

高级计量经济学

理论经济学博士课程 2023-2024

Lecture 6: Statistical Properties of OLS

Davidson, R. & MacKinnon, J. (2009). *Econometrics Theory and Methods*. Oxford University Press.

黄嘉平

工学博士 经济学博士
深圳大学中国经济特区研究中心 讲师

办公室 粤海校区汇文楼1510
E-mail huangjp@szu.edu.cn
Website <https://huangjp.com>

多变量分布

Multivariate Distribution

$\mathbf{x} = (X_1, X_2, \dots, X_n)^\top$ 为随机向量, 其联合密度函数为 $f(x_1, x_2, \dots, x_n)$

- 期望值

$$E[\mathbf{x}] = \begin{bmatrix} E[X_1] \\ E[X_2] \\ \vdots \\ E[X_n] \end{bmatrix} = \begin{bmatrix} \mu_{X_1} \\ \mu_{X_2} \\ \vdots \\ \mu_{X_n} \end{bmatrix} = \boldsymbol{\mu}$$

- 方差-协方差矩阵

$$\text{Var}[\mathbf{x}] = \begin{bmatrix} \sigma_{X_1}^2 & \sigma_{X_1 X_2} & \cdots & \sigma_{X_1 X_n} \\ \sigma_{X_2 X_1} & \sigma_{X_2}^2 & \cdots & \sigma_{X_2 X_n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{X_n X_1} & \sigma_{X_n X_2} & \cdots & \sigma_{X_n}^2 \end{bmatrix} = E[\mathbf{x}\mathbf{x}^\top] - \boldsymbol{\mu}\boldsymbol{\mu}^\top$$

OLS 假设

线性模型与 OLS 假设

independently and identically distributed
在给出 X 的条件下, u_t 服从独立同分布,
其均值为 0, 方差为 σ^2 , 任意 u_s 与 u_t 的协方差为 0。

完整的线性回归模型可以写成

$$y = X\beta + u, \quad u | X \sim \text{IID}(\mathbf{0}, \sigma^2 I)$$

这个简洁的表述中包含以下假设:

1. 线性 linearity: y 是 β 的线性函数
2. 解释变量 X 可以是随机的, 此时 $(y_t, x_{t1}, \dots, x_{tk})$ 为同一个联合分布的独立样本
3. 外生性 exogeneity: $E[u_t | X] = 0 \Rightarrow E[y | X] = X\beta$
4. 同方差性 homoskedasticity: $\text{Var}[u_t | X] = \sigma^2$
无自相关 non-autocorrelation: $\text{Cov}[u_s, u_t | X] = 0$

\Rightarrow 在外生性成立时,
 $\text{Var}[u_t | X] = \sigma^2 \Leftrightarrow E[u_t^2 | X] = \sigma^2$
 $\text{Cov}[u_s, u_t | X] = 0 \Leftrightarrow E[u_s u_t | X] = 0$

为了保证 OLS 估计量 $\hat{\beta} = (X^T X)^{-1} X^T y$ 存在, 我们还需要假设:

右侧两个条件可以简洁地
写成 $E[uu^T | X] = \sigma^2 I$

5. X 列满秩: $\text{rank}(X) = k$

外生性条件

The Exogeneity Condition

$E[u_t | X] = 0$ 被称为严格外生性 (strict exogeneity)，在横截面数据中容易被满足，但在时间序列数据中往往过强。一个较弱的条件是前定性 (predeterminedness) $E[u_t | X_t] = 0$ 。

由严格外生性可以导出：

$$1. E[u_t] = E_X[E[u_t | X]] = 0$$

$$2. E[uu^\top] = E_X[E[uu^\top | X]] = E[\sigma^2 I] = \sigma^2 I$$

此处假设同方差性和无自相关

因此， X 为随机变量时的结论也适用于 X 为固定变量时。

(条件) 均值独立

(Conditional) Mean Independence

均值独立性

当 $E[Y | X] = E[Y]$ 时，称 Y 均值独立于 X 。

由定义可知，均值独立性不是对称的概念， Y 均值独立于 X 不等于 X 均值独立于 Y 。

定理：

$$\begin{array}{ccc} X \text{ 与 } Y \text{ 独立} & \Rightarrow & X \text{ 与 } Y \text{ 均值独立} & \Rightarrow & X \text{ 与 } Y \text{ 线性不相关} \\ X \perp\!\!\!\perp Y & \Leftrightarrow & E[X | Y] = E[X], & \Leftrightarrow & \text{Cov}[X, Y] = 0 \Leftrightarrow E[XY] = E[X]E[Y] \\ & & E[Y | X] = E[Y] & & \end{array}$$

详细信息可参考 <https://www.econometrics.blog/post/why-econometrics-is-confusing-part-ii-the-independence-zoo/>

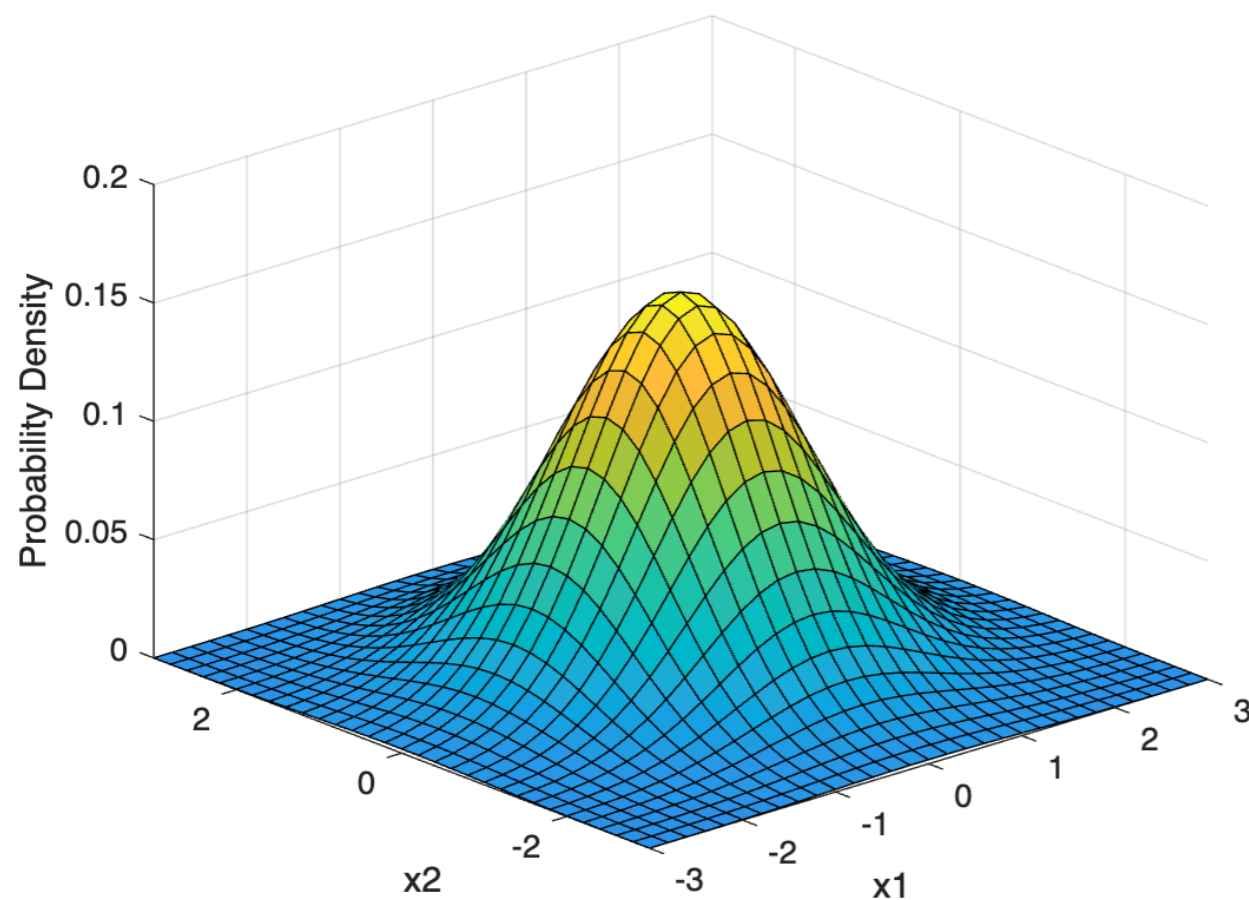
若严格外生性 $E[u_t | X] = 0$ 成立，则误差项 u 均值独立于解释变量 X 。

若 $X \perp\!\!\!\perp \mu_t$ ，则 $E[u_t | X] = E[u_t]$ 。此时可以直接假设 $E[u_t] = 0$ 。

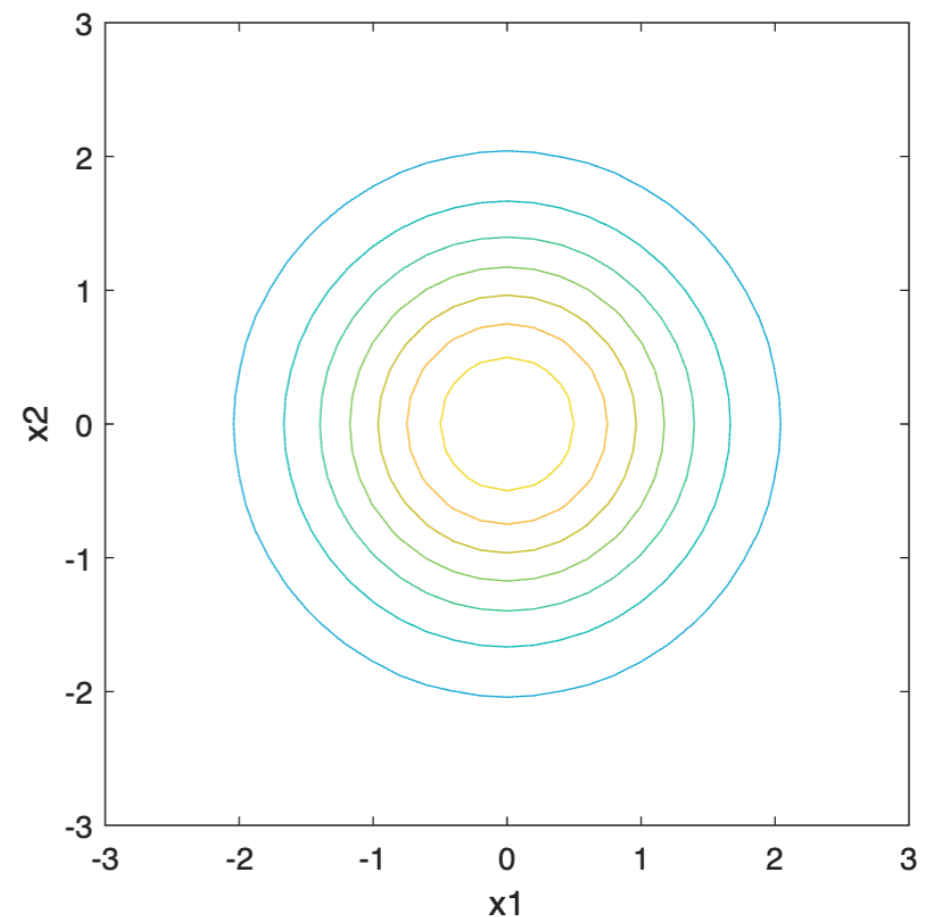
同方差无自相关性条件

The Homoskedasticity & Non-autocorrelation Condition

$\text{Var}[\mathbf{u} | \mathbf{X}] = \sigma^2 \mathbf{I}$ 常被称为 spherical error, 或称误差项服从 spherical distribution (球面分布)。



pdf of bivariate $N(\mathbf{0}, \mathbf{I})$



contour plot

2D \rightarrow 3D: circle \rightarrow sphere

$\hat{\beta}$ 的非偏性和一致性

非偏性

Unbiasedness

非偏性

若参数 θ 的真实值是 θ_0 , $\hat{\theta}$ 是 θ 的一个估计量。则称 $E[\hat{\theta}] - \theta_0$ 为估计量 $\hat{\theta}$ 的偏差 (bias) 或偏误。当 $E[\hat{\theta}] - \theta_0 = 0$, 即

$$E[\hat{\theta}] = \theta_0$$

时, $\hat{\theta}$ 是非偏的。

估计量的一致性

若参数 θ 的真实值是 θ_0 , $\hat{\theta}^n$ 是基于含有 n 个观测值样本的估计量。则当

$$\text{plim}_{n \rightarrow \infty} \hat{\theta}^n = \theta_0 \quad \text{即当 } n \rightarrow \infty \text{ 时 } \hat{\theta} \xrightarrow{p} \theta$$

时, $\hat{\theta}^n$ 是一致的。

非偏性与一致性

(x_1, x_2, \dots, x_n) 为独立随机样本，且每个 x_i 的总体均值为 $E[x_i] = \mu$ 。

已知在适当的条件下（例如总体方差有限，或服从同分布），样本均值 \bar{x} 是 μ 的非偏统计量和一致估计量（大数法则）。

考虑下面三个 μ 的估计量：

$$\hat{\mu}_1 = \frac{1}{n+1} \sum_{i=1}^n x_i = \frac{n}{n+1} \bar{x} \xrightarrow{P} \mu$$

$$\hat{\mu}_2 = \frac{1.01}{n} \sum_{i=1}^n x_i = 1.01 \bar{x} \xrightarrow{P} 1.01\mu$$

$$\hat{\mu}_3 = 0.01x_1 + \frac{0.99}{n-1} \sum_{i=2}^n x_i = 0.01x_1 + 0.99\bar{x} \xrightarrow{P} 0.01x_1 + 0.99\mu$$

其中仅 $\hat{\mu}_1$ 满足一致性，仅 $\hat{\mu}_3$ 满足非偏性。

$\hat{\beta}$ 的非偏性

假设 β 的真实值是 β_0 。 β 的 OLS 估计量 $\hat{\beta}$ 为

$$\begin{aligned}\hat{\beta} &= (X^T X)^{-1} X^T y \\ &= (X^T X)^{-1} X^T (X\beta_0 + u) \\ &= \beta_0 + (X^T X)^{-1} X^T u\end{aligned}$$

$$\Rightarrow E[\hat{\beta}] = \beta_0 + E[(X^T X)^{-1} X^T u]$$

根据迭代期望法则，右侧第二项可以写成

$$\begin{aligned}E[(X^T X)^{-1} X^T u] &= E_X[E[(X^T X)^{-1} X^T u \mid X]] \\ &= E_X[(X^T X)^{-1} X^T E[u \mid X]]\end{aligned}$$

因此，当外生性条件 $E[u \mid X] = \mathbf{0}$ 成立时，

$$E[(X^T X)^{-1} X^T u] = E_X[\mathbf{0}] = \mathbf{0}$$

$$\Rightarrow E[\hat{\beta}] = \beta_0 \quad (\hat{\beta} \text{ 是 } \beta \text{ 的非偏估计量})$$

时间序列数据往往不满足外生性，
而只满足前定性。此时 $\hat{\beta}$ 有偏。
详见 D&M (2021) pp.90-92

$\hat{\beta}$ 的一致性

$$\text{已知 } \hat{\beta} = \beta_0 + (X^\top X)^{-1} X^\top u$$

我们的目标是把右侧第二项的概率极限分解为

$$\left(\text{plim}_{n \rightarrow \infty} \frac{1}{n} X^\top X \right)^{-1} \left(\text{plim}_{n \rightarrow \infty} \frac{1}{n} X^\top u \right)$$

$X^\top X$ 是 $k \times k$ 矩阵, 其 (i, j) 要素是

$$\mathbf{x}_i^\top \mathbf{x}_j = \sum_{t=1}^n x_{ti} x_{tj}$$

一般情况下不存在 $\text{plim}_{n \rightarrow \infty} \mathbf{x}_i^\top \mathbf{x}_j$, 但有可能存在 $\text{plim}_{n \rightarrow \infty} \frac{1}{n} \mathbf{x}_i^\top \mathbf{x}_j$ 。

因为 $X^\top X = \sum_{t=1}^n X_t^\top X_t$, 则根据大数定律, $\text{plim}_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n X_t^\top X_t = E[X_t^\top X_t]$, 因此

需要假设 X_t 是 IID 的随机样本

$$\text{plim}_{n \rightarrow \infty} \frac{1}{n} X^\top X = E[X_t^\top X_t] \equiv Q \Rightarrow \text{plim}_{n \rightarrow \infty} \frac{1}{n} (X^\top X)^{-1} = Q^{-1}$$

Theorem D.14 (Greene, 2020)

积的概率极限 = 概率极限之积, 即

$$\text{plim}_{n \rightarrow \infty} X^n = A, \text{plim}_{n \rightarrow \infty} Y^n = B \Rightarrow \text{plim}_{n \rightarrow \infty} X^n Y^n = AB$$

Theorem D.14 (Greene, 2020)

逆矩阵的概率极限

$$\text{plim}_{n \rightarrow \infty} W^n = \Omega \Rightarrow \text{plim}_{n \rightarrow \infty} (W^n)^{-1} = \Omega^{-1}$$

书中将 Q 记为 $S_{X^\top X}$

$\hat{\beta}$ 的一致性

接下来我们讨论 $\text{plim}_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}^\top \mathbf{u}$ 。

因为 $\mathbf{X}^\top \mathbf{u} = \sum_{t=1}^n \mathbf{X}_t^\top u_t$ ，则根据大数法则， $\text{plim}_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \mathbf{X}_t^\top u_t = E[\mathbf{X}_t^\top u_t]$ ，因此

$$\text{plim}_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}^\top \mathbf{u} = E[\mathbf{X}_t^\top u_t]$$

如果假设前定性 $E[u_t | \mathbf{X}_t] = 0$ ，则根据 LIE，

$$\begin{aligned} E[\mathbf{X}_t^\top u_t] &= E_{\mathbf{X}_t} [E[\mathbf{X}_t^\top u_t | \mathbf{X}_t]] = E_{\mathbf{X}_t} [\mathbf{X}_t^\top \cdot E[u_t | \mathbf{X}_t]] \\ &= E_{\mathbf{X}_t} [\mathbf{X}_t^\top \cdot 0] = \mathbf{0} \end{aligned}$$

因此， $\text{plim}_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}^\top \mathbf{u} = E[\mathbf{X}_t^\top u_t] = \mathbf{0}$ 。

当样本为 **IID**，且满足前定性条件时（不需要外生性条件）， $\hat{\beta} = \beta_0 + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u}$ 右侧第二项的概率极限为 $\mathbf{0}$ ，即 $\hat{\beta}$ 是 β 的一致估计量。

关于非偏性与一致性的一些讨论

- 非偏性适用于任何大小的样本，而一致性仅适用于大样本。
- 保证 $\hat{\beta}$ 的非偏性需要假设外生性，而保证 $\hat{\beta}$ 的一致性只需要假设前定性。虽然前定性条件弱于外生性条件，但这并不代表非偏性可以推导出一致性，它们是两个不同的概念。非偏性只和期望值有关，而一致性是分布上的性质。
- 解释变量中包含被解释变量的滞后项 (lag) 的时间序列模型不满足外生性，但满足前定性。因此其 OLS 估计量有偏但一致。
- 非偏性并不一定是最好的性质。有时我们需要在非偏但方差很大的估计量和略偏但方差很小的估计量中进行选择。
- 不满足一致性有两种情况：一种是概率极限不存在（不收敛），另一种是概率极限存在但不收敛于真实的参数值。

$\hat{\beta}$ 的方差与协方差

随机向量的协方差矩阵

Covariance Matrix of Random Vectors

$\mathbf{x} = (X_1, X_2, \dots, X_n)^\top$ 为随机向量时, 其方差-协方差矩阵 (或习惯性称为协方差矩阵) 是

$$\begin{aligned}\text{Var}[\mathbf{x}] &= \begin{bmatrix} \sigma_{X_1}^2 & \sigma_{X_1 X_2} & \cdots & \sigma_{X_1 X_n} \\ \sigma_{X_2 X_1} & \sigma_{X_2}^2 & \cdots & \sigma_{X_2 X_n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{X_n X_1} & \sigma_{X_n X_2} & \cdots & \sigma_{X_n}^2 \end{bmatrix} \\ &= E[(\mathbf{x} - E[\mathbf{x}])(\mathbf{x} - E[\mathbf{x}])^\top] \\ &= E[\mathbf{x}\mathbf{x}^\top] - \boldsymbol{\mu}\boldsymbol{\mu}^\top\end{aligned}$$

若 $\boldsymbol{\mu} = E[\mathbf{x}] = \mathbf{0}$, 则 $\text{Var}[\mathbf{x}] = E[\mathbf{x}\mathbf{x}^\top]$ 。

协方差矩阵是半正定矩阵

Covariance Matrices are Positive Semidefinite

$k \times k$ 矩阵 A 是正定 (半正定) 矩阵 $\Leftrightarrow \underset{(\geq)}{x^\top A x} > 0$ for all $x \neq \mathbf{0}$.

- 对任意矩阵 B , $B^\top B$ 是半正定矩阵。如果 B 为列满秩, 则 $B^\top B$ 是正定矩阵。
- 如果 A 是正定矩阵, 则 $B^\top A B$ 是半正定。如果 B 为列满秩, 则 $B^\top A B$ 是正定。
- 如果 A 是对称正定矩阵, 则存在 $k \times k$ 矩阵 B 满足 $A = B^\top B$, 但 B 不是唯一的。(Cholesky 分解)

考虑随机向量 x 的任意线性结合 $w^\top x$:

$$\begin{aligned}\text{Var}[w^\top x] &= E[w^\top x x^\top w] - E[w^\top x]E[x^\top w] \\ &= w^\top E[x x^\top] w - w^\top E[x]E[x^\top] w \\ &= w^\top \text{Var}[x] w \geq 0\end{aligned}$$

w 非随机且已知

$w^\top x$ 为 scalar

因此, $\text{Var}[x]$ 是半正定。

大多数情况下 $\text{Var}[u]$ 是正定矩阵, 只有当 $w^\top x = 0$ 时等号才会成立。

$\hat{\beta}$ 的协方差矩阵

假设同方差性和无自相关性 $\text{Var}[u | X] = \sigma^2 I$ 。当外生性 $E[u | X] = \mathbf{0}$ 成立时， $E[uu^\top | X] = \sigma^2 I$ 。

$\hat{\beta}$ 的协方差矩阵是

$$\begin{aligned}\text{Var}[\hat{\beta}] &= E[(\hat{\beta} - E[\hat{\beta}])(\hat{\beta} - E[\hat{\beta}])^\top] \\ &= E[(\hat{\beta} - \beta_0)(\hat{\beta} - \beta_0)^\top] \\ &= E[(X^\top X)^{-1} X^\top uu^\top X (X^\top X)^{-1}] && \hat{\beta} = \beta_0 + (X^\top X)^{-1} X^\top u \\ &= E_X \left[E[(X^\top X)^{-1} X^\top uu^\top X (X^\top X)^{-1} | X] \right] \\ &= E_X [(X^\top X)^{-1} X^\top E[uu^\top | X] X (X^\top X)^{-1}] \\ &= \sigma^2 (X^\top X)^{-1}\end{aligned}$$

以上结论也可以针对 $\text{Var}[\hat{\beta} | X]$ 导出。

如果 σ^2 已知且其真实值是 σ_0^2 ，则 $\text{Var}[\hat{\beta}] = \sigma_0^2 (X^\top X)^{-1}$ 。如果 σ^2 未知，则需要对其进行估计。

预测误差的方差

令 $\gamma = \omega^\top \beta$, ω 已知。 γ 的估计量为 $\hat{\gamma} = \omega^\top \hat{\beta}$ 。

$\hat{\gamma}$ 的协方差是

$$\begin{aligned}\text{Var}[\hat{\gamma}] &= \text{Var}[\omega^\top \hat{\beta}] = E[\omega^\top (\hat{\beta} - \beta_0)(\hat{\beta} - \beta_0)^\top \omega] \\ &= \omega^\top E[(\hat{\beta} - \beta_0)(\hat{\beta} - \beta_0)^\top] \omega \\ &= \omega^\top \text{Var}[\hat{\beta}] \omega \\ &= \omega^\top \sigma_0^2 (X^\top X)^{-1} \omega\end{aligned}$$

上面的结论可以用来推导预测误差的方差。假设我们用样本外的观测值 X_s 预测 y_s , 此时预测值为 $\hat{y}_s = X_s \hat{\beta}$ 。当外生性成立时, $\hat{\beta}$ 为非偏, 因此 \hat{y}_s 也非偏。预测误差的方差是

$$\begin{aligned}\text{Var}[y_s - \hat{y}_s] &= E[(y_s - X_s \hat{\beta})^2] = E[(X_s \beta_0 + u_s - X_s \hat{\beta})^2] \\ &= E[u_s^2] + E[(X_s \beta_0 - X_s \hat{\beta})^2] \quad \text{因假设了外生性} \\ \text{Var}[X] &= E[X^2] - E[X]^2 \\ &= \sigma_0^2 + \text{Var}[X_s \hat{\beta}] \\ &= \sigma_0^2 + X_s \text{Var}[\hat{\beta}] X_s^\top \\ &= \sigma_0^2 + \sigma_0^2 X_s (X^\top X)^{-1} X_s^\top\end{aligned}$$

$\hat{\beta}$ 的有效性

估计量的有效性

Efficiency of Estimator

对同一参数可以有不同的估计量，且估计量的估计精确度一般也不同。如果估计量 a 的精确度大于估计量 b ，我们说 a 比 b 更有效 (more efficient)。有效一词可以理解为更有效地利用样本所提供的信息进行估计。

估计量 a 的估计精确度和 $\text{Var}[a]^{-1}$ 成正比，因此 a 比 b 更有效可以表达为 $\text{Var}[a]^{-1} > \text{Var}[b]^{-1}$ ，或 $\text{Var}[a] < \text{Var}[b]$ 。

- 标量： a 比 b 更有效 $\Leftrightarrow \text{Var}[a] < \text{Var}[b]$ 。
- 向量： a 比 b 更有效 $\Leftrightarrow (\text{Var}[b] - \text{Var}[a])$ 是非零半正定矩阵。

a 比 b 更有效意味着 a 中的每个要素，以及任意要素的任意线性结合都至少和 b 中的一样有效。即对任意的权重 ω ， $\gamma_a = \omega^\top a$ 应当至少和 $\gamma_b = \omega^\top b$ 一样有效，即

$$\begin{aligned}\text{Var}[\gamma_a] \leq \text{Var}[\gamma_b] &\Leftrightarrow \omega^\top \text{Var}[a] \omega \leq \omega^\top \text{Var}[b] \omega \\ &\Leftrightarrow \omega^\top (\text{Var}[b] - \text{Var}[a]) \omega \geq 0\end{aligned}$$

$\hat{\beta}$ 的有效性

已知在**外生性**条件下， β 的 OLS 估计量 $\hat{\beta}$ 是非偏的。同时 $\hat{\beta}$ 是 y 的要素的线性结合，因此我们说 $\hat{\beta}$ 是线性估计量。

这里考虑 β 的另一个线性估计量 $\tilde{\beta} = Ay$ 。（ A 是 X 的函数）

如果定义 $C = A - (X^T X)^{-1} X^T$ ，则有

$$\tilde{\beta} = Ay = (X^T X)^{-1} X^T y + Cy = \hat{\beta} + Cy$$

如果我们要求 $\tilde{\beta}$ 为非偏估计量，则需要保证 $E[\tilde{\beta}] = \beta_0$ ，即

$$\begin{aligned} E[\tilde{\beta}] &= E[Ay] = E[A(X\beta_0 + u)] \\ &= E[AX\beta_0] + E[Au] = \beta_0 \end{aligned}$$

因为 A 是 X 的函数，根据外生性和 LIE 可得 $E[Au] = \mathbf{0}$ 。因此 $\tilde{\beta}$ 非偏需要保证 $AX = I$ ，或等价的 $CX = \mathbf{0}$ 。

$\hat{\beta}$ 的有效性

当满足外生性和 $CX = \mathbf{0}$ 时, $\hat{\beta}$ 和 $\tilde{\beta} = \hat{\beta} + Cy$ 都是 β 的线性非偏估计量。

$$\begin{aligned}\text{Var}[\tilde{\beta}] &= \text{Var}[\hat{\beta} + Cy] \\ &= \text{Var}[\hat{\beta}] + \text{Var}[Cy] + 2\text{Cov}[\hat{\beta}, Cy]\end{aligned}$$

$CX = \mathbf{0} \Rightarrow Cy = Cu \Rightarrow E[Cy] = \mathbf{0}$ 。因此,

$$\begin{aligned}\text{Cov}[\hat{\beta}, Cy] &= E[(\hat{\beta} - \beta_0)(Cy)^\top] \\ &= E[(X^\top X)^{-1}X^\top uu^\top C^\top] \\ &= (X^\top X)^{-1}X^\top E[uu^\top]C^\top \\ &= (X^\top X)^{-1}X^\top \sigma_0^2 IC^\top = \mathbf{0}\end{aligned}$$

此处使用了条件
 $E[uu^\top | X] = \sigma^2 I$

因此, $\text{Var}[\tilde{\beta}] - \text{Var}[\hat{\beta}] = \text{Var}[Cy]$ 。已知协方差矩阵是半正定, 因此 $\hat{\beta}$ 比 $\tilde{\beta}$ 更有效。

Gauss-Markov Theorem

当线性回归模型 $y = X\beta + u$ 满足外生性条件 $E[u | X] = \mathbf{0}$ 和同方差无自相关条件 $E[uu^\top | X] = \sigma^2 I$ 时，OLS 估计量 $\hat{\beta}$ 是线性非偏估计量当中最有效的。即对任意线性非偏估计量 $\tilde{\beta}$ ，协方差矩阵之差

$$\text{Var}[\tilde{\beta}] - \text{Var}[\hat{\beta}]$$

是半正定矩阵。

$\hat{\beta}$ 也被称为 BLUE (Best Linear Unbiased Estimator)。

残差 \hat{u} 与误差项 u

误差项的估计量

Estimator of the Error Term

在线性模型 $y = X\beta + u$ 中， β 和 u 是未知的。我们可以将 OLS 残差 $\hat{u} = y - X\hat{\beta}$ 当作误差项 u 的估计量，并用它来估计 σ^2 。

下面是 \hat{u} 的一些统计学性质：

- 从 $\hat{\beta}$ 的一致性可以推出 \hat{u} 的一致性，即 $\text{plim}_{n \rightarrow \infty} \hat{u} = u$ 。
- 从外生性可以推出 $E[\hat{u}_t | X] = \mathbf{0}$ 。
- 从同方差无自相关性可以推出 $\text{Var}[\hat{u}_t | X] < \text{Var}[u_t | X]$ 。

$$E[\hat{u}_t | X] = \mathbf{0}$$

从残差的定义可得 $\hat{u} = M_X y = M_X X \beta_0 + M_X u = M_X u$ ，因此

$$\begin{aligned}\hat{u}_t &= u_t - X_t (X^\top X)^{-1} X^\top u && M_X u \text{ 的第 } t \text{ 行} \\ &= u_t - \sum_{s=1}^n X_t (X^\top X)^{-1} X_s^\top u_s\end{aligned}$$

因此，

$$\begin{aligned}E[\hat{u}_t | X] &= E[u_t | X] - \sum_{s=1}^n X_t (X^\top X)^{-1} X_s^\top E[u_s | X] \\ &= 0\end{aligned}$$

$$\text{Var}[\hat{u}_t | X] < \text{Var}[u_t | X]$$

从 \hat{u} 的协方差矩阵可得

$$\begin{aligned}\text{Var}[\hat{u} | X] &= \text{Var}[M_X u | X] = E[M_X u u^\top M_X | X] \\ &= M_X E[u u^\top | X] M_X = M_X \sigma^2 I M_X \\ &= \sigma^2 M_X\end{aligned}$$

令 e_t 为第 t 要素为 1 其他要素都为 0 的向量。则 P_X 的第 t 对角要素可以表达为

$$h_t = e_t^\top P_X e_t = e_t^\top P_X P_X e_t = \|P_X e_t\|^2 \geq 0$$

如果回归模型包含常数项，则 $h_t > 0$ 。又因为 $I = P_X + M_X$,

$$e_t = P_X e_t + M_X e_t \Rightarrow h_t = \|P_X e_t\|^2 \leq \|e_t\|^2 = 1$$

$\text{Var}[\hat{u}_t | X]$ 是 $\sigma^2 M_X$ 的第 t 对角要素，因此可以表达为

$$\text{Var}[\hat{u}_t | X] = \sigma^2(1 - h_t) < \sigma^2 = \text{Var}[u_t | X]$$

σ^2 的估计

Estimating σ^2

如果误差项的方差 σ^2 未知，我们就需要去估计它。根据矩估计法，可以用误差项的样本方差 $\frac{1}{n} \sum_{t=1}^n u_t^2$ 作为估计量。但是误差项无法观察，我们只能考虑用残差替代。

最简单的MM估计量是 $\hat{\sigma}^2 = \frac{1}{n} \sum_{t=1}^n \hat{u}_t^2$ 。

$$\begin{aligned} E[\hat{\sigma}^2 | \mathbf{X}] &= \frac{1}{n} \sum_{t=1}^n E[\hat{u}_t^2 | \mathbf{X}] \\ &= \frac{1}{n} \sum_{t=1}^n \text{Var}[\hat{u}_t | \mathbf{X}] = \frac{1}{n} \sum_{t=1}^n (1 - h_t) \sigma^2 \end{aligned}$$

因为 $\sum_{t=1}^n h_t = k$ (S.2.6, p.80)，可得

$$E[\hat{\sigma}^2 | \mathbf{X}] = \frac{n-k}{n} \sigma^2 < \sigma^2$$

因此， $\hat{\sigma}^2$ 有偏，不适合作为 σ^2 的估计量。而 $s^2 = \frac{1}{n-k} \sum_{t=1}^n \hat{u}_t^2$ 则是一个非偏估计量。 s 被称作回归标准误 (the standard error of regression)。我们可以用 s^2 估计 $\hat{\boldsymbol{\beta}}$ 的协方差矩阵，即

$$\widehat{\text{Var}}[\hat{\boldsymbol{\beta}}] = s^2 (\mathbf{X}^\top \mathbf{X})^{-1}$$

均方误差

均方误差

Mean Squared Error

我们在讨论 $\hat{\beta}$ 的估计精确度时假设了 $\hat{\beta}$ 的非偏性。如果一个估计量 $\tilde{\beta}$ 有偏，那就无法用它的协方差矩阵 $\text{Var}[\tilde{\beta}]$ 作为估计精确度的指标。

更具普遍性的估计精确度指标是均方误差 (mean squared error, MSE) :

$$\begin{aligned}\text{MSE}(\tilde{\beta}) &= E[(\tilde{\beta} - \beta_0)(\tilde{\beta} - \beta_0)^\top] \\ &= E[(\tilde{\beta} - E[\tilde{\beta}] + E[\tilde{\beta}] - \beta_0)(\tilde{\beta} - E[\tilde{\beta}] + E[\tilde{\beta}] - \beta_0)^\top] \\ &= E[(\tilde{\beta} - E[\tilde{\beta}])(\tilde{\beta} - E[\tilde{\beta}])^\top] + E[(E[\tilde{\beta}] - \beta_0)(E[\tilde{\beta}] - \beta_0)^\top] \\ &= \text{Var}[\tilde{\beta}] + (E[\tilde{\beta}] - \beta_0)(E[\tilde{\beta}] - \beta_0)^\top\end{aligned}$$

因此，对于单一参数估计量， $\text{MSE} = \text{Var} + \text{bias}^2$ 。当估计量是非偏时， $\text{MSE} = \text{Var}$ ，因此我们可以用协方差矩阵衡量非偏估计量的估计精确度。对于一致但有偏的估计量，则应该用 MSE。