

去除模型部分自变量的影响

Remark: 任意、维度相同的方阵 M_1, M_2 。写法 $M_1 > (\geq) M_2$ 表示 $M_1 - M_2$ (半)

正定

一、全模型 & 选模型

Full & Reduced Model

1.1 Full & Reduced Model

对于原始线性模型：

$$E(y) = X\beta \quad \text{cov}(y) = \sigma^2 I$$

满足 Gauss-Markov Assumption & $\underline{\text{rank}(X) = p}, X \in \mathbb{R}^{n \times p}$

对模型分划： $S+Q = P$

$$X = (S, Q), \quad S \in \mathbb{R}^{n \times s}, Q \in \mathbb{R}^{n \times q}$$

$$\beta = (\beta_S^\top, \beta_Q^\top)^\top, \quad \beta_S \in \mathbb{R}^s, \quad \beta_Q \in \mathbb{R}^q$$

则有原模型的 OLSE 为

$$\hat{\beta} = (\hat{\beta}_S^\top, \hat{\beta}_Q^\top)^\top = (X^\top X)^{-1} X^\top y$$

而针对 β_S 的 OLSE 为

$$\tilde{\beta}_S = (S^\top S)^{-1} S^\top y$$

$\hat{\beta}$ 为 β 的全模型的 OLSE

$\tilde{\beta}_S$ 为 β_S 的选模型的 OLSE

$$\tilde{\beta}_S \rightarrow \begin{bmatrix} \beta_S \\ \beta_Q \end{bmatrix} = \beta \leftarrow \hat{\beta} = \begin{bmatrix} \hat{\beta}_S \\ \hat{\beta}_Q \end{bmatrix}$$

1.2 Full & Reduced : $\hat{\beta}_s$ 与 $\tilde{\beta}_s$ 的关系

由分块矩阵求逆公式：

$$\begin{aligned} (X^T X)^{-1} &= \begin{bmatrix} S^T S & S^T Q \\ Q^T S & Q^T Q \end{bmatrix}^{-1} \\ &= \begin{bmatrix} (S^T S)^{-1} + A D A^T & -A D \\ - (A D)^T & D \end{bmatrix} \end{aligned}$$

其中 $A = (S^T S)^{-1} S^T Q$

$$D = \{Q^T Q - Q^T S (S^T S)^{-1} S^T Q\}^{-1}$$

$$\hat{\beta} = \begin{bmatrix} (S^T S)^{-1} + A D A^T & -A D \\ - (A D)^T & D \end{bmatrix} \begin{bmatrix} S^T y \\ Q^T y \end{bmatrix}$$

又因为 $X^T y = \begin{bmatrix} S^T y \\ Q^T y \end{bmatrix}$, 于是可求 $\hat{\beta}$:

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_s \\ \hat{\beta}_Q \end{bmatrix} = \begin{bmatrix} (S^T S)^{-1} S^T y + A D (S A - Q)^T y \\ -(S A D)^T y + D Q^T y \end{bmatrix}$$

综上可知： $\hat{\beta}_s$ 与 $\tilde{\beta}_s$ 的关系

$$\hat{\beta}_s = \tilde{\beta}_s + A D (S A - Q)^T y$$

故一般： $\hat{\beta}_s \neq \tilde{\beta}_s$

其中 $\left\{ \begin{array}{l} A = (S^T S)^{-1} S^T Q \\ D = \{Q^T Q - Q^T S (S^T S)^{-1} S^T Q\}^{-1} \end{array} \right.$

1.3 $\hat{\beta}_s$ 与 $\tilde{\beta}_s$ 何时相等

由 $\hat{\beta}_s = \tilde{\beta}_s + AD(SA - Q)^T y$ 可知, $\hat{\beta}_s = \tilde{\beta}_s$ 当:

(a). $AD = 0$, 此时有 $\hat{\beta}_s = \tilde{\beta}_s$

Remark: 由 $(X^T X)^{-1}$ 的形式, 当 $AD = 0$ 时, $X^T X$ 为对角阵

$\Rightarrow S^T Q = 0$, 即 S 与 Q 正交, 此时

$$\begin{aligned}\hat{\beta} &= (X^T X)^{-1} X^T y = \text{diag}\{(S^T S)^{-1}, (Q^T Q)^{-1}\} \cdot X^T y \\ &= \begin{bmatrix} (S^T S)^{-1} S^T y \\ (Q^T Q)^{-1} Q^T y \end{bmatrix}\end{aligned}$$

说明 $\hat{\beta} = (\hat{\beta}_s^\top, \hat{\beta}_Q^\top)^\top$ 中两部分互不相关

(或: 假设 X 每列已中心化, $S^T Q$ 表示两部分自变量的协方差)

(b) $SA - Q = 0$, 此时有 $\hat{\beta}_s = \tilde{\beta}_s$

Remark: 由 $A = (S^T S)^{-1} S^T Q$ 知, $SA - Q = 0$

$$\Rightarrow Q - SA = Q - S(S^T S)^{-1} S^T Q = \{I - S(S^T S)^{-1} S^T\} \cdot Q = 0$$

由 Null-Rank Theorem (秩零定理):

记 $\mathcal{V} = \{v \in \mathbb{R}^n : [I - S(S^T S)^{-1} S^T] \cdot v = 0\}$ 有:

$$\dim(\mathcal{V}) + \text{rank}([I - S(S^T S)^{-1} S^T]) = n$$

但 $[I - S(S^T S)^{-1} S^T]$ 类 $I - H$, 为幂等阵 \Rightarrow 它的特征值为 0 或 1

$$\therefore \text{rank}([I - S(S^T S)^{-1} S^T]) = \sum \lambda_i = \text{tr}([I - S(S^T S)^{-1} S^T])$$

$$= n - \text{tr}(S(S^T S)^{-1} S^T) = n - \text{tr}(S S^T (S^T S)^{-1}) = n - \text{tr}(S^T S (S^T S)^{-1})$$

$$= n - \text{Tr}(I_s) = n-s \quad (S \in \mathbb{R}^s)$$

$$\Rightarrow \dim(V) = s$$

由 $Q \in \mathbb{R}^q$, 若 $\text{rank}(Q) = q > s = \dim(V)$, 则 $Q - SA = 0$ 必矛盾!

二、模型选择对估计的影响

Effect of model selection on estimation

考虑, $\hat{\beta}_s$ 与 $\tilde{\beta}_s$ 的期望, 协方差矩阵, 均方误差矩阵

阵 (mean squared error matrix, MSE), 以及对应的 σ^2 的估计量, 比较二者。

2.1 E & Cov

定理 1 β_s 选模型的估计 $\hat{\beta}_s = (S^T S)^{-1} S^T y$ 满足:

$$(a) \quad E(\hat{\beta}_s) = \beta_s + A\beta_Q$$

$$(b) \quad \text{cov}(\hat{\beta}_s) = \sigma^2 (S^T S)^{-1}$$

$$(c) \quad \text{cov}(\hat{\beta}_s) \geq \text{cov}(\tilde{\beta}_s)$$

$$\text{其中 } A = (S^T S)^{-1} S^T Q \quad \hat{\beta} = (X^T X)^{-1} X^T y = (\hat{\beta}_s^T, \hat{\beta}_Q^T)^T \quad \beta = (\beta_s^T, \beta_Q^T)^T$$

证明见习作

(1) 定理 1 (a): $\tilde{\beta}_s$ 一般是 β_s 的有偏估计

当 $\beta_Q = 0$ 或 $A = 0$ 时, 才有无偏性 $E(\tilde{\beta}_s) = \beta_s$

(i) $\beta_Q = 0$ 则 $y = X\beta + \varepsilon = S\beta + \varepsilon$, 显然删去 Q 合理

(ii) $A = (S^T S)^{-1} S^T Q = 0 \Rightarrow S^T Q = 0$ 说明 $\hat{\beta} = (\hat{\beta}_S^\top, \hat{\beta}_Q^\top)^\top$ 中

两部分互不相关, 可分别回归

(2) 定理 1 (c): $\hat{\beta}_S$ 相较于 $\hat{\beta}$ 的波动更小,

Remark: It's also a trade-off between biasedness and variance

2.2 MSEM

定义 MSEM 对于 b 的任意估计量 \hat{b} , 定义其 MSEM 为

$$\begin{aligned} MSEM(\hat{b}) &= E\{(\hat{b}-b)(\hat{b}-b)^\top\} \\ &= \text{cov}(\hat{b}) + [E(\hat{b})-b]\{E(\hat{b})-b\}^\top \end{aligned}$$

相比于 $MSE(\hat{b}) = E\{(\hat{b}-b)^\top(\hat{b}-b)\}$ 为标量, $MSEM(\hat{b})$ 为矩阵

$$\begin{aligned} \text{Proof: } E\{(\hat{b}-b)(\hat{b}-b)^\top\} &= E\{(\hat{b}-E(\hat{b})+E(\hat{b})-b)(\hat{b}-E(\hat{b})+E(\hat{b})-b)^\top\} \\ &= E\{(\hat{b}-E(\hat{b}))(\hat{b}-E(\hat{b}))^\top + (E(\hat{b})-b)(E(\hat{b})-b)^\top\} + 0 + 0 \\ &= \text{cov}(\hat{b}) + [E(\hat{b})-b]\{E(\hat{b})-b\}^\top \quad \square \end{aligned}$$

推论: MSEM 比 MSE 结论更强。已知 \hat{a}, \hat{b} 是对 a, b 的估计。若有 $MSEM(\hat{a}) - MSEM(\hat{b}) \geq 0$ (PSD)

\Rightarrow 则有 $MSE(\hat{a}) - MSE(\hat{b}) \geq 0$ 。反之不成立

Proof: 引理:

$$MSE(\hat{a}) - MSE(\hat{b}) = \text{tr} \{ MSEM(\hat{a}) - MSEM(\hat{b}) \}$$

$$\begin{aligned} MSE(\hat{a}) - MSE(\hat{b}) &= E\{(\hat{a}-a)^\top(\hat{a}-a)\} - E\{(\hat{b}-b)^\top(\hat{b}-b)\} \\ &= E\text{tr}\{(\hat{a}-a)^\top(\hat{a}-a)\} - E\text{tr}\{(\hat{b}-b)^\top(\hat{b}-b)\} \\ &= E\text{tr}\{(\hat{a}-a)(\hat{a}-a)^\top\} - E\text{tr}\{(\hat{b}-b)(\hat{b}-b)^\top\} \\ &= \text{tr} E\{(\hat{a}-a)(\hat{a}-a)^\top\} - \text{tr} E\{(\hat{b}-b)(\hat{b}-b)^\top\} \\ &= \text{tr} \{ MSEM(\hat{a}) - MSEM(\hat{b}) \} \end{aligned}$$

故, $MSEM(\hat{a}) - MSEM(\hat{b}) \geq 0$ (PSD) 时, 有 $\lambda_i \geq 0$

于是 $\text{tr} \{ MSEM(\hat{a}) - MSEM(\hat{b}) \} = \sum \lambda_i \geq 0 \quad \square$

2.3 一定条件下, 选模型优化了 MSEM

定理 2

若 $\text{cov}(\hat{\beta}_Q) \geq \beta_Q \beta_Q^\top$, 则 $MSEM(\hat{\beta}_S) \geq MSEM(\tilde{\beta}_S)$

Proof: 由定理 1 (a) (b) 知

$$E(\tilde{\beta}_S) = \beta_S + A\beta_Q$$

$$\text{cov}(\tilde{\beta}_S) = \sigma^2 (S^\top S)^{-1}$$

$$\text{故 } MSEM(\tilde{\beta}_S) = \text{cov}(\tilde{\beta}_S) + [E(\tilde{\beta}_S) - \beta_S][E(\tilde{\beta}_S) - \beta_S]^\top$$

$$= \underline{\sigma^2 (S^\top S)^{-1} + A\beta_Q \beta_Q^\top A^\top}$$

$$\text{由 } \text{cov}(\hat{\beta}) = \text{cov} \left[\begin{bmatrix} \hat{\beta}_S \\ \hat{\beta}_Q \end{bmatrix} \right] = \sigma^2 (X^\top X)^{-1}$$

$$= \sigma^2 \begin{bmatrix} (S^T S)^{-1} + ADA^T & -AD \\ - (AD)^T & D \end{bmatrix}$$

$$\therefore \text{cov}(\hat{\beta}_s) = \sigma^2 \cdot \{(S^T S)^{-1} + ADA^T\}$$

$$\text{cov}(\hat{\beta}_Q) = \sigma^2 \cdot D$$

$$\text{从而 } \text{MSEM}(\hat{\beta}_s) = \text{cov}(\hat{\beta}_s) + 0 = \underline{\sigma^2 \cdot \{(S^T S)^{-1} + ADA^T\}}$$

$$\begin{aligned} \Rightarrow \text{MSEM}(\hat{\beta}_s) - \text{MSEM}(\tilde{\beta}_s) &= \sigma^2 ADA^T - A\beta_Q\beta_Q^T A^T \\ &= A \text{cov}(\hat{\beta}_Q) A^T - A\beta_Q\beta_Q^T A^T \\ &= A \{ \text{cov}(\hat{\beta}_Q) - \beta_Q\beta_Q^T \} A^T \geq 0 \end{aligned}$$

当 $\text{cov}(\hat{\beta}_Q) - \beta_Q\beta_Q^T \geq 0$ 时 \square

Remark: 定理2 说明 Q 对 y 无影响时, ($\beta_Q = 0$), $\tilde{\beta}_s$ 在 MSEM 标准下优于 $\hat{\beta}_s$

定理2 说明即使 Q 对 y 有影响, 当 $\hat{\beta}_Q$ 波动很大 ($\text{cov}(\hat{\beta}_Q) \geq \beta_Q\beta_Q^T$), $\tilde{\beta}_s$ 在 MSEM 标准下优于 $\hat{\beta}_s$

2.3 对 σ^2 的估计

基于 Reduced Model 对 σ^2 估计:

$$\begin{aligned} \tilde{\sigma}^2 &= \frac{1}{n-s} \|y - s\tilde{\beta}_s\|^2 \\ &= \frac{1}{n-s} \cdot y^T \{I - s(S^T S)^{-1} S^T\} y \end{aligned}$$

下面定理说明 $\tilde{\sigma}^2$ 是 σ^2 的有偏估计：

定理 3 选模型 (Reduced Model) 的 σ^2 估计量 $\tilde{\sigma}^2$ 满足：

$$E(\tilde{\sigma}^2) = \sigma^2 + \frac{1}{n-s} \beta_Q^\top D^{-1} \beta_Q$$

证明见习作。

回顾： $D = \{Q^\top Q - Q^\top S(S^\top S)^{-1} S^\top Q\}^{-1}$

$$\begin{aligned} S(S^\top S)^{-1} S^\top = I &\Rightarrow E(y^\top (I-H)y) = E(y^\top (I-H)) \cdot E(y) + \text{tr}\{Cov(I-H)y\} \\ \tilde{\sigma}^2 = y^\top (I-H)y / (n-s) &= \beta^\top X^\top (I-H)X\beta + \frac{\text{tr}\{(I-H)\cdot \sigma^2 \cdot I \cdot (I-H)\}}{(n-s)\cdot \sigma^2} \end{aligned}$$

$$\beta = \begin{bmatrix} \beta_S \\ \beta_Q \end{bmatrix} X = (S, Q)$$

三、模型选择对预测的影响

$$\beta^\top X^\top (I-H)X\beta$$

Effect of Model Selection on Prediction

11

$$(\beta_S^\top S^\top + \beta_Q^\top Q^\top)(I-H)(S\beta_S + Q\beta_Q)$$

假设新观测值 (x_0, y_0) 满足：

$$y_0 \perp \{y_i : i=1, 2, \dots, n\}$$

$$E(y_0) = x_0^\top \beta$$

$$\text{var}(y_0) = \sigma^2$$

$$\text{而 } \beta^\top S^\top (I-H)S\beta_S = 0$$

$$\therefore \beta^\top = \beta_Q^\top Q^\top (I-H)Q\beta_Q$$

$$\begin{aligned} D^{-1} &= Q^\top Q - Q^\top S(S^\top S)^{-1} S^\top Q \\ &= Q^\top Q - Q^\top I - Q \end{aligned}$$

基于全模型和选模型，分别得到对 y_0 的两个预测：

$$x_0^\top \hat{\beta}$$

$$x_0^\top \tilde{\beta}_S$$

其中 $x_0 = (x_S^\top, x_Q^\top)^\top$

记两个模型的预测误差为 u 和 \tilde{u} ：

$$u = y_0 - x_0^\top \hat{\beta}$$

$$\tilde{u} = y_0 - x_0^\top \tilde{\beta}_S$$

下面定理比较 u, \tilde{u} 的质量：

定理 4 预测误差 u 与 \tilde{u} 满足：

(a). $E(\tilde{u}) = x_Q^\top \beta_Q - x_s^\top A \beta_Q$ 而 $E(u) = 0$

(b). $\text{var}(u) \geq \text{var}(\tilde{u})$

(c). 若 $\text{cov}(\hat{\beta}_Q) \geq \beta_Q \beta_Q^\top$, 则 $E(u^2) \geq E(\tilde{u}^2)$

证明见习作

(1) 定理 4 (a) : $x_s^\top \tilde{\beta}_s$ 一般是 $E(\beta_0)$ 的有偏估计。

若 $\beta_Q = 0$ 或 $A = 0$, 同上讨论 \Rightarrow 此时 $E(x_s^\top \tilde{\beta}_s) = E(y_0)$

(2) 定理 4 (b) : 相比于 $x_0^\top \hat{\beta}$, 选模型 $x_s^\top \tilde{\beta}_s$ 的误差波动更小

(3) 定理 4 (c) : 综合考量