

# 第30讲 基于子模型的预测

2020.6.5

*Essentially, all models are wrong, but some are useful*

- Box

# 线性模型中的预测问题

训练数据:  $(\mathbf{x}_i, y_i), i = 1, \dots, n$ ,  $\mathbf{x}_i$  为  $p \times 1$  自变量,  $y_i$  为响应变量。

模型:  $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i, \varepsilon_i \sim (0, \sigma^2)$

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim (0, \sigma^2 I_n),$$

对  $\boldsymbol{\beta}$  的任一估计  $\tilde{\boldsymbol{\beta}}$ ,  $\mathbf{y}$  的拟合值向量  $\tilde{\mathbf{y}} = X\tilde{\boldsymbol{\beta}}$ .

待预测数据:  $(\mathbf{x}_0, y_0)$ ,  $\mathbf{x}_0$  所对应的  $y_0$  待预测。  $y_0$  与  $y_1, \dots, y_n$  独立。

假设同样的模型:  $y_0 = \mathbf{x}_0^\top \boldsymbol{\beta} + \varepsilon_0, \varepsilon_0 \sim (0, \sigma^2)$

若  $\mathbf{x}_0$  等于某个  $\mathbf{x}_i, i \in \{1, 2, \dots, n\}$ , 称为 in-sample 问题;  
若  $\mathbf{x}_0$  不等于任何  $\mathbf{x}_i$ , 称为 out-of-sample 问题。

预测统计量: 设  $\tilde{\boldsymbol{\beta}}$  为  $\boldsymbol{\beta}$  的一个估计, 以  $\tilde{y}_0 = \mathbf{x}_0^\top \tilde{\boldsymbol{\beta}}$  预测  $y_0$ .

命题1. 假设  $\tilde{\boldsymbol{\beta}}$  为  $\boldsymbol{\beta}$  的任一估计, 记其MSE矩阵  $M(\tilde{\boldsymbol{\beta}}) = E(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta})(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top$ 。

我们以  $\tilde{y}_0 = \mathbf{x}_0^\top \tilde{\boldsymbol{\beta}}$  预测  $y_0 = \mathbf{x}_0^\top \boldsymbol{\beta} + \varepsilon_0$ , 则  $\tilde{y}_0$  预测误差

$$e(\tilde{y}_0) = E(\tilde{y}_0 - y_0)^2 = m(\tilde{y}_0) + \sigma^2 = \mathbf{x}_0^\top M(\tilde{\boldsymbol{\beta}}) \mathbf{x}_0 + \sigma^2,$$

其中均方误差  $m(\tilde{y}_0) = E(\tilde{y}_0 - \mathbf{x}_0^\top \boldsymbol{\beta})^2 = \mathbf{x}_0^\top M(\tilde{\boldsymbol{\beta}}) \mathbf{x}_0$ 。

证: 被预测变量  $y_0 = \mathbf{x}_0^\top \boldsymbol{\beta} + \varepsilon_0$  的期望  $\theta = E(y_0) = \mathbf{x}_0^\top \boldsymbol{\beta}$ , 则

$$\begin{aligned} m(\tilde{y}_0) &= E(\tilde{y}_0 - \theta)^2 = E(\mathbf{x}_0^\top \tilde{\boldsymbol{\beta}} - \mathbf{x}_0^\top \boldsymbol{\beta})^2 = E(\mathbf{x}_0^\top (\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta})(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbf{x}_0) \\ &= \mathbf{x}_0^\top E((\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta})(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top) \mathbf{x}_0 = \mathbf{x}_0^\top M(\tilde{\boldsymbol{\beta}}) \mathbf{x}_0. \end{aligned}$$

推论1. 设  $\hat{\boldsymbol{\beta}} = (X^\top X)^{-1} X^\top \mathbf{y}$  为LS估计,  $y_0$  的预测量取为  $\hat{y}_0 = \mathbf{x}_0^\top \hat{\boldsymbol{\beta}}$ ,

则预测误差  $e(\hat{y}_0) = (1 + \mathbf{x}_0^\top (X^\top X)^{-1} \mathbf{x}_0) \sigma^2$ ,

均方误差  $m(\hat{y}_0) = \mathbf{x}_0^\top M(\hat{\boldsymbol{\beta}}) \mathbf{x}_0 = \mathbf{x}_0^\top (X^\top X)^{-1} \mathbf{x}_0 \sigma^2$ 。

证:  $\text{bias}(\hat{\boldsymbol{\beta}}) = 0$ ,  $M(\hat{\boldsymbol{\beta}}) = \text{var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (X^\top X)^{-1}$ 。

令 $\tilde{\mathbf{y}} = X\tilde{\boldsymbol{\beta}}$ 为原数据 $\mathbf{y}$ 的拟合( $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ), 其MSE矩阵为:

$$M(\tilde{\mathbf{y}}) = E(\tilde{\mathbf{y}} - X\boldsymbol{\beta})(\tilde{\mathbf{y}} - X\boldsymbol{\beta})^\top = E(X\tilde{\boldsymbol{\beta}} - X\boldsymbol{\beta})(X\tilde{\boldsymbol{\beta}} - X\boldsymbol{\beta})^\top = XM(\tilde{\boldsymbol{\beta}})X^\top.$$

推论2. 在关于设计阵 $X$ 通常的假设下( $n > p$ 及 $X$ 列满秩),  $\tilde{y}_0 = \mathbf{x}_0^\top \tilde{\boldsymbol{\beta}}$ 的均方误差和预测误差与均方误差矩阵 $M(\tilde{\mathbf{y}}) = XM(\tilde{\boldsymbol{\beta}})X^\top$ 有关。

证明:  $X$  列满秩 $\text{Rank}(X) = p \Rightarrow L(X^\top) = R^p$ , 对任何 $\mathbf{x}_0 \in R^p$ , 存在 $\mathbf{a} \in R^n$ , 使得 $\mathbf{x}_0 = X^\top \mathbf{a}$ , 此时 $\tilde{y}_0 = \mathbf{x}_0^\top \tilde{\boldsymbol{\beta}} = \mathbf{a}^\top X\tilde{\boldsymbol{\beta}} = \mathbf{a}^\top \tilde{\mathbf{y}}$ 的MSE:

$$m(\tilde{y}_0) = \mathbf{x}_0^\top M(\tilde{\boldsymbol{\beta}})\mathbf{x}_0 = \mathbf{a}^\top XM(\tilde{\boldsymbol{\beta}})X^\top \mathbf{a} = \mathbf{a}^\top M(\tilde{\mathbf{y}})\mathbf{a}.$$

$\tilde{y}_0 = \mathbf{x}_0^\top \tilde{\boldsymbol{\beta}}$ 的预测误差:  $e(\tilde{y}_0) = \mathbf{a}^\top M(\tilde{\mathbf{y}})\mathbf{a} + \sigma^2$

推论2说明, 在新数据上的预测效果主要由训练数据拟合值向量的均方误差 $M(\tilde{\mathbf{y}})$ 决定, 所以后续内容我们将主要考察 $M(\tilde{\mathbf{y}})$ .

常用的有偏（压缩）预测方法有：

- 变量选择：选取部分变量(其它变量的回归系数压缩为0);
- 压缩估计方法(shrinkage)：比如： $\tilde{\boldsymbol{\beta}} = \Lambda \hat{\boldsymbol{\beta}}_{LS}$ ,  $\Lambda < I_p$
- 规则化方法 / 惩罚最小二乘：对回归系数大小进行约束
- Bayes方法：通过假设参数的随机性，实现压缩、减少参数的效果。

注：这些方法界限不一定分明，比如（1）如果压缩估计把某些估计压缩为0,则达到了选择变量的效果;(2) 规则化方法可理解为Bayes方法。

# 基于子模型的预测

全模型（真模型）： $\beta_1 : k \times 1, \beta_2 : (p-k) \times 1$

$$y = X\beta + \varepsilon = X_1\beta_1 + X_2\beta_2 + \varepsilon, \quad \varepsilon \sim (0, \sigma^2 I_n)$$

LS估计  $\hat{\beta} = (X^\top X)^{-1} X^\top y$ ，拟合值： $\hat{y} = P_X y$

子模型（选变量模型, 错误模型）：

$$y = X_1\beta_1 + \delta, \quad \delta \sim (0, \tau^2 I_n),$$

LS估计  $\tilde{\beta}_1 = (X_1^\top X_1)^{-1} X_1^\top y$ ，全模型中  $\beta$  的估计取为  $\tilde{\beta} = (\tilde{\beta}_1^\top, 0)^\top$ ，

拟合值： $\tilde{y} = X\tilde{\beta} = X_1\tilde{\beta}_1 = P_{X_1} y$

我们假设全模型是真模型，而子模型是错误的模型。在一定条件下，基于简单的较少变量的子模型比基于更多变量的真模型预测效果更好。

为了衡量子模型的预测效果，我们需要计算  $M(\tilde{\beta})$  或  $M(\tilde{y})$ 。

引理1. 设  $A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} > 0$ , 记  $A_{11\bullet 2} = A_{11} - A_{12}A_{22}^{-1}A_{21}$ ,  $A_{22\bullet 1} = A_{22} - A_{21}A_{11}^{-1}A_{12}$ , 则

(1)  $A_{11\bullet 2}^{-1} = A_{11}^{-1} + A_{11}^{-1}A_{12}A_{22\bullet 1}^{-1}A_{21}A_{11}^{-1}$ , 同样  $A_{22\bullet 1}^{-1}$  有类似表达.

(2)  $A_{11\bullet 2}^{-1}A_{12}A_{22}^{-1} = (A_{22\bullet 1}^{-1}A_{21}A_{11}^{-1})^T$ .

(3)  $A^{-1} = \begin{pmatrix} A_{11\bullet 2}^{-1} & -A_{11\bullet 2}^{-1}A_{12}A_{22}^{-1} \\ -A_{22\bullet 1}^{-1}A_{21}A_{11}^{-1} & A_{22\bullet 1}^{-1} \end{pmatrix} \geq \begin{pmatrix} A_{11}^{-1} & 0 \\ 0 & 0 \end{pmatrix}$

证:(1),(2)容易验证。由第4讲命题6,  $A^{-1} = \begin{pmatrix} A_{11\bullet 2}^{-1} & -A_{11\bullet 2}^{-1}A_{12}A_{22}^{-1} \\ -A_{22\bullet 1}^{-1}A_{21}A_{11}^{-1} & A_{22\bullet 1}^{-1} \end{pmatrix}$

$$= \begin{pmatrix} A_{11}^{-1} + A_{11}^{-1}A_{12}A_{22\bullet 1}^{-1}A_{21}A_{11}^{-1} & -A_{11}^{-1}A_{12}A_{22\bullet 1}^{-1} \\ -A_{22\bullet 1}^{-1}A_{21}A_{11}^{-1} & A_{22\bullet 1}^{-1} \end{pmatrix}$$

$$= \begin{pmatrix} A_{11}^{-1} & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} A_{11}^{-1}A_{12}A_{22\bullet 1}^{-1}A_{21}A_{11}^{-1} & -A_{11}^{-1}A_{12}A_{22\bullet 1}^{-1} \\ -A_{22\bullet 1}^{-1}A_{21}A_{11}^{-1} & A_{22\bullet 1}^{-1} \end{pmatrix}$$

$$= \begin{pmatrix} A_{11}^{-1} & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} A_{11}^{-1}A_{12} \\ -I \end{pmatrix} A_{22\bullet 1}^{-1} (A_{21}A_{11}^{-1} - I) \geq \begin{pmatrix} A_{11}^{-1} & 0 \\ 0 & 0 \end{pmatrix}$$

引理2: 对任何  $\mathbf{x} \in R^n$ ,  $\mathbf{x}^T \mathbf{x} \leq \lambda \Leftrightarrow \mathbf{x} \mathbf{x}^T \leq \lambda I_n$

证明: 注意到  $\mathbf{x} \mathbf{x}^T$  是秩为1的矩阵, 其唯一非0特征根为  $\mathbf{x}^T \mathbf{x}$ ,  
 $\mathbf{x} \mathbf{x}^T \leq \lambda I_n \Leftrightarrow \mathbf{x} \mathbf{x}^T$  的所有特征根  $\leq \lambda \Leftrightarrow \mathbf{x}^T \mathbf{x} \leq \lambda$

引理3: 实数  $\lambda > 0$ , 矩阵  $A_{n \times n} > 0$ , 向量  $\mathbf{x} \in R^n$ , 则

$$\mathbf{x}^T A^{-1} \mathbf{x} \leq \lambda \Leftrightarrow \mathbf{x} \mathbf{x}^T \leq \lambda A$$

证明: 若  $\mathbf{x}^T A^{-1} \mathbf{x} \leq \lambda$ , 即  $(A^{-1/2} \mathbf{x})^T (A^{-1/2} \mathbf{x}) \leq \lambda$

$$\Rightarrow (A^{-1/2} \mathbf{x})(A^{-1/2} \mathbf{x})^T \leq \lambda I_n, \text{ 即 } A^{-1/2} \mathbf{x} \mathbf{x}^T A^{-1/2} \leq \lambda I_n$$

$$\Rightarrow \mathbf{x} \mathbf{x}^T \leq \lambda A$$

反之, 若  $\lambda A \geq \mathbf{x} \mathbf{x}^T$ , 则  $\lambda I_n \geq A^{-1/2} \mathbf{x} \mathbf{x}^T A^{-1/2} = (A^{-1/2} \mathbf{x})(A^{-1/2} \mathbf{x})^T$ ,

$$\Rightarrow \lambda \geq (A^{-1/2} \mathbf{x})^T (A^{-1/2} \mathbf{x}) = \mathbf{x}^T A^{-1} \mathbf{x}$$



引理4: 假设模型  $\mathbf{y}_{n \times 1} = X_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1} = X_1 \boldsymbol{\beta}_1 + X_2 \boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim (0, \sigma^2 I_n)$ ,

其中  $\boldsymbol{\beta}_1$  和  $\boldsymbol{\beta}_2$  的长度分别为  $k$  和  $p-k$ , 记  $\tilde{\boldsymbol{\beta}} = \begin{pmatrix} \tilde{\boldsymbol{\beta}}_1 \\ \mathbf{0} \end{pmatrix}$ , 其中  $\tilde{\boldsymbol{\beta}}_1 = (X_1^\top X_1)^{-1} X_1^\top \mathbf{y}$ .

记  $\tilde{\mathbf{y}} = X \tilde{\boldsymbol{\beta}} = X_1 \tilde{\boldsymbol{\beta}}_1$ ,  $A = (X_1^\top X_1)^{-1} X_1^\top X_2$ , 则

$$(1) M(\tilde{\boldsymbol{\beta}}) = \begin{pmatrix} \sigma^2 (X_1^\top X_1)^{-1} & 0 \\ 0 & 0^\top \end{pmatrix} + \begin{pmatrix} A \\ -I \end{pmatrix} \boldsymbol{\beta}_2 \boldsymbol{\beta}_2^\top (A^\top, -I).$$

$$(2) M(\tilde{\mathbf{y}}) = \sigma^2 P_{X_1} + X_2^\perp \boldsymbol{\beta}_2 \boldsymbol{\beta}_2^\top X_2^{\perp \top}, \quad m(\tilde{\mathbf{y}}) = k\sigma^2 + \|X_2^\perp \boldsymbol{\beta}_2\|^2.$$

$$k = p \text{ 时, } \hat{\mathbf{y}} = X \hat{\boldsymbol{\beta}}, \quad M(\hat{\mathbf{y}}) = \sigma^2 P_X, \quad m(\hat{\mathbf{y}}) = p\sigma^2$$

$$M(\hat{\mathbf{y}}) = \text{var}(\hat{\mathbf{y}}) = X \text{var}(\hat{\boldsymbol{\beta}}) X^\top = \sigma^2 X (X^\top X)^{-1} X^\top$$

$$m(\hat{\mathbf{y}}) = \text{tr} M(\hat{\mathbf{y}}) = p\sigma^2$$

证明: (1)  $E(\tilde{\boldsymbol{\beta}}_1) = E(X_1^\top X_1)^{-1} X_1^\top \mathbf{y} = (X_1^\top X_1)^{-1} X_1^\top (X_1 \boldsymbol{\beta}_1 + X_2 \boldsymbol{\beta}_2) = \boldsymbol{\beta}_1 + A \boldsymbol{\beta}_2$   
 $\Rightarrow \text{Bias}(\tilde{\boldsymbol{\beta}}_1) = E(\tilde{\boldsymbol{\beta}}_1) - \boldsymbol{\beta}_1 = A \boldsymbol{\beta}_2,$

$$\tilde{\boldsymbol{\beta}} \text{ 的偏差 } \mathbf{b} \triangleq E(\tilde{\boldsymbol{\beta}}) - \boldsymbol{\beta} = E \begin{pmatrix} \tilde{\boldsymbol{\beta}}_1 \\ 0 \end{pmatrix} - \begin{pmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{pmatrix} = \begin{pmatrix} A \boldsymbol{\beta}_2 \\ -\boldsymbol{\beta}_2 \end{pmatrix} = \begin{pmatrix} A \\ -I \end{pmatrix} \boldsymbol{\beta}_2$$

另外,  $\text{var}(\tilde{\boldsymbol{\beta}}_1) = \text{var}((X_1^\top X_1)^{-1} X_1^\top \mathbf{y}) = \sigma^2 (X_1^\top X_1)^{-1}$ , 所以 MSE 矩阵

$$M(\tilde{\boldsymbol{\beta}}) = \text{var}(\tilde{\boldsymbol{\beta}}) + \mathbf{b} \mathbf{b}^\top = \begin{pmatrix} \sigma^2 (X_1^\top X_1)^{-1} & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} A \\ -I \end{pmatrix} \boldsymbol{\beta}_2 \boldsymbol{\beta}_2^\top (A^\top, -I)$$

(2) 因为  $\tilde{\mathbf{y}} = X_1 \tilde{\boldsymbol{\beta}}_1 = X_1 (X_1^\top X_1)^{-1} X_1^\top \mathbf{y} = P_{X_1} \mathbf{y}$ , 所以,

$$\text{偏差 } \mathbf{b}(\tilde{\mathbf{y}}) = E(\tilde{\mathbf{y}}) - X \boldsymbol{\beta} = P_{X_1} X \boldsymbol{\beta} - X \boldsymbol{\beta} = -X_2^\perp \boldsymbol{\beta}_2,$$

$$\text{var}(\tilde{\mathbf{y}}) = \text{var}(P_{X_1} \mathbf{y}) = \sigma^2 P_{X_1},$$

$$\Rightarrow M(\tilde{\mathbf{y}}) = \text{var}(\tilde{\mathbf{y}}) + \mathbf{b}(\tilde{\mathbf{y}}) \mathbf{b}(\tilde{\mathbf{y}})^\top = \sigma^2 P_{X_1} + X_2^\perp \boldsymbol{\beta}_2 (X_2^\perp \boldsymbol{\beta}_2)^\top.$$

(2)也可由(1)得到:

由  $A = (X_1^\top X_1)^{-1} X_1^\top X_2$  及(1)的结果

$$\begin{aligned}
 M(\tilde{\mathbf{y}}) &= XM(\tilde{\boldsymbol{\beta}})X^\top = (X_1, X_2) \begin{pmatrix} \sigma^2 (X_1^\top X_1)^{-1} & 0 \\ 0 & 0' \end{pmatrix} \begin{pmatrix} X_1^\top \\ X_2^\top \end{pmatrix} \\
 &\quad + (X_1, X_2) \begin{pmatrix} A \\ -I \end{pmatrix} \boldsymbol{\beta}_2 \boldsymbol{\beta}_2^\top (A^\top, -I) \begin{pmatrix} X_1^\top \\ X_2^\top \end{pmatrix} = \sigma^2 P_{X_1} + X_2^\perp \boldsymbol{\beta}_2 \boldsymbol{\beta}_2^\top X_2^{\perp\top} \\
 \Rightarrow m(\tilde{\mathbf{y}}) &= \text{tr}(M(\tilde{\mathbf{y}})) = \sigma^2 \text{tr} P_{X_1} + \text{tr}(X_2^\perp \boldsymbol{\beta}_2 \boldsymbol{\beta}_2^\top X_2^{\perp\top}) = k\sigma^2 + \|X_2^\perp \boldsymbol{\beta}_2\|^2.
 \end{aligned}$$

命题2:假设模型  $\mathbf{y}_{n \times 1} = X_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1} = X_1 \boldsymbol{\beta}_1 + X_2 \boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$ ,  $\boldsymbol{\varepsilon} \sim (0, \sigma^2 I_n)$ ,

其中  $X_1$  和  $X_2$  分别为  $n \times k$  和  $n \times (p-k)$ 。令  $\tilde{\mathbf{y}} = X\tilde{\boldsymbol{\beta}} = X_1 \tilde{\boldsymbol{\beta}}_1$ , 其中  $\tilde{\boldsymbol{\beta}} = \begin{pmatrix} \tilde{\boldsymbol{\beta}}_1 \\ \mathbf{0} \end{pmatrix}$ ,

$\tilde{\boldsymbol{\beta}}_1 = (X_1^\top X_1)^{-1} X_1^\top \mathbf{y}$ 。令  $\hat{\boldsymbol{\beta}}$  为LS估计,  $\hat{\mathbf{y}} = X\hat{\boldsymbol{\beta}}$ 。则

$$(1) \quad \text{var}(\tilde{\boldsymbol{\beta}}_1) \leq \text{var}(\hat{\boldsymbol{\beta}}_1), \quad \text{var}(\tilde{\boldsymbol{\beta}}) \leq \text{var}(\hat{\boldsymbol{\beta}}).$$

$$(2) \quad \text{若 } \|X_2^\perp \boldsymbol{\beta}_2\| \leq \sqrt{p-k} \sigma, \text{ 则 } m(\tilde{\mathbf{y}}) \leq m(\hat{\mathbf{y}}) = p\sigma^2.$$

$$(3) \quad \text{若 } \|X_2^\perp \boldsymbol{\beta}_2\| \leq \sigma, \text{ 则 } M(\tilde{\mathbf{y}}) \leq M(\hat{\mathbf{y}}) = \sigma^2 P_X, \text{ 且 } M(\tilde{\boldsymbol{\beta}}) \leq M(\hat{\boldsymbol{\beta}}).$$

注1: (1)与GM定理并不矛盾, 因为 $\tilde{\boldsymbol{\beta}}_1, \tilde{\boldsymbol{\beta}}$ 是线性有偏估计。

$$\text{注2: } H_0: \boldsymbol{\beta}_2 = \mathbf{0} \text{ 的检验 } F = \frac{\|X_2^\perp \hat{\boldsymbol{\beta}}_2\|^2}{(p-k)\hat{\sigma}^2}.$$

条件(2):  $\|X_2^\perp \boldsymbol{\beta}_2\| \leq \sqrt{p-k} \sigma$  大致等价于  $F \leq 1$ ;

条件(3):  $\|X_2^\perp \boldsymbol{\beta}_2\| \leq \sigma$  大致等价于  $F \leq \frac{1}{p-k}$ 。

证明：(1) 所有期望计算都是给定自变量条件下，

$$\text{由引理1, } \text{var}(\tilde{\boldsymbol{\beta}}) = \begin{pmatrix} \sigma^2 (X_1^\top X_1)^{-1} & 0 \\ 0 & 0 \end{pmatrix} \leq \sigma^2 (X^\top X)^{-1} = \text{var}(\hat{\boldsymbol{\beta}}) \Rightarrow \text{var}(\tilde{\boldsymbol{\beta}}_1) = \text{var}(\hat{\boldsymbol{\beta}}_1).$$

$$\text{或由 } \text{var}(\tilde{\boldsymbol{\beta}}_1) = \text{var}\left((X_1^\top X_1)^{-1} X_1^\top \mathbf{y}\right) = \sigma^2 (X_1^\top X_1)^{-1} \leq \sigma^2 (X_1^{\perp\top} X_1^\perp)^{-1} = \text{var}(\hat{\boldsymbol{\beta}}_1).$$

(2) 若  $\|X_2^\perp \boldsymbol{\beta}_2\| \leq \sqrt{p-k}\sigma$ , 则

$$m(\tilde{\mathbf{y}}) = k\sigma^2 + \|X_2^\perp \boldsymbol{\beta}_2\|^2 \leq k\sigma^2 + (p-k)\sigma^2 = p\sigma^2 = m(\hat{\mathbf{y}}).$$

(3) 因为  $M(\tilde{\mathbf{y}}) = E(X\tilde{\boldsymbol{\beta}} - X\boldsymbol{\beta})(X\tilde{\boldsymbol{\beta}} - X\boldsymbol{\beta})^\top = XM(\tilde{\boldsymbol{\beta}})X^\top$ ,

$$M(\hat{\mathbf{y}}) = E(X\hat{\boldsymbol{\beta}} - X\boldsymbol{\beta})(X\hat{\boldsymbol{\beta}} - X\boldsymbol{\beta})^\top = XM(\hat{\boldsymbol{\beta}})X^\top,$$

$$\text{所以 } M(\tilde{\boldsymbol{\beta}}) \leq M(\hat{\boldsymbol{\beta}}) \Leftrightarrow XM(\tilde{\boldsymbol{\beta}})X^\top \leq XM(\hat{\boldsymbol{\beta}})X^\top \Leftrightarrow M(\tilde{\mathbf{y}}) \leq M(\hat{\mathbf{y}}).$$

下面我们仅需证明给定条件下,  $M(\tilde{\mathbf{y}}) \leq M(\hat{\mathbf{y}})$

由引理4(2),  $M(\tilde{\mathbf{y}}) = \sigma^2 P_{X_1} + X_2^\perp \boldsymbol{\beta}_2 \boldsymbol{\beta}_2^\top X_2^{\perp\top}$ ,  $M(\hat{\mathbf{y}}) = \sigma^2 P_X$

$$M(\tilde{\mathbf{y}}) \leq M(\hat{\mathbf{y}}) \Leftrightarrow X_2^\perp \boldsymbol{\beta}_2 \boldsymbol{\beta}_2^\top X_2^{\perp\top} \leq \sigma^2 (P_X - P_{X_1}) = \sigma^2 P_{X_2^\perp} = \sigma^2 X_2^\perp (X_2^{\perp\top} X_2^\perp)^{-1} X_2^{\perp\top}$$

$$\Leftrightarrow \boldsymbol{\beta}_2 \boldsymbol{\beta}_2^\top \leq \sigma^2 (X_2^{\perp\top} X_2^\perp)^{-1}$$

(两边同时左乘 $X_2^{\perp\top}$ ,右乘 $X_2^\perp$ )

引理3

$$\Leftrightarrow \boldsymbol{\beta}_2^\top X_2^{\perp\top} X_2^\perp \boldsymbol{\beta}_2 \leq \sigma^2$$

注意到最后一式即条件  $\|X_2^\perp \boldsymbol{\beta}_2\| \leq \sigma$ , 证毕。

$M(\tilde{\beta}) \leq M(\hat{\beta})$ 的另外一个直接证明(利用命题1给出的 $M(\tilde{\beta})$ 表达式):

$$\text{由于 } M(\tilde{\beta}) = \begin{pmatrix} \sigma^2(X_1^\top X_1)^{-1} & 0 \\ 0 & 0' \end{pmatrix} + \begin{pmatrix} A \\ -I \end{pmatrix} \beta_2 \beta_2^\top (A^\top, -I), \text{ 其中 } A = (X_1^\top X_1)^{-1} X_1^\top X_2.$$

$$\begin{aligned} \text{记 } B = (X_2^{\perp\top} X_2^\perp)^{-1}, \hat{\beta} \text{ 的 MSE 矩阵: } M(\hat{\beta}) &= \sigma^2(X^\top X)^{-1} = \sigma^2 \begin{pmatrix} (X_1^\top X_1)^{-1} + ABA^\top & -AB \\ -BA^\top & B \end{pmatrix} \\ &= \sigma^2 \begin{pmatrix} (X_1^\top X_1)^{-1} & 0 \\ 0 & 0' \end{pmatrix} + \sigma^2 \begin{pmatrix} ABA^\top & -AB \\ -BA^\top & B \end{pmatrix} \end{aligned}$$

$$\begin{aligned} \text{所以 } M(\hat{\beta}) - M(\tilde{\beta}) &= \sigma^2 \begin{pmatrix} ABA^\top & -AB \\ -BA^\top & B \end{pmatrix} - \begin{pmatrix} A \\ -I \end{pmatrix} \beta_2 \beta_2^\top (A^\top, -I) \\ &= \sigma^2 \begin{pmatrix} A \\ -I \end{pmatrix} B (A^\top, -I) - \begin{pmatrix} A \\ -I \end{pmatrix} \beta_2 \beta_2^\top (A^\top, -I). \end{aligned}$$

由引理3知, 条件  $\|X_2^\perp \beta_2\|^2 = \beta_2^\top X_2^{\perp\top} X_2^\perp \beta_2 \leq \sigma^2 \Leftrightarrow \beta_2 \beta_2^\top \leq \sigma^2 (X_2^{\perp\top} X_2^\perp)^{-1} = \sigma^2 B$ ,  
所以  $M(\hat{\beta}) - M(\tilde{\beta}) \geq 0$ .