

KEYPOINT SUMMARY II

I. PRELIMINARIES

1. Biasness: $\hat{\theta}$ is unbiased if $\text{Bias}(\hat{\theta}) = \mathbb{E}(\hat{\theta}) - \theta = 0$.
2. Efficiency: $\hat{\theta}$ is efficient if $\text{Var}(\hat{\theta}) \leq \text{Var}(\tilde{\theta})$ for all unbiased estimator $\tilde{\theta}$.
3. Consistency: $\hat{\theta}$ is consistent if $\hat{\theta} \xrightarrow{P} \theta$.
4. Mean square error: $\text{MSE}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \theta)^2]$
5. Mean square convergence: $\hat{\theta} \xrightarrow{m.s.} \theta$ if $\text{MSE}(\hat{\theta}) \rightarrow 0$ as $n \rightarrow \infty$.
 - $\text{MSE}(\hat{\theta}) = \text{Var}(\hat{\theta}) + (\text{Bias}(\hat{\theta}, \theta))^2$
 - $\hat{\theta} \xrightarrow{m.s.} \theta$ iff $\text{Bias}(\hat{\theta}) \rightarrow 0$ and $\text{Var}(\hat{\theta}) \rightarrow 0$
6. Small o (Convergence in probability): $X_n = o_p(n^k)$ if $\forall \epsilon > 0, \lim_{n \rightarrow \infty} \mathbb{P}(|\frac{X_n}{n^k}| > \epsilon) = 0$
7. Big O (Stochastic boundedness): $X_n = O_p(n^k)$ if $\forall \epsilon > 0, \exists K > 0, N > 0$ s.t. $\forall n > N, \mathbb{P}(|\frac{X_n}{n^k}| > K) < \epsilon$

II. SIMPLE LINEAR REGRESSION

Let $y_i = \beta x_i + \epsilon_i$. How to find an estimator for β ?

Derive an estimator by solving

$$\min_{\beta} \sum_{i=1}^n (y_i - \beta x_i)^2$$

The first-order condition gives

$$\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} = \beta + \frac{\frac{1}{n} \sum_{i=1}^n x_i \epsilon_i}{\frac{1}{n} \sum_{i=1}^n x_i^2}$$

By Law of Large Numbers (LLN),

$$\begin{aligned} & \cdot \frac{1}{n} \sum_{i=1}^n x_i \epsilon_i \xrightarrow{P} \mathbb{E}(x_i \epsilon_i) \\ & \cdot \frac{1}{n} \sum_{i=1}^n x_i^2 \xrightarrow{P} \mathbb{E}(x_i^2) \end{aligned}$$

$\hat{\beta}$ is consistent if $\mathbb{E}(x_i \epsilon_i) = 0$ and $\mathbb{E}(x_i^2) \neq 0$.

To derive the *asymptotic normality*, consider

$$\sqrt{n}(\hat{\beta} - \beta) = \sqrt{n} \cdot \frac{\frac{1}{n} \sum_{i=1}^n x_i \epsilon_i}{\frac{1}{n} \sum_{i=1}^n x_i^2}$$

Assume $\mathbb{E}(x_i^2 \epsilon_i^2)$ is finite, by Central Limit Theorem,

$$\cdot \sqrt{n} \cdot \frac{1}{n} \sum_{i=1}^n x_i \epsilon_i \xrightarrow{d} \mathcal{N}(0, \text{Var}(x_i \epsilon_i))$$

Since $\frac{1}{n} \sum_{i=1}^n x_i^2 \xrightarrow{P} \mathbb{E}(x_i^2)$, by Slutsky's Theorem,

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} \mathcal{N}\left(0, \frac{\text{Var}(x_i \epsilon_i)}{[\mathbb{E}(x_i^2)]^2}\right)$$

III. MULTIPLE LINEAR REGRESSION

A. Notation

Suppose the model is specified by

$$y_i = \beta_0 + x_{i1}\beta_1 + \cdots + x_{ik}\beta_k + \epsilon_i$$

Let $\mathbf{x}_i, \boldsymbol{\beta}$ be $k \times 1$ vectors,

$$\mathbf{x}_i = \begin{bmatrix} 1 \\ x_{i1} \\ \vdots \\ x_{ik} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}$$

Then the model can be written as $y_i = \mathbf{x}_i' \boldsymbol{\beta} + \epsilon_i$.

Let $\mathbf{X}, \mathbf{Y}, \boldsymbol{\epsilon}$ be matrices containing all observations,

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1' \\ \vdots \\ \mathbf{x}_k' \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nk} \end{bmatrix}, \quad \mathbf{Y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Therefore, $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$.

B. Assumptions

(A1) Linearity: $y_i = \mathbf{x}_i' \boldsymbol{\beta} + \epsilon_i = x_{i1}\beta_1 + \cdots + x_{ik}\beta_k + \epsilon_i$

(A2) Full rank: $\mathbb{E}(\mathbf{x}_i \mathbf{x}_i')$ is nonsingular.

(A3) Exogeneity: $\mathbb{E}(\epsilon_i | \mathbf{x}_j) = 0$, for $i, j = 1, \dots, n$

(A4) Homoskedasticity and nonautocorrelation:
 $\text{Var}(\boldsymbol{\epsilon} | \mathbf{X}) = \sigma^2 \mathbf{I}$ ($\text{Var}(\epsilon_i) = \sigma^2$ and $\text{Cov}(\epsilon_i, \epsilon_j) = 0$)

(A5) Independent and identical data: $\{(y_i, \mathbf{x}_i)\}$ are i.i.d.

If the regressors can be treated as nonstochastic (as they would be in an experiment situation in which the analyst choose the values in \mathbf{X}), \mathbf{X} can be treated as constant matrix. If \mathbf{X} is stochastic (random variables), the analysis should be done conditioned on the observed \mathbf{X} . For notation simplicity, in the following text, \mathbf{X} is treated as constant wherever possible.

C. The OLS Estimator

Minimizing the sum of squared residuals:

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^n (y_i - \mathbf{x}_i' \boldsymbol{\beta})^2$$

or equivalently,

$$\min_{\beta} (Y - X\beta)(Y - X\beta)'$$

The first-order condition gives

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

D. Properties of OLS Estimator

1. Unbiasedness

$$\mathbb{E}[\hat{\beta}] = \beta + \mathbb{E}[(X'X)^{-1}X'\epsilon] = \beta + (X'X)^{-1}X'\mathbb{E}(\epsilon) = \beta$$

2. Efficiency

OLS is the most efficient unbiased linear estimator.

Proof. Let $\tilde{\beta} = C'y$ be another unbiased linear estimator. Since $\hat{\beta}$ is unbiased, $\mathbb{E}(\tilde{\beta}) = \mathbb{E}(C'y) = \mathbb{E}(C'(X\beta + \epsilon)) = C'X\beta = \beta$, which implies $C'X = I$. Therefore, $\tilde{\beta} = C'Xb + C'\epsilon = \beta + C'\epsilon$.

$$\begin{aligned} \text{Var}(\tilde{\beta}) &= \text{Var}(C'\epsilon) = \mathbb{E}(C'\epsilon\epsilon'C) = \sigma^2 C'C \\ &= \sigma^2 (X'X)^{-1} + \sigma^2 Z Z' \\ &\geq \sigma^2 (X'X)^{-1} \end{aligned}$$

where $Z = C' - (X'X)^{-1}X'$.

Gauss-Markov Theorem: the least square estimator $\hat{\beta}$ is the minimal variance (most efficient) linear unbiased estimator.

3. Consistency

$$\begin{aligned} \text{Var}(\hat{\beta}) &= \mathbb{E}[(\hat{\beta} - \beta)(\hat{\beta} - \beta)'] \\ &= \mathbb{E} \left[\left(\sum_i \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\sum_i \mathbf{x}_i \epsilon_i \right) \left(\sum_i \mathbf{x}_i \epsilon_i \right)' \left(\sum_i \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \right] \\ &= \mathbb{E} \left[\left(\sum_i \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\sum_i \sum_j \mathbf{x}_i \epsilon_i \epsilon_j \mathbf{x}_j' \right) \left(\sum_i \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \right] \\ &= \mathbb{E} \left[\left(\sum_i \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\sum_i \mathbf{x}_i \mathbf{x}_i' \sigma^2 \right) \left(\sum_i \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \right] \\ &= \sigma^2 \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} = \frac{\sigma^2}{n} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \\ &\rightarrow \mathbf{0} \text{ as } n \rightarrow \infty \end{aligned}$$

Therefore, $\hat{\beta} \xrightarrow{m.s.} \beta$ and $\hat{\beta} \xrightarrow{p} \beta$.

4. Asymptotic normality

Multiply by the “stabler” \sqrt{n} ,

$$\sqrt{n}(\hat{\beta} - \beta) = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\sqrt{n} \cdot \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \epsilon_i \right)$$

The following properties hold as $n \rightarrow \infty$,

$$\begin{aligned} &\cdot \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \xrightarrow{p} [\mathbb{E}(\mathbf{x}_i \mathbf{x}_i')]^{-1} \text{ by LLN;} \\ &\cdot \sqrt{n} \cdot \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \epsilon_i \xrightarrow{d} \mathcal{N}(\mathbf{0}, \text{Var}(\mathbf{x}_i \epsilon_i)) \text{ by CLT;} \\ &\cdot \text{Var}(\mathbf{x}_i \epsilon_i) = \mathbb{E}(\mathbf{x}_i \mathbf{x}_i' \epsilon_i^2) = \sigma^2 \mathbb{E}(\mathbf{x}_i \mathbf{x}_i') \end{aligned}$$

Therefore, $\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sigma^2 [\mathbb{E}(\mathbf{x}_i \mathbf{x}_i')]^{-1})$.

In practice, $\mathbb{E}(\mathbf{x}_i \mathbf{x}_i')$ is estimated by $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'$, and σ^2 is estimated by either

$$\begin{aligned} &\cdot \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n e_i^2, \text{ or} \\ &\cdot s^2 = \frac{e'e}{n-k}. \end{aligned}$$

where e_i and e both stand for residuals.

E. Violation of Assumptions

1. Multicollinearity

If $X'X$ is “close” to singular, i.e. $\det(X'X) \approx 0$, then $\text{Var}(\hat{\beta}) = \sigma^2 (X'X)^{-1}$ will be very large, which leads to imprecise estimator.

2. Heteroskedasticity

If $\mathbb{E}(\epsilon_i^2 | \mathbf{x}_i) = \sigma_i^2$ different for each i . Assume $\mathbb{E}(\epsilon_i \epsilon_j) = 0$ for $i \neq j$. Reevaluate the properties of OLS estimator:

$$\begin{aligned} \mathbb{E}(\hat{\beta}) &= \beta + \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \sum_{i=1}^n \mathbf{x}_i \mathbb{E}(\epsilon_i) = \beta \\ \text{Var}(\hat{\beta}) &= \frac{1}{n} \left(\frac{1}{n} \sum_i \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\frac{1}{n} \sum_i \mathbf{x}_i \mathbf{x}_i' \sigma_i^2 \right) \left(\frac{1}{n} \sum_i \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \\ &\rightarrow \mathbf{0} \text{ as } n \rightarrow \infty \end{aligned}$$

$\hat{\beta}$ is still unbiased and consistent.

Asymptotic normality:

$$\begin{aligned} \sqrt{n}(\hat{\beta} - \beta) &= \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\sqrt{n} \cdot \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \epsilon_i \right) \\ &\cdot \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \xrightarrow{p} [\mathbb{E}(\mathbf{x}_i \mathbf{x}_i')]^{-1} \text{ by LLN;} \end{aligned}$$

- Though ϵ_i is heteroskedastic, we still have $\sqrt{n} \cdot \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \epsilon_i \xrightarrow{d} \mathcal{N}(0, \frac{1}{n} \sum_{i=1}^n \text{Var}(\mathbf{x}_i \epsilon_i))$ under some conditions.

If we define:

- $\mathbf{Q} = \mathbb{E}(\mathbf{x}_i \mathbf{x}_i')$
- $\mathbf{R} = \frac{1}{n} \sum_i \text{Var}(\mathbf{x}_i \epsilon_i) = \frac{1}{n} \sum_i \mathbb{E}(\mathbf{x}_i \mathbf{x}_i') \sigma_i^2$

Then, $\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{Q}^{-1} \mathbf{R} \mathbf{Q}^{-1})$.

Therefore, in heteroskedastic case, $\hat{\beta}$ still conforms to asymptotic normality. But the variance is no longer $\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$. So traditional statistical inference based on $s^2(\mathbf{X}'\mathbf{X})^{-1}$ will be misleading.

Robust standard error:

$$\tilde{\mathbf{R}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' e_i^2$$

where e_i is the residual. It can be shown, under some conditions, $\tilde{\mathbf{R}}_n \xrightarrow{p} \mathbf{R}$.

3. Problem of dependency

If $\{\mathbf{x}_i\}$ are not independent, $\hat{\beta}$ is still unbiased, but it might not be consistent, because LLN and CLT no longer hold.

4. Endogeneity

Endogeneity problem rises when $\mathbb{E}(\mathbf{x}_i \epsilon_i) \neq \mathbf{0}$. Two sources of endogeneity: (a) measurement error; (b) omitted variable.

a. Measurement error

$$y_i = x_i^* \beta + \epsilon_i$$

Suppose x_i^* stands for the measurement error free value of x_i . Let $x_i^u = x_i^* + v_i$ where v_i is the measurement error. For simplicity, assume $\mathbb{E}(x_i^*) = \mathbb{E}(v_i) = 0$, $\text{Var}(v_i) = \sigma_v^2$. And assume the best scenario when there is a measurement error: $x_i^* \perp \epsilon_i$, $x_i^* \perp v_i$, $v_i \perp \epsilon_i$.

$$y_i = (x_i^u - v_i) \beta + \epsilon_i^u = x_i^u \beta + \epsilon_i - \beta v_i = x_i^u \beta + \epsilon_i^u$$

where $\epsilon_i^u = \epsilon_i - \beta v_i$.

$$\mathbb{E}(x_i^u \epsilon_i^u) = \mathbb{E}((x_i^* + v_i)(\epsilon_i - \beta v_i)) = \mathbb{E}(-\beta v_i^2) = -\beta \sigma_v^2 \neq 0$$

So we have endogeneity problem. If we regress y_i on x_i^u ,

$$\begin{aligned} \hat{\beta} &= \frac{\sum_{i=1}^n x_i^u y_i}{\sum_{i=1}^n (x_i^u)^2} = \beta + \frac{\sum_{i=1}^n x_i^u \epsilon_i^u}{\sum_{i=1}^n (x_i^u)^2} \\ &\xrightarrow{p} \beta + \frac{\mathbb{E}(x_i^u \epsilon_i^u)}{\text{Var}(x_i^u)} = \beta \left(1 - \frac{\sigma_v^2}{\text{Var}(x_i^*) + \sigma_v^2} \right) \end{aligned}$$

Several observations:

1. If $\sigma_v^2 = 0$, $\hat{\beta}$ is consistent.
2. The larger the measurement error σ_v^2 , the larger the bias.
3. $\hat{\beta}$ always has the same sign as β .
4. Attenuation bias: $|\hat{\beta}| \leq |\beta|$. Therefore, if the result is significant in measurement error cases, it is also significant in measurement-error free case.

b. Omitted variable

$$y_i = x_i \beta + z_i \gamma + \epsilon_i$$

Suppose $\mathbb{E}(x_i \epsilon_i) = 0$, $\text{Cov}(x_i, z_i) \neq 0$. If the variable z_i is omitted in the model,

$$y_i = x_i \beta + \delta_i$$

where $\delta_i = z_i \gamma + \epsilon_i$. Then δ_i is correlated with x_i .

To resolve the omitted variable bias, we need to use instrumental variable (IV) estimation.

IV. IV ESTIMATION

A. The IV Estimator

Suppose the structural model is

$$y_i = \mathbf{x}_i' \beta + \epsilon_i$$

where ϵ_i is correlated with \mathbf{x}_i .

Suppose \mathbf{z}_i are instruments satisfying:

- \mathbf{z}_i has the same dimension as \mathbf{x}_i ;
- $\mathbb{E}(\mathbf{z}_i \mathbf{x}_i')$ has full rank;
- $\mathbb{E}(\mathbf{z}_i \epsilon_i) = 0$.

Then we have

$$\mathbb{E}(\mathbf{z}_i \epsilon_i) = \mathbb{E}(\mathbf{z}_i (y_i - \mathbf{x}_i' \beta)) = \mathbb{E}(\mathbf{z}_i y_i) - \mathbb{E}(\mathbf{z}_i \mathbf{x}_i') \beta = 0$$

Therefore, $\beta = (\mathbb{E}(\mathbf{z}_i \mathbf{x}_i'))^{-1} \mathbb{E}(\mathbf{z}_i y_i)$.

Define the IV estimator:

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

1. Consistency

$$\hat{\beta}_{IV} = \beta + \left(\frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{x}_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \epsilon_i \right)$$

By Law of Large Numbers,

$$\cdot \left(\frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{x}_i' \right)^{-1} \xrightarrow{P} [\mathbb{E}(\mathbf{z}_i \mathbf{x}_i')]^{-1};$$

$$\cdot \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \epsilon_i \xrightarrow{P} \mathbb{E}(\mathbf{z}_i \epsilon_i) = 0.$$

Therefore, $\hat{\beta}_{IV} \xrightarrow{P} \beta$.

Note: IV estimator is generally biased, but consistent.

2. Asymptotic normality

$$\sqrt{n}(\hat{\beta}_{IV} - \beta) = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{x}_i' \right)^{-1} \left(\sqrt{n} \cdot \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \epsilon_i \right)$$

$$\cdot \left(\frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{x}_i' \right)^{-1} \xrightarrow{P} [\mathbb{E}(\mathbf{z}_i \mathbf{x}_i')]^{-1} \text{ by LLN};$$

$$\cdot \sqrt{n} \cdot \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \epsilon_i \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbb{E}(\mathbf{z}_i \mathbf{z}_i' \epsilon_i^2)) \text{ by CLT.}$$

Therefore, by Slutsky's Theorem,

$$\sqrt{n}(\hat{\beta}_{IV} - \beta) \xrightarrow{d} \mathcal{N}(\mathbf{0}, [\mathbb{E}(\mathbf{x}_i \mathbf{z}_i')]^{-1} \mathbb{E}(\mathbf{z}_i \mathbf{z}_i' \epsilon_i^2) [\mathbb{E}(\mathbf{z}_i \mathbf{x}_i')]^{-1})$$

B. 2SLS Estimator

Suppose in a more general case,

$$\mathbf{Y} = \mathbf{X} \beta + \epsilon$$

where $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_k]$, and $\mathbb{E}(\epsilon|\mathbf{X}) \neq 0$.

Suppose \mathbf{Z} are instruments, where $l \geq k$, i.e. there could be more instruments than independent variables.

The following two procedures are equivalent:

a. *IV estimation*

- 1) Regress \mathbf{X} on \mathbf{Z} , get fitted $\hat{\mathbf{X}} = \mathbf{Q}$;
- 2) Regress \mathbf{Y} on \mathbf{X} using \mathbf{Q} as the instrument.

b. *2SLS estimation*

- 1) Regress \mathbf{X} on \mathbf{Z} , get fitted $\hat{\mathbf{X}} = \mathbf{Q}$;
- 2) Regress \mathbf{Y} on $\hat{\mathbf{X}}$.

Proof. Regressing \mathbf{X} on \mathbf{Z} :

$$\begin{aligned} \hat{\gamma}_1 &= (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X}_1, & \mathbf{Q}_1 &= \mathbf{Z} \hat{\gamma}_1 \\ \hat{\gamma}_2 &= (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X}_2, & \mathbf{Q}_2 &= \mathbf{Z} \hat{\gamma}_2 \\ &\vdots & &\vdots \\ \hat{\gamma}_k &= (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X}_k, & \mathbf{Q}_k &= \mathbf{Z} \hat{\gamma}_k \end{aligned}$$

Then,

$$\begin{aligned} \hat{\mathbf{X}} &= \mathbf{Q} = [\mathbf{Q}_1 \quad \mathbf{Q}_2 \quad \dots \quad \mathbf{Q}_k] \\ &= [\mathbf{Z}(\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X}_1 \quad \dots \quad \mathbf{Z}(\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X}_k] \\ &= \mathbf{Z}(\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' [\mathbf{X}_1 \quad \dots \quad \mathbf{X}_k] \\ &= \mathbf{Z}(\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X} \end{aligned}$$

If we regress \mathbf{Y} on \mathbf{X} using \mathbf{Q} as IV, then

$$\begin{aligned} \hat{\beta}_{IV} &= (\mathbf{Q}' \mathbf{X})^{-1} (\mathbf{Q}' \mathbf{Y}) \\ &= (\mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X})^{-1} (\mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{Y}) \end{aligned}$$

If we regress \mathbf{Y} directly on $\hat{\mathbf{X}}$,

$$\begin{aligned} \hat{\beta}_{2SLS} &= (\hat{\mathbf{X}}' \hat{\mathbf{X}})^{-1} (\hat{\mathbf{X}}' \mathbf{Y}) \\ &= (\mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \cdot \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X})^{-1} \\ &\quad (\mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{Y}) \\ &= (\mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X})^{-1} (\mathbf{X}' \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{Y}) \end{aligned}$$

Therefore, $\hat{\beta}_{IV} = \hat{\beta}_{2SLS}$.

V. MAXIMUM LIKELIHOOD EXTIMATION

A. The Likelihood Function

Let $\{z_i\}$ be i.i.d. $f(z_i|\boldsymbol{\theta})$ is the *pdf* for z_i conditioned on a set of parameters $\boldsymbol{\theta}$.

The likelihood function is the joint density function:

$$L(\boldsymbol{\theta}|\mathbf{Z}) = f(z_1, \dots, z_n|\boldsymbol{\theta}) = \prod_{i=1}^n f(z_i|\boldsymbol{\theta})$$

It is usually simpler to work with the log of the likelihood function:

$$\ln L(\boldsymbol{\theta}|\mathbf{Z}) = \sum_{i=1}^n \ln f(z_i|\boldsymbol{\theta})$$

The maximum likelihood estimator (MLE) is

$$\hat{\boldsymbol{\theta}}_{ML} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \sum_{i=1}^n \ln f(z_i|\boldsymbol{\theta})$$

B. The Likelihood Inequality

Suppose $z_i \sim f(z|\boldsymbol{\theta}_0)$ where $\boldsymbol{\theta}_0$ is the true parameter. Then for any $\boldsymbol{\theta}$, the following inequality holds

$$\mathbb{E}[\ln f(z|\boldsymbol{\theta}_0)] \geq \mathbb{E}[\ln f(z|\boldsymbol{\theta})]$$

Proof.

$$\begin{aligned} \mathbb{E}[\ln f(z|\boldsymbol{\theta})] - \mathbb{E}[\ln f(z|\boldsymbol{\theta}_0)] &= \mathbb{E}[\ln f(z|\boldsymbol{\theta}) - \ln f(z|\boldsymbol{\theta}_0)] = \\ \mathbb{E} \left[\ln \frac{f(z|\boldsymbol{\theta})}{f(z|\boldsymbol{\theta}_0)} \right] &\leq \ln \mathbb{E} \left[\frac{f(z|\boldsymbol{\theta})}{f(z|\boldsymbol{\theta}_0)} \right] = \ln \int \frac{f(z|\boldsymbol{\theta})}{f(z|\boldsymbol{\theta}_0)} f(z|\boldsymbol{\theta}_0) dz = 0 \end{aligned}$$

C. Assumptions

- (A1) $\{z_i\}$ are i.i.d.
 (A2) θ_0 is the true parameter, Θ is a compact set;
 (A3) $\text{Var}(\nabla_{\theta} \ln f(z_i|\theta_0))$ is nonsingular;
 (A4) First, second and third own and cross derivatives of $\ln f(z_i|\theta)$ with respect to θ are all bounded;
 (A5) Let Ω_Z be the support of Z , then either
 (a) Ω_Z does not depend on θ , or
 (b) $f(Z|\theta) = 0$ on the boundary of Ω_Z .

D. Score Function

Define the score function:

$$s(z_i|\theta) = \nabla_{\theta} \ln f(z_i|\theta)$$

Properties of the score function:

1. $\{s(z_i|\theta)\}$ are i.i.d. because $\{z_i\}$ are i.i.d.
2. If (A5) holds, then $\mathbb{E}[s(z_i|\theta_0)] = 0$.

Proof. By definition,

$$\int_{\Omega_Z} f(z_i|\theta) dz = 1$$

Therefore, $\frac{\partial}{\partial \theta} \int_{\Omega_Z} f(z_i|\theta) dz = 0$. By Leibniz rule,

$$\int_{A(\theta)}^{B(\theta)} \frac{\partial}{\partial \theta} f(z_i|\theta) dz + B'(\theta) f(B(\theta)|\theta) - A'(\theta) f(A(\theta)|\theta) = 0$$

Under (A5), either

1. $A'(\theta) = B'(\theta) = 0$, or
2. $f(A(\theta)|\theta) = f(B(\theta)|\theta) = 0$

In either case, we conclude

$$\int_{A(\theta)}^{B(\theta)} \frac{\partial}{\partial \theta} f(z_i|\theta) dz = 0$$

Therefore,

$$\begin{aligned} \mathbb{E}[s(z_i|\theta_0)] &= \int_{\Omega_Z} s(z_i|\theta_0) f(z_i|\theta_0) dz \\ &= \int_{\Omega_Z} \frac{\ln f(z_i|\theta_0)}{\partial \theta} f(z_i|\theta_0) dz \\ &= \int_{\Omega_Z} \frac{1}{f(z_i|\theta_0)} \frac{f(z_i|\theta_0)}{\partial \theta} f(z_i|\theta_0) dz \\ &= \int_{\Omega_Z} \frac{\partial}{\partial \theta} f(z_i|\theta_0) dz = 0 \end{aligned}$$

E. Properties of MLE

1. Consistency

Under some regularity conditions, we have $\hat{\theta} \rightarrow \theta_0$.

2. Asymptotic normality

$$\hat{\theta}_{ML} = \underset{\theta \in \Theta}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \ln f(z_i|\theta)$$

The first-order condition is

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial \ln f(z_i|\hat{\theta})}{\partial \theta} = 0 \Leftrightarrow \frac{1}{n} \sum_{i=1}^n s(z_i|\hat{\theta}) = 0$$

$$\text{Let } \mathbf{H}(z_i|\theta) = \frac{\partial s(z_i|\theta)}{\partial \theta} = \frac{\partial^2 s(z_i|\theta)}{\partial \theta \partial \theta'}.$$

Taylor expand the FOC around θ_0 :

$$\frac{1}{n} \sum_{i=1}^n s(z_i|\theta_0) + \frac{1}{n} \sum_{i=1}^n \mathbf{H}(z_i|\theta_0)(\hat{\theta} - \theta_0) = 0$$

Therefore,

$$\sqrt{n}(\hat{\theta} - \theta_0) = \left(-\frac{1}{n} \sum_{i=1}^n \mathbf{H}(z_i|\theta_0) \right)^{-1} \left(\sqrt{n} \cdot \frac{1}{n} \sum_{i=1}^n s(z_i|\theta_0) \right)$$

· By Law of Large Numbers,

$$\begin{aligned} \left(-\frac{1}{n} \sum_{i=1}^n \mathbf{H}(z_i|\theta_0) \right)^{-1} &\xrightarrow{p} \mathbb{E}[-\mathbf{H}(z_i|\theta_0)]^{-1} \\ &\rightarrow \mathbb{E}[s(z_i|\theta_0)s(z_i|\theta_0)']^{-1} \end{aligned}$$

· By Central Limit Theorem,

$$\begin{aligned} \sqrt{n} \cdot \frac{1}{n} \sum_{i=1}^n s(z_i|\theta_0) &\xrightarrow{d} \mathcal{N}(\mathbf{0}, \text{Var}(s(z_i|\theta_0))) \\ &\rightarrow \mathcal{N}(\mathbf{0}, \mathbb{E}[s(z_i|\theta_0)s(z_i|\theta_0)']) \end{aligned}$$

By Slutsky's Theorem,

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbb{E}[s(z_i|\theta_0)s(z_i|\theta_0)']^{-1})$$

3. Asymptotic efficiency

$\hat{\theta}$ is asymptotic efficient, meaning $\hat{\theta}$ has an asymptotic covariance matrix that is not larger than the asymptotic covariance of any other consistent, asymptotically normally distributed estimator.

4. Invariance

If $\hat{\theta}$ is the MLE of θ , $g(\hat{\theta})$ is the MLE for $g(\theta)$ for g a continuously differentiable function.

F. Delta Method

If there is a sequence of random variables $\{x_n\}$ satisfying

$$\sqrt{n}(x_n - \theta) \xrightarrow{d} \mathcal{N}(0, \sigma^2)$$

Then for any g satisfying the property that $g'(\theta)$ exists and non-zero valued, we have

$$\sqrt{n}(g(x_n) - g(\theta)) \xrightarrow{d} \mathcal{N}(0, \sigma^2[g'(\theta)]^2)$$

G. Quasi-MLE

Consider the maximum likelihood estimation of the linear model

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + \epsilon_i$$

Assume $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ (ϵ_i may not be normal, but we assume it anyway). Then the log likelihood function is

$$\begin{aligned} L(\boldsymbol{\beta}, \sigma^2) &= \ln f(y_1, \mathbf{x}_1, \dots, y_n, \mathbf{x}_n | \boldsymbol{\beta}, \sigma^2) \\ &= \ln f(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\beta}, \sigma^2) \\ &\quad + \ln f(\mathbf{x}_1, \dots, \mathbf{x}_n | \boldsymbol{\beta}, \sigma^2) \end{aligned}$$

We assume $\boldsymbol{\beta}, \sigma^2$ do not depend on $\mathbf{x}_1, \dots, \mathbf{x}_n$, therefore $f(\mathbf{x}_1, \dots, \mathbf{x}_n | \boldsymbol{\beta}, \sigma^2) = f(\mathbf{x}_1, \dots, \mathbf{x}_n)$. We can drop it from the likelihood function without altering the result.

$$\begin{aligned} (\boldsymbol{\beta}, \sigma^2) &= \operatorname{argmax}_{\boldsymbol{\beta}, \sigma^2} \ln f(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\beta}, \sigma^2) \\ &= \operatorname{argmax}_{\boldsymbol{\beta}, \sigma^2} \sum_{i=1}^n \ln f(y_i | \mathbf{x}_i, \boldsymbol{\beta}, \sigma^2) \end{aligned}$$

Since $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$, $y_i | \mathbf{x}_i, \boldsymbol{\beta}, \sigma^2 \sim \mathcal{N}(\mathbf{x}_i' \boldsymbol{\beta}, \sigma^2)$.

$$f(y_i | \mathbf{x}_i, \boldsymbol{\beta}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \mathbf{x}_i' \boldsymbol{\beta})^2}{2\sigma^2}\right)$$

Therefore, the likelihood function can be instantiated as

$$L(\boldsymbol{\beta}, \sigma^2) = \sum_{i=1}^n \left[-\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\sigma^2) - \frac{(y_i - \mathbf{x}_i' \boldsymbol{\beta})^2}{2\sigma^2} \right]$$

The first-order condition gives

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\sum_{i=1}^n \mathbf{x}_i y_i \right) \\ \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}})^2 \end{aligned}$$

which is exactly the same as the OLS estimator.