $\mathbf{x}_i = (1, x_{i1}, \dots, x_{ik})'$ denote the i^{th} column of \mathbf{X}' , we can write the linear specification $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$ as:

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + u_i, \quad i = 1, \dots, n.$$

It follows that the OLS estimator $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{y})$ can be expressed as:

$$\hat{\boldsymbol{\beta}} = \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\sum_{i=1}^n \mathbf{x}_i y_i \right).$$

Modern Assumption I

The random variables y_i are $y_i = \mathbf{x}_i' \mathbf{b}_o + \varepsilon_i$, $i = 1, \dots, n$, such that

- (i) \mathbf{x}_i , $i = 1, \dots, n$, are i.i.d. random vectors with finite second moment;
- (ii) ε_i , $i = 1, \dots, n$, are i.i.d. with $\mathbb{E}(\varepsilon_i) = 0$ and $\mathbb{E}(\mathbf{x}_i \varepsilon_i) = \mathbf{0}$.

Given Modern Assumption I,

$$\begin{aligned}\hat{\beta} &= \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\sum_{i=1}^n \mathbf{x}_i (\mathbf{x}_i' \mathbf{b}_o + \varepsilon_i) \right) \\ &= \mathbf{b}_o + \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\sum_{i=1}^n \mathbf{x}_i \varepsilon_i \right).\end{aligned}$$

Clearly, $\hat{\beta} \xrightarrow{\mathbb{P}} \mathbf{b}_o$ when the second term converges to zero. Writing

$$\hat{\beta} = \mathbf{b}_o + \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i \right),$$

we have from the WLLN that $n^{-1} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i \xrightarrow{\mathbb{P}} \mathbb{E}(\mathbf{x}_i \varepsilon_i) = \mathbf{0}$, and

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \xrightarrow{\mathbb{P}} \mathbb{E}(\mathbf{x}_i \mathbf{x}_i') =: \mathbf{M}_{xx}.$$

It follows that

$$\hat{\beta} \xrightarrow{\mathbb{P}} \mathbf{b}_o + \mathbf{M}_{xx}^{-1} \mathbf{0} = \mathbf{b}_o.$$

OLS Consistency

Given the specification $y_i = \mathbf{x}'_i \beta + u_i$, suppose that Modern Assumption 1 holds. Then, $\hat{\beta} \xrightarrow{\mathbb{P}} \mathbf{b}_o$; that is, $\hat{\beta}$ is weakly consistent for \mathbf{b}_o .

- Aside from the WLLN effect, consistency hinges on the condition $\mathbb{E}(\mathbf{x}_i \varepsilon_i) = \mathbf{0}$, i.e., ε_i and \mathbf{x}_i are uncorrelated. Thus, the linear specification suffices to capture the systematic part of y_i ; no other information could be of help in a linear fashion.
- Consistency ensures that, when more information (a larger sample) become available, it is more likely that the OLS estimates will be close to the true parameter values. Note that consistency is a property that depends on sample size n , but unbiasedness is not.

Consistency and Finite Samples

Q1: Why is consistency relevant if we only have finite samples in practice?

Ans: It is true that we could never have an infinite sample, but consistency ensures that an estimator can approximate the true parameters **arbitrarily well** when enough information are available.

Q2: For what sample size should we expect estimates to well approximate the true parameters?

Ans: It depends on the stochastic properties of the data. If the data are independent, it typically requires a larger sample to achieve the same degree of accuracy when the data distributions have fatter tails than does the normal distribution. If the data are dependent (correlated), an even larger sample may be needed. We must emphasize that there is **no** “magic number” (such as 30, 50 or 100) for the desired sample size.

Exclusion of Important Variables

Given the specification $y_i = \mathbf{x}_i' \boldsymbol{\beta} + u_i$, suppose y_i actually depend on more variables \mathbf{z}_i such that $y_i = \mathbf{x}_i' \mathbf{b}_o + \mathbf{z}_i' \mathbf{d}_o + \varepsilon_i$, where ε_i satisfy Modern Assumption I(ii) with $\mathbb{E}(\mathbf{x}_i \varepsilon_i) = \mathbf{0}$ and $\mathbb{E}(\mathbf{z}_i \varepsilon_i) = \mathbf{0}$. That is, our specification excludes the important variables \mathbf{z} that should have been included in the model. In this case,

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i y_i \right) \\ &= \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i (\mathbf{x}_i' \mathbf{b}_o + \mathbf{z}_i' \mathbf{d}_o + \varepsilon_i) \right) \\ &= \mathbf{b}_o + \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left[\left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{z}_i' \right) \mathbf{d}_o + \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i \right) \right].\end{aligned}$$

By the WLLN,

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{z}'_i \xrightarrow{\mathbb{P}} \mathbb{E}(\mathbf{x}_i \mathbf{z}'_i) =: \mathbf{M}_{\mathbf{xz}}, \quad \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i \xrightarrow{\mathbb{P}} \mathbb{E}(\mathbf{x}_i \varepsilon_i) = \mathbf{0}.$$

It follows that $\hat{\boldsymbol{\beta}}$ is **inconsistent** for \mathbf{b}_o :

$$\hat{\boldsymbol{\beta}} \xrightarrow{\mathbb{P}} \mathbf{b}_o + \mathbf{M}_{\mathbf{xx}}^{-1} \mathbf{M}_{\mathbf{xz}} \mathbf{d}_o + \mathbf{M}_{\mathbf{xx}}^{-1} \cdot \mathbf{0} = \mathbf{b}_o + \mathbf{M}_{\mathbf{xx}}^{-1} \mathbf{M}_{\mathbf{xz}} \mathbf{d}_o.$$

Note that consistency would still obtain if the vector of the omitted variables \mathbf{z} is **orthogonal** to \mathbf{x} , i.e., $\mathbf{M}_{\mathbf{xz}} = \mathbb{E}(\mathbf{x}_i \mathbf{z}'_i) = \mathbf{0}$. This shows that inconsistency arises because the included variables in \mathbf{x} are correlated with the omitted variables in \mathbf{z} . Given the specification $y_i = \mathbf{x}'_i \boldsymbol{\beta} + u_i$, the correlation between \mathbf{x} and \mathbf{z} can also be understood as $\mathbb{E}(\mathbf{x}_i u_i) \neq \mathbf{0}$, because u_i now includes the omitted variables \mathbf{z}_i .

Recall that the dummy variable model for estimating the average treatment effect on the outcome variable Y is:

$$y_i = \alpha + \mathbf{x}_i' \boldsymbol{\beta} + D_i \gamma + u_i,$$

where D is the treatment indicator. D is **endogenous** when D and the outcome Y depend on some unobserved factors. For example, in the study of “return to college education”, D is the dummy for college attendance, and Y is the earning. Here, D and Y may depend on the factors that characterize “ability”, such that people with better “ability” are more likely to go to college and earn a higher pay. Thus, y_i actually follow:

$$y_i = \alpha_o + \mathbf{x}_i' \boldsymbol{\beta}_o + D_i \gamma_o + (A_i \delta_o + \varepsilon_i),$$

where A_i denote the unobserved “ability” index. The regression error u_i now includes $A_i \delta_o$ and ε_i , and $\mathbb{E}(D_i u_i) \neq \mathbf{0}$ because D_i also depend on A_i . As a result, the OLS estimators are inconsistent.

Inclusion of Irrelevant Variables

Suppose we specify multiple linear regression: $y_i = \mathbf{x}_i' \boldsymbol{\beta} + \mathbf{z}_i' \boldsymbol{\delta} + u_i$, but y_i depend only on \mathbf{x}_i : $y_i = \mathbf{x}_i' \mathbf{b}_o + \varepsilon_i$, where $\mathbb{E}(\mathbf{x}_i \varepsilon_i) = \mathbf{0}$. That is, we include the irrelevant variables \mathbf{z} in our specification. As discussed before, y_i can be written as

$$y_i = \mathbf{x}_i' \mathbf{b}_o + \mathbf{z}_i' \cdot \mathbf{0} + \varepsilon_i,$$

with $\mathbf{0}$ as the true parameter associated with \mathbf{z}_i , and u_i agree with ε_i . Then, as long as $\mathbb{E}(\mathbf{z}_i \varepsilon_i) = 0$, $\hat{\boldsymbol{\beta}}$ is consistent for \mathbf{b}_o , and $\hat{\boldsymbol{\delta}}$ is consistent for $\mathbf{0}$. This shows that a more general model specification does not affect consistency.

Asymptotic Normality

In addition to Modern Assumption I(i) and (ii), we also postulate the following condition.

Modern Assumption I(iii)

The random variables y_i are $y_i = \mathbf{x}_i' \mathbf{b}_o + \varepsilon_i$, $i = 1, \dots, n$, such that

(iii) $\mathbf{x}_i \varepsilon_i$ obey a multivariate CLT:

$$\mathbf{V}^{-1/2} \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i \xrightarrow{D} \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

where for $\mathbf{V} = \text{var}(\mathbf{x}_i \varepsilon_i)$ is a positive definite matrix.

Recall

$$\hat{\beta} - \mathbf{b}_o = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i \right).$$

We can then write

$$\sqrt{n}(\hat{\beta} - \mathbf{b}_o) = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \underbrace{\mathbf{V}^{1/2} \mathbf{V}^{-1/2}}_{=I} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i \varepsilon_i \right).$$

Modern Assumption I(iii) and the fact that $n^{-1} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \xrightarrow{\mathbb{P}} \mathbf{M}_{xx}$ lead to:

$$\sqrt{n}(\hat{\beta} - \mathbf{b}_o) \xrightarrow{D} \mathbf{M}_{xx}^{-1} \mathbf{V}^{1/2} \mathcal{N}(\mathbf{0}, I),$$

where the limit is simply a linear transformation of $\mathcal{N}(\mathbf{0}, I)$.

Note that $\mathbf{x}_i \mathbf{x}_i'$ is symmetric, and so is $\mathbb{E}(\mathbf{x}_i \mathbf{x}_i') = \mathbf{M}_{xx}$. Also, \mathbf{V} is a covariance matrix and must also be symmetric. It follows that the limit is a multivariate normal distribution with mean $\mathbf{0}$ and covariance matrix:

$$(\mathbf{M}_{xx}^{-1} \mathbf{V}^{1/2})' (\mathbf{M}_{xx}^{-1} \mathbf{V}^{1/2}) = \mathbf{M}_{xx}^{-1} \mathbf{V}^{1/2} \mathbf{V}^{1/2} \mathbf{M}_{xx}^{-1} = \mathbf{M}_{xx}^{-1} \mathbf{V} \mathbf{M}_{xx}^{-1}.$$

We have established the following result.

OLS Asymptotic Normality

Given the specification $y_i = \mathbf{x}_i' \boldsymbol{\beta} + u_i$, suppose that Modern Assumption I(i), (ii) and (iii) hold. Then,

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \mathbf{b}_o) \xrightarrow{D} \mathcal{N}(\mathbf{0}, \mathbf{M}_{xx}^{-1} \mathbf{V} \mathbf{M}_{xx}^{-1}),$$

where $\mathbf{M}_{xx}^{-1} \mathbf{V} \mathbf{M}_{xx}^{-1}$ is known as the **asymptotic covariance matrix**.

When ε_i are **conditionally homoskedastic** such that $\mathbb{E}(\varepsilon_i^2 | \mathbf{x}_i) = \sigma_o^2$,

$$\mathbf{V} = \mathbb{E}(\varepsilon_i^2 \mathbf{x}_i \mathbf{x}_i') = \mathbb{E}[\mathbb{E}(\varepsilon_i^2 | \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i'] = \sigma_o^2 \mathbb{E}(\mathbf{x}_i \mathbf{x}_i'),$$

where the 2nd equality follows from the **law of iterated expectations**. In this case, the asymptotic covariance matrix simplifies to

$$\mathbf{M}_{xx}^{-1} \mathbf{V} \mathbf{M}_{xx}^{-1} = \mathbf{M}_{xx}^{-1} (\sigma_o^2 \mathbf{M}_{xx}) \mathbf{M}_{xx}^{-1} = \sigma_o^2 \mathbf{M}_{xx}^{-1},$$

and

$$\sqrt{n}(\hat{\beta} - \mathbf{b}_o) \xrightarrow{D} \mathcal{N}(\mathbf{0}, \sigma_o^2 \mathbf{M}_{xx}^{-1}).$$

When $\mathbb{E}(\varepsilon_i^2 | \mathbf{x}_i)$ depend on \mathbf{x}_i , i.e., ε_i are **conditionally heteroskedastic**, the asymptotic covariance matrix has the *sandwich* form and cannot be simplified.

Some Remarks

- OLS consistency and asymptotic normality would hold provided that data satisfy proper WLLN and CLT. We do **not** require \mathbf{x} to be non-random, **nor** do we impose normality and (conditional) homoskedasticity on y_i .
- Under these weaker conditions, we cannot derive the **exact** distribution of $\hat{\beta}$, but are able to derive the limiting distribution of $\sqrt{n}(\hat{\beta} - \mathbf{b}_o)$. The limiting normal distribution serves as an **approximation** to the unknown, exact distribution.
- For a given sample size, whether the approximation to normality is good depends on the stochastic properties of the data. A larger sample typically yields better approximation to normality, but, again, there is **no** “magic number” for the sample size (such as 30, 50 or 100) that assures good approximation.

Estimation of Asymptotic Covariance Matrix

Let $\mathbf{D}_o := \mathbf{M}_{xx}^{-1} \mathbf{V} \mathbf{M}_{xx}^{-1}$. Then, $\mathbf{D}_o^{-1/2} \sqrt{n}(\hat{\boldsymbol{\beta}} - \mathbf{b}_o) \xrightarrow{D} \mathcal{N}(\mathbf{0}, \mathbf{I})$. By the WLLN, $\mathbf{M}_{xx} = \mathbb{E}(\mathbf{x}_i \mathbf{x}_i')$ can be consistently estimated by its sample counterpart:

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i',$$

and $\mathbf{V} = \mathbb{E}(\varepsilon_i^2 \mathbf{x}_i \mathbf{x}_i')$ can be consistently estimated by (proof omitted):

$$\hat{\mathbf{V}} = \frac{1}{n} \sum_{i=1}^n \hat{u}_i^2 \mathbf{x}_i \mathbf{x}_i',$$

where \hat{u}_i are the OLS residuals. This $\hat{\mathbf{V}}$ is known as a **heteroskedasticity-consistent** covariance matrix estimator because \mathbf{V} permits conditional heteroskedasticity of an unknown form.

A weakly consistent estimator of \mathbf{D}_o is thus

$$\hat{\mathbf{D}} = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \hat{u}_i^2 \mathbf{x}_i \mathbf{x}_i' \right) \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1};$$

this is known as the **Eicker-White** estimator. Replacing \mathbf{D}_o by its consistent estimator $\hat{\mathbf{D}}$ we obtain:

$$\hat{\mathbf{D}}^{-1/2} \sqrt{n}(\hat{\beta} - \mathbf{b}_o) = \underbrace{\hat{\mathbf{D}}^{-1/2} \mathbf{D}_o^{1/2}}_{\xrightarrow{P} \mathbf{I}} \underbrace{\mathbf{D}_o^{-1/2} \sqrt{n}(\hat{\beta} - \mathbf{b}_o)}_{\xrightarrow{D} \mathcal{N}(\mathbf{0}, \mathbf{I})} \xrightarrow{D} \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

In practice, we use the following covariance matrix estimator:

$$\tilde{\mathbf{D}} = \hat{\mathbf{D}}/n = \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\sum_{i=1}^n \hat{u}_i^2 \mathbf{x}_i \mathbf{x}_i' \right) \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1},$$

such that

$$\tilde{\mathbf{D}}^{-1/2} \sqrt{n}(\hat{\beta} - \mathbf{b}_o) = \hat{\mathbf{D}}^{-1/2} \sqrt{n}(\hat{\beta} - \mathbf{b}_o) \xrightarrow{D} \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

Let \tilde{d}_{jj} denote the j^{th} diagonal element of $\tilde{\mathbf{D}}$; its square root, $\sqrt{\tilde{d}_{jj}}$, is also referred to as the **Eicker-White standard error** for $\hat{\beta}_{j-1}$.

Special case: When there is conditional homoskedasticity: $\mathbb{E}(\varepsilon_i^2 | \mathbf{x}_i) = \sigma_o^2$, $\mathbf{D}_o = \sigma_o^2 \mathbf{M}_{xx}^{-1}$ can be consistently estimated by

$$\hat{\mathbf{D}} = \hat{\sigma}^2 \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1}, \quad \text{with} \quad \hat{\sigma}^2 = \frac{1}{n - k - 1} \sum_{i=1}^n \hat{u}_i^2.$$

In this case, $\tilde{\mathbf{D}} = \hat{\mathbf{D}}/n = \hat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1}$, which is exactly the estimator $\widehat{\text{var}(\hat{\beta})}$ obtained earlier under Classical Assumption.

Remark: It is important to note that the estimator $\hat{\sigma}^2 (n^{-1} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i')^{-1}$ would **not** be consistent for \mathbf{D}_o when conditional heteroskedasticity is present. Consequently, we would lose asymptotic normality if $\hat{\beta}$ is normalized using this inconsistent covariance matrix estimator.

Testing A Single Hypothesis

Suppose we are testing only one parameter: $\mathbf{R}\mathbf{b}_o = b_2 = c$, with

$$\mathbf{R} = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 & \cdots & 0 \end{pmatrix}.$$

Recall that by asymptotic normality,

$$\tilde{\mathbf{D}}^{-1/2}(\hat{\boldsymbol{\beta}} - \mathbf{b}_o) \xrightarrow{D} \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

where $\tilde{\mathbf{D}} = (\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i')^{-1} (\sum_{i=1}^n \hat{u}_i^2 \mathbf{x}_i \mathbf{x}_i') (\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i')^{-1}$. It is easy to see that

$$(\mathbf{R}\tilde{\mathbf{D}}\mathbf{R}')^{-1/2}[\mathbf{R}(\hat{\boldsymbol{\beta}} - \mathbf{b}_o)] \xrightarrow{D} \mathcal{N}(0, 1).$$

In this case, $\mathbf{R}\tilde{\mathbf{D}}\mathbf{R}' = \tilde{d}_{33}$, the third diagonal element of $\tilde{\mathbf{D}}$.

Then under the null hypothesis, the statistic

$$(\mathbf{R}\tilde{\mathbf{D}}\mathbf{R}')^{-1/2}(\mathbf{R}\hat{\boldsymbol{\beta}} - c) = \frac{\hat{\beta}_2 - c}{\text{EW-se}(\hat{\beta}_2)} \xrightarrow{D} \mathcal{N}(0, 1),$$

where $\text{EW-se}(\hat{\beta}_2) = \sqrt{\tilde{d}_{33}}$ is the Eicker-White standard error of $\hat{\beta}_2$. Note that this is the t statistic standardized by the Eicker-White standard error and has the limiting distribution $\mathcal{N}(0, 1)$.

Suppose the hypothesis involves two parameters: $\mathbf{R}\mathbf{b}_o = b_2 - b_3 = c$, with

$$\mathbf{R} = \begin{pmatrix} 0 & 0 & 1 & -1 & 0 & \cdots & 0 \end{pmatrix}.$$

In this case, one can verify $\mathbf{R}\tilde{\mathbf{D}}\mathbf{R}' = \tilde{d}_{33} + \tilde{d}_{44} - 2\tilde{d}_{34}$, where \tilde{d}_{ij} the (i, j) th term of $\tilde{\mathbf{D}}$. Then,

$$(\mathbf{R}\tilde{\mathbf{D}}\mathbf{R}')^{-1/2}(\mathbf{R}\hat{\boldsymbol{\beta}} - c) = \frac{\hat{\beta}_2 - \hat{\beta}_3 - c}{\text{EW-se}(\hat{\beta}_2 - \hat{\beta}_3)} \xrightarrow{D} \mathcal{N}(0, 1),$$

where $\text{EW-se}(\hat{\beta}_2 - \hat{\beta}_3) = \sqrt{\tilde{d}_{33} + \tilde{d}_{44} - 2\tilde{d}_{34}}$.

Limiting Distribution of the t Statistic

Under the null hypothesis $\mathbf{R}\mathbf{b}_o = c$, where \mathbf{R} is $1 \times (k + 1)$, we have the t statistic:

$$\frac{\mathbf{R}\hat{\beta} - c}{\text{EW-se}(\mathbf{R}\hat{\beta})} \xrightarrow{D} \mathcal{N}(0, 1),$$

where $\text{EW-se}(\mathbf{R}\hat{\beta}) = (\mathbf{R}\tilde{\mathbf{D}}\mathbf{R}')^{1/2}$ denotes the Eicker-White standard error.

Remark: As $\hat{\mathbf{D}}$ is consistent for \mathbf{D}_o when conditional heteroskedasticity is present, the t statistic based on the standard errors from $\tilde{\mathbf{D}}$ is said to be **robust to conditional heteroskedasticity**. Such t statistic uses the robust, Eicker-White standard error and is valid even when conditional heteroskedasticity is present with an unknown form. This is the statistic we should use in empirical studies with cross-section data.

A Special Case

When there is conditional homoskedasticity: $\mathbb{E}(\varepsilon_i^2 | \mathbf{x}_i) = \sigma_o^2$,
 $\mathbf{D}_o = \sigma_o^2 \mathbf{M}_{xx}^{-1}$. In this case, the Eicker-White covariance matrix estimator remains consistent but is not needed. Instead, we may use $\tilde{\mathbf{D}} = \hat{\sigma}^2 (\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i')^{-1}$, so that the t statistic is:

$$\frac{\mathbf{R}\hat{\boldsymbol{\beta}} - c}{\text{se}(\mathbf{R}\hat{\boldsymbol{\beta}})} = \frac{\mathbf{R}\hat{\boldsymbol{\beta}} - c}{\hat{\sigma} \sqrt{\mathbf{R}(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i')^{-1} \mathbf{R}'}} \xrightarrow{D} \mathcal{N}(0, 1),$$

where $\text{se}(\mathbf{R}\hat{\boldsymbol{\beta}})$ is the conventional OLS standard error of $\mathbf{R}\hat{\boldsymbol{\beta}}$.

Remark: This t statistic is valid only when conditional homoskedasticity is present. This result is analogous to that under Classical Assumption.

Testing Multiple Hypotheses

Consider now the joint hypothesis $\mathbf{R}\mathbf{b}_o = \mathbf{c}$, where \mathbf{R} is $q \times (k + 1)$. Then under the null hypothesis,

$$(\mathbf{R}\tilde{\mathbf{D}}\mathbf{R}')^{-1/2}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}) \xrightarrow{D} \mathcal{N}(\mathbf{0}, \mathbf{I}_q),$$

so that the limit contains q independent $\mathcal{N}(0, 1)$ variables. The **Wald statistic** is the inner product of the left-hand side:

$$\begin{aligned}\mathcal{W} &= (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c})'(\mathbf{R}\tilde{\mathbf{D}}\mathbf{R}')^{-1/2}(\mathbf{R}\tilde{\mathbf{D}}\mathbf{R}')^{-1/2}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}) \\ &= (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c})'(\mathbf{R}\tilde{\mathbf{D}}\mathbf{R}')^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}),\end{aligned}$$

and its limiting distribution is therefore the sum of q independent, squared $\mathcal{N}(0, 1)$ variables, i.e., the $\chi^2(q)$ distribution.

We have established the following result.

Limiting Distribution of the Wald Statistic

Under the null hypothesis $\mathbf{R}\mathbf{b}_0 = \mathbf{c}$, where \mathbf{R} is $q \times (k + 1)$, we have the Wald statistic:

$$(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c})'(\mathbf{R}\tilde{\mathbf{D}}\mathbf{R}')^{-1}(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}) \xrightarrow{D} \chi^2(q),$$

where $\tilde{\mathbf{D}} = (\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i')^{-1} (\sum_{i=1}^n \hat{u}_i^2 \mathbf{x}_i \mathbf{x}_i') (\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i')^{-1}$.

Remark: The Wald statistic is **robust to conditional heteroskedasticity** when the Eicker-White covariance matrix estimator $\tilde{\mathbf{D}}$ is used in the statistic.

Example: Consider the hypothesis of 3 parameters being zero:

$$\mathbf{R} = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 1 & 0 & \cdots & 0 \end{pmatrix}, \quad \mathbf{R}\mathbf{b}_o = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix},$$

The Wald statistic is then

$$\mathcal{W} = \begin{pmatrix} \hat{\beta}_1 & \hat{\beta}_2 & \hat{\beta}_3 \end{pmatrix} (\mathbf{R}\tilde{\mathbf{D}}\mathbf{R}')^{-1} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \end{pmatrix} \xrightarrow{D} \chi^2(3).$$

You can easily check $\mathbf{R}\tilde{\mathbf{D}}\mathbf{R}'$ is a 3×3 matrix but not a diagonal matrix in this case; write down this matrix using the notations d_{ij} (homework).

A Special Case

When there is conditional homoskedasticity: $\mathbb{E}(\varepsilon^2|\mathbf{x}_i) = \sigma_o^2$, $\mathbf{D}_o = \sigma_o^2 \mathbf{M}_{xx}^{-1}$. In this case, $\tilde{\mathbf{D}} = \hat{\sigma}^2 (\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i')^{-1}$, so that

$$\mathcal{W} = (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c})' \left[\mathbf{R} \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \mathbf{R}' \right]^{-1} \frac{\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{c}}{\hat{\sigma}^2} \xrightarrow{D} \chi^2(q).$$

- This Wald statistic is **not robust** to conditional heteroskedasticity, because $\hat{\sigma}^2 (\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i')^{-1}$ is not a consistent estimator when conditional heteroskedasticity is present.
- It can be seen that \mathcal{W}/q is nothing but the conventional F statistic. Thus, F statistics are **not robust** to conditional heteroskedasticity. In practice, we may replace the F tests in Lecture 3 with the corresponding, robust Wald tests.

Example: Wage Regressions

The estimated wage model based on Taiwan's 2010 male data (11561 obs):

	3.8939	+ 0.0800 educ	+ 0.0166 exper	
OLS-se	(0.0198)	(0.0012)	(0.0003)	
<i>t</i>	(197.05)	(65.41)	(50.45)	
EW-se	(0.0208)	(0.0013)	(0.0004)	
<i>t</i>	(186.94)	(61.63)	(43.71)	
	$\bar{R}^2 = 0.2893$	$\hat{\sigma} = 0.3595$	Reg $F = 1992.4$	
	3.7902	+ 0.0779 educ	+ 0.0365 exper	− 0.0005 exper ²
OLS-se	(0.0199)	(0.0012)	(0.0009)	(0.00002)
<i>t</i>	(190.60)	(64.77)	(38.72)	(−22.47)
EW-se	(0.0207)	(0.0013)	(0.0010)	(0.00003)
<i>t</i>	(183.53)	(61.04)	(35.05)	(−18.63)
	$\bar{R}^2 = 0.319$	$\hat{\sigma} = 0.3519$	Reg $F = 1600.2$	

Regression F statistics are based on White's standard errors.

Tests of Conditional Homoskedasticity

When there is **no** conditional heteroskedasticity, the asymptotic covariance matrix of the OLS estimator has a simpler form and can be consistently estimated by the conventional OLS covariance matrix estimator. It is thus desirable to first check if the data exhibit conditional homoskedasticity. Rejecting this hypothesis suggests that the robust, Eicker-White standard errors should be employed.

The well known **Breusch-Pagan (BP) test** is to test the null hypothesis of conditional homoskedasticity: $\mathbb{E}(\varepsilon_i^2 | \mathbf{x}_i, \mathbf{w}_i) = \sigma_o^2$, where the variables in \mathbf{w}_i are from the information set but different from those of \mathbf{x}_i . The BP test is designed to test against the alternative that $\mathbb{E}(\varepsilon_i^2 | \mathbf{x}_i, \mathbf{w}_i)$ is an **unknown function of \mathbf{z}_i** and hence changes with i , where \mathbf{z}_i contains the elements of \mathbf{x}_i and \mathbf{w}_i that may affect the conditional variance of ε_i .

Breusch-Pagan Test

The Breusch-Pagan statistic is computed as follows:

- 1 Regress y_i on \mathbf{x}_i to obtain the OLS residuals \hat{u}_i .
- 2 Regress \hat{u}_i^2 on constant one and \mathbf{z}_i and obtain R_{aux}^2 .
- 3 The BP statistic is $n R_{\text{aux}}^2 \xrightarrow{D} \chi^2(m)$, where m is the number of elements in \mathbf{z}_i .

Remark:

- Intuitively, the residuals \hat{u}_i approximate the true errors. When there is homoskedasticity, we expect that \hat{u}_i^2 do not depend on \mathbf{z}_i , so that the auxiliary regression in Step 2 above would have a small R_{aux}^2 (even after multiplied by n). Otherwise, R_{aux}^2 would be large, suggesting the null hypothesis is false.
- A drawback of the BP test is that the choice of \mathbf{z}_i is **subjective**.

White Test

As for the choice of \mathbf{z}_i , the **White test** of White (1980) proposes to choose all non-constant regressors of \mathbf{x}_i , and their respective squares and pairwise products. The resulting test is also known as a test of conditional heteroskedasticity of unknown form. For example, when \mathbf{x}_i contains the constant one and 3 variables x_{i1}, x_{i2}, x_{i3} ,

$$\mathbf{z}_i = (x_{i1}, x_{i2}, x_{i3}, x_{i1}^2, x_{i2}^2, x_{i3}^2, x_{i1}x_{i2}, x_{i1}x_{i3}, x_{i2}x_{i3})'.$$

Then $n R_{\text{aux}}^2$ is distributed as $\chi^2(9)$ in the limit. A difficulty of the White test is that the number of variables in \mathbf{z}_i may be too large; this may happen even when the original regression contains only a moderate number of regressors. Note that the White test was originally derived as an **information matrix test**.

Endogeneity

In linear regressions, a regressor is said to be **endogenous** if it is correlated with the regression error.

- Treatment effect model: The treatment indicator and outcome variable may both depend on some unobserved factors and hence are correlated. In this case, the OLS estimator is inconsistent
- **Measurement error problem:** For $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$, the regressor x_i can only be observed with random error: $w_i = x_i + v_i$. The operational model is:

$$y_i = \beta_0 + \beta_1 w_i + (\varepsilon_i - \beta_1 v_i),$$

with $\varepsilon_i - \beta_1 v_i$ the regression error. Clearly, w_i must be correlated with the regression error. Here, the random error v_i plays the same role as the unobserved factor in the case above.

- **Simultaneous equation system:** The following examples are taken from Chapter 3 of Hayashi (2000).

- 1 Consider the following demand-supply system:

$$q_i = \alpha_0 + \alpha_1 p_i + u_i, \quad (\text{demand})$$

$$q_i = \beta_0 + \beta_1 p_i + v_i. \quad (\text{supply})$$

Solving this system for p_i and q_i , we find p_i depend on both u_i and v_i and hence are endogenous in both equations.

- 2 A simple economic system with the GNP identity: $y_i = c_i + g_i$, where c is consumption and g is investment, and the consumption function:

$$c_i = \alpha_0 + \alpha_1 y_i + u_i, \quad 0 < \alpha_1 < 1.$$

It is easily calculated that

$$y_i = \frac{\alpha_0}{1 - \alpha_1} + \frac{1}{1 - \alpha_1} g_i + \frac{1}{1 - \alpha_1} u_i,$$

so that y_i are correlated with u_i .

Instrumental Variable Estimation

Suppose that $y_i = \mathbf{x}_i' \beta_o + u_i$, where \mathbf{x}_i is $(k + 1) \times 1$ and contains some endogenous regressors such that $\mathbb{E}(\mathbf{x}_i u_i) \neq \mathbf{0}$. Let \mathbf{z}_i be the $(k + 1) \times 1$ vector that is correlated with \mathbf{x}_i but uncorrelated with u_i , i.e., $\mathbb{E}(\mathbf{z}_i u_i) = \mathbf{0}$. The variables in \mathbf{z}_i are referred to as the **instrumental variables** (IV) or simply the **instruments**. Instead of solving the normal equations $\mathbf{X}'(\mathbf{y} - \mathbf{X}\beta) = \mathbf{0}$ for β , the IV method solves:

$$\mathbf{Z}'(\mathbf{y} - \mathbf{X}\beta) = \mathbf{0}.$$

This system is **just identified** in the sense that it contains $k + 1$ equations with $k + 1$ unknowns. Given that $\mathbf{Z}'\mathbf{X}$ is of full rank (so that it is invertible), the solution is the **IV estimator**:

$$\hat{\beta}_{IV} = (\mathbf{Z}'\mathbf{X})^{-1} \mathbf{Z}'\mathbf{y} = \left(\sum_{i=1}^n \mathbf{z}_i \mathbf{x}_i' \right)^{-1} \left(\sum_{i=1}^n \mathbf{z}_i y_i \right).$$

It is easy to see that the IV estimator is consistent for β_o :

$$\begin{aligned}\hat{\beta}_{IV} &= \left(\sum_{i=1}^n \mathbf{z}_i \mathbf{x}'_i \right)^{-1} \left(\sum_{i=1}^n \mathbf{z}_i (\mathbf{x}'_i \beta_o + u_i) \right) \\ &= \beta_o + \left(\frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{x}'_i \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{z}_i u_i \right) \\ &\xrightarrow{\mathbf{P}} \beta_o + [\mathbb{E}(\mathbf{z}_i \mathbf{x}'_i)]^{-1} \mathbb{E}(\mathbf{z}_i u_i) = \beta_o,\end{aligned}$$

because $\mathbb{E}(\mathbf{z}_i u_i) = \mathbf{0}$. As for the OLS asymptotics, we assume the CLT:

$$n^{-1/2} \sum_{i=1}^n \mathbf{z}_i u_i \xrightarrow{D} \mathcal{N}(\mathbf{0}, \mathbf{V}), \quad \mathbf{V} = \mathbb{E}(u_i^2 \mathbf{z}_i \mathbf{z}'_i).$$

The asymptotic distribution result for $\hat{\beta}_{IV}$ is immediate:

$$\sqrt{n}(\hat{\beta}_{IV} - \beta_o) \xrightarrow{D} \mathcal{N}(\mathbf{0}, \mathbf{M}_{ZX}^{-1} \mathbf{V} \mathbf{M}_{XZ}^{-1}),$$

where $\mathbf{M}_{ZX} := \mathbb{E}(\mathbf{z}_i \mathbf{x}'_i) = \mathbf{M}'_{XZ}$.

The asymptotic covariance matrix $\mathbf{D}_o = \mathbf{M}_{zx}^{-1} \mathbf{V} \mathbf{M}_{xz}^{-1}$ can be estimated using the Eicker-White-type estimator:

$$\hat{\mathbf{D}} = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{x}_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \hat{u}_i^2 \mathbf{z}_i \mathbf{z}_i' \right) \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{z}_i' \right)^{-1},$$

where $\hat{u}_i = y_i - \mathbf{x}_i' \hat{\beta}_{IV}$ are the IV residuals. As before,

$$\tilde{\mathbf{D}} = \hat{\mathbf{D}}/n = \left(\sum_{i=1}^n \mathbf{z}_i \mathbf{x}_i' \right)^{-1} \left(\sum_{i=1}^n \hat{u}_i^2 \mathbf{z}_i \mathbf{z}_i' \right) \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{z}_i' \right)^{-1},$$

so that

$$\tilde{\mathbf{D}}^{-1/2} (\hat{\beta}_{IV} - \beta_o) = \hat{\mathbf{D}}^{-1/2} \sqrt{n} (\hat{\beta}_{IV} - \beta_o) \xrightarrow{D} \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

The Eicker-White standard errors for $\hat{\beta}_{IV}$ are taken from $\tilde{\mathbf{D}}$ and robust to potential conditional heteroskedasticity. The Wald tests can be conducted as in the OLS context.

Remarks:

- In practice, \mathbf{x}_i may contain only one endogenous variable. For example, write $\mathbf{x}_i = (\boldsymbol{\xi}'_i \ x_{i1})'$, where $\boldsymbol{\xi}_i = (1 \ x_{i2} \ \dots \ x_{ik})'$ contains all exogenous regressors such that $\mathbb{E}(\boldsymbol{\xi}_i u_i) = \mathbf{0}$, and x_{i1} is endogenous such that $\mathbb{E}(x_{i1} u_i) \neq 0$. We can set $\mathbf{z}_i = (\boldsymbol{\xi}'_i \ \zeta_i)$, where ζ_i is an IV with $\mathbb{E}(\zeta_i u_i) = 0$. That is, only one IV is needed for this endogenous variable. Finding a proper IV is practically challenging, however.
- The condition that $\mathbf{Z}'\mathbf{X}$ is of full rank requires that the variables in \mathbf{Z} must be correlated with the variables in \mathbf{X} . In other words, we cannot use an arbitrary but irrelevant variable as the instrument.
- There may exist multiple instruments for a single endogenous variable. In that case, \mathbf{Z} is $n \times \ell$ with $\ell > k + 1$, and the system $\mathbf{Z}'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{0}$ is **over-identified**, and the IV estimator cannot be computed. We will discuss this estimation problem later.

An Alternative Expression

Consider the simple linear regression: $y_i = \beta_0 + \beta_1 x_i + u_i$, with x_i endogenous. It is readily verified that the IV estimator for β_1 is (check):

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

The IV estimator can be expressed as the **ratio** of two regression coefficients:

$$\hat{\beta}_{1,IV} = \frac{\sum_{i=1}^n (z_i - \bar{z})(y_i - \bar{y}) / \sum_{i=1}^n (z_i - \bar{z})^2}{\sum_{i=1}^n (z_i - \bar{z})(x_i - \bar{x})' / \sum_{i=1}^n (z_i - \bar{z})^2},$$

where the numerator is the slope of the regression of y on 1 and z , and the denominator is the slope of the regression of x on 1 and z .

To see why it is the case, write the regression for the denominator term of $\hat{\beta}_{1,IV}$ as:

$$x_i = \gamma_0 + \gamma_1 z_i + v_i.$$

Plugging this into the original model: $y_i = \beta_0 + \beta_1 x_i + u_i$, we obtain the so-called “reduced form” (model without endogenous regressor):

$$y_i = (\beta_0 + \beta_1 \gamma_0) + \beta_1 \gamma_1 z_i + (\beta_1 v_i + u_i).$$

which is the regression for the numerator term of $\hat{\beta}_{1,IV}$. While the numerator of the IV estimator provides an estimate of $\beta_1 \gamma_1$, the denominator yields an estimate of γ_1 . Their ratio is thus an estimate of β_1 , as it ought to be.

Two-Stage Least Squares Estimation

Again consider the case $y_i = \mathbf{x}_i' \boldsymbol{\beta}_o + u_i$, where $\mathbf{x}_i = (\boldsymbol{\xi}_i' x_{i1})'$ with x_{i1} endogenous and $\boldsymbol{\xi}_i$ exogenous, and $\boldsymbol{\beta}_o = (\mathbf{b}' \beta_1)'$. We may first regress x_{i1} on $\boldsymbol{\xi}_i$ and z_i to obtain the fitted values \hat{x}_{i1} ; here, z_i may contain multiple instruments. Plugging the fitted values into the original model we have

$$y_i = \mathbf{x}_i' \boldsymbol{\beta}_o + u_i = \boldsymbol{\xi}_i' \mathbf{b} + \beta_1 \hat{x}_{i1} + [\beta_1(x_{i1} - \hat{x}_{i1}) + u_i].$$

This suggests the second-stage regression of y_i on $\boldsymbol{\xi}_i$ and \hat{x}_{i1} ; note that $\boldsymbol{\xi}_i$ are uncorrelated with u_i , $\sum_{i=1}^n \boldsymbol{\xi}_i'(x_{i1} - \hat{x}_{i1}) = 0$, and $\sum_{i=1}^n \hat{x}_{i1}(x_{i1} - \hat{x}_{i1}) = 0$. The resulting estimators of \mathbf{b} and β_1 are known as the **two-stage least squares (2SLS)** estimators.

We now focus on the case that \mathbf{Z} is the $n \times \ell$ matrix of instruments with $\ell > k + 1$, i.e., the number of instruments is greater than the number of parameters. The system, $\mathbf{Z}'(\mathbf{y} - \mathbf{X}\beta) = \mathbf{0}$, is over-identified and cannot be solved. Instead, we solve the following just-identified system:

$$\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'(\mathbf{y} - \mathbf{X}\beta) = \mathbf{0},$$

and obtain the 2SLS estimator:

$$\hat{\beta}_{2SLS} = [\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X}]^{-1}[\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{y}].$$

Let $\mathbf{P}_Z = \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'$, an orthogonal projection matrix. The first-stage regression of \mathbf{X} on \mathbf{Z} leads to the matrix of fitted values:

$\hat{\mathbf{X}} = \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X} = \mathbf{P}_Z\mathbf{X}$. Then, the second-stage regression of \mathbf{y} on $\hat{\mathbf{X}}$ yields the estimator:

$$(\hat{\mathbf{X}}'\hat{\mathbf{X}})^{-1}\hat{\mathbf{X}}'\mathbf{y} = (\mathbf{X}'\mathbf{P}_Z\mathbf{X})^{-1}\mathbf{X}'\mathbf{P}_Z\mathbf{y},$$

which is exactly the 2SLS estimator $\hat{\beta}_{2SLS}$.

Asymptotic Properties of the 2SLS Estimator

Given $\mathbf{y} = \mathbf{X}\beta_o + \mathbf{u}$, we have

$$\begin{aligned}\hat{\beta}_{2SLS} &= [\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X}]^{-1}[\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{y}] \\ &= [\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X}]^{-1}[\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'(\mathbf{X}\beta_o + \mathbf{u})] \\ &= \beta_o + [(\mathbf{X}'\mathbf{Z}/n)(\mathbf{Z}'\mathbf{Z}/n)^{-1}(\mathbf{Z}'\mathbf{X}/n)]^{-1} \\ &\quad [(\mathbf{X}'\mathbf{Z}/n)(\mathbf{Z}'\mathbf{Z}/n)^{-1}(\mathbf{Z}'\mathbf{u}/n)].\end{aligned}$$

By suitable WLLN, we have seen that $\mathbf{X}'\mathbf{Z}/n \xrightarrow{\mathbf{P}} \mathbb{E}(\mathbf{x}_i\mathbf{z}'_i) =: \mathbf{M}_{xz}$, $\mathbf{Z}'\mathbf{X}/n \xrightarrow{\mathbf{P}} \mathbf{M}'_{xz} =: \mathbf{M}_{zx}$, $\mathbf{Z}'\mathbf{Z}/n \xrightarrow{\mathbf{P}} \mathbb{E}(\mathbf{z}_i\mathbf{z}'_i) =: \mathbf{M}_{zz}$, and

$$\mathbf{Z}'\mathbf{u}/n = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i u_i \xrightarrow{\mathbf{P}} \mathbb{E}(\mathbf{z}_i u_i) = \mathbf{0}.$$

It follows that $\hat{\beta}_{2SLS} \xrightarrow{\mathbf{P}} \beta_o$.

Assuming the CLT: $n^{-1/2} \sum_{i=1}^n \mathbf{z}_i u_i \xrightarrow{D} \mathcal{N}(\mathbf{0}, \mathbf{V})$, it is readily verified that (check):

$$\sqrt{n}(\hat{\beta}_{2SLS} - \beta_o) \xrightarrow{D} \mathcal{N}(\mathbf{0}, \mathbf{D}_o),$$

with $\mathbf{D}_o = (\mathbf{M}_{xz} \mathbf{M}_{zz}^{-1} \mathbf{M}_{zx})^{-1} (\mathbf{M}_{xz} \mathbf{M}_{zz}^{-1} \mathbf{V} \mathbf{M}_{zz}^{-1} \mathbf{M}_{zx}) (\mathbf{M}_{xz} \mathbf{M}_{zz}^{-1} \mathbf{M}_{zx})^{-1}$. The asymptotic covariance matrix can be estimated by its sample counterpart $\hat{\mathbf{D}}$. The robust standard errors for $\hat{\beta}_{2SLS}$ are taken from:

$$\begin{aligned} \tilde{\mathbf{D}} = \hat{\mathbf{D}}/n &= [\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X}]^{-1} \\ &\quad [\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\hat{\mathbf{V}}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X}][\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X}]^{-1}, \end{aligned}$$

where $\hat{\mathbf{V}} = \sum_{i=1}^n \hat{u}_i^2 \mathbf{z}_i \mathbf{z}_i'$, with $\hat{u}_i = y_i - \mathbf{x}_i' \hat{\beta}_{2SLS}$.

Remark: The asymptotic covariance matrix must be estimated directly from the formula above; the covariance matrix estimator from the second-stage regression is **not** a correct estimator for \mathbf{D}_o .