

课程主页: <http://staff.ustc.edu.cn/~ynyang/2022>

第十二讲 最小二乘法 II

2022.11.18

误差同分布假设是对所有数据之间的共性(统计意义上相同), 在数学上给出的一种方便的描述, 使得我们可以在平均意义下(期望)评价数据分析方法的效果。

Recap

- $\mathbf{y} \in R^n$ 在 $V \subset R^n$ 上的正交投影 $\hat{\mathbf{y}} = P_V \mathbf{y} \in V$ 满足
$$\mathbf{y} - P_V \mathbf{y} \perp V \Leftrightarrow \text{对任何 } \mathbf{x} \in V, (\mathbf{x}, \mathbf{y}) = (\mathbf{x}, P_V \mathbf{y})$$
- 若 $V = C(A)$, 则 $P_V = A(A^T A)^{-1} A^T$
- 对于矩阵 $B = (\mathbf{b}_1, \dots, \mathbf{b}_k)$, $P_V B = (P_V \mathbf{b}_1, \dots, P_V \mathbf{b}_k)$ 为各列的投影组成的矩阵。

A, B 为矩阵或向量,
 $A^\perp = A - P_C A, B^\perp = B - P_C B,$
则 $A^{\perp T} B^\perp = A^{\perp T} B$

$$\Sigma(x_i - \bar{x})(y_i - \bar{y}) = \Sigma(x_i - \bar{x})y_i$$

线性回归模型 $\mathbf{y}_{n \times 1} = X_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1}, \boldsymbol{\varepsilon} \sim (0, \sigma^2 I_n)$, 正则方程:

$$X^T \boldsymbol{\varepsilon} = X^T (\mathbf{y} - X \boldsymbol{\beta}) = 0,$$

(X 各列与 $\boldsymbol{\varepsilon}$ 正交, 各个变量与误差不相关)

解得 $\boldsymbol{\beta}$ 的 LS 估计

$$\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y},$$

误差方差的LS估计

拟合值 残差

- 投影 $\hat{\mathbf{y}} = P_X \mathbf{y} = X\hat{\boldsymbol{\beta}} = X(X^T X)^{-1} X^T \mathbf{y}$, 称为拟合值向量。
- $\mathbf{e} = \mathbf{y}^\perp = \mathbf{y} - \hat{\mathbf{y}} = (I_n - P_X)\mathbf{y}$ 称为残差向量。

$\hat{\boldsymbol{\beta}}$ 满足正则方程: $X^T(\mathbf{y} - X\hat{\boldsymbol{\beta}}) = X^T \mathbf{e} = 0$, 即 $\mathbf{e} \perp C(X)$,
特别地 $\mathbf{e} \perp \mathbf{1}$, 即 $\bar{e} = \mathbf{1}^T \mathbf{e} / n = 0$ (样本均值为0)

误差方 差估计

残差平方和: $RSS = \|\mathbf{e}\|^2 = \|\mathbf{y} - X\hat{\boldsymbol{\beta}}\|^2 = \mathbf{y}^T (I_n - P_X) \mathbf{y}$

σ^2 的LS估计 定义为: $\hat{\sigma}^2 = \frac{RSS}{n-p} = \frac{1}{n-p} \|\mathbf{e}\|^2$ 。

拟合优度：决定系数 R^2

模型： $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon} = \mathbf{1}\beta_0 + Z\mathbf{b} + \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon} \sim (0, \sigma^2 I_n)$

拟合值/投影： $\hat{\mathbf{y}} = P_X \mathbf{y}$; 残差/投影： $\mathbf{e} = \mathbf{y}^\perp = \mathbf{y} - \hat{\mathbf{y}} = (I_n - P_X) \mathbf{y}$

注意到 $P_X \mathbf{1} = \mathbf{1}$, 所以 $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n)^T$ 的样本均值:

$$\bar{\hat{y}} = \mathbf{1}^T \hat{\mathbf{y}} / n = \mathbf{1}^T P_X \mathbf{y} / n = \mathbf{1}^T \mathbf{y} / n = \bar{y}$$

所有残差的均值: $\bar{e} = \mathbf{1}^T \mathbf{e} / n = 0$

正交分解及
平方和分解

正交分解: $(\mathbf{y} - \mathbf{1}\bar{y}) = (\hat{\mathbf{y}} - \mathbf{1}\bar{y}) + \mathbf{e}$, $\mathbf{e} \perp \mathbf{1}, \hat{\mathbf{y}}$

平方和分解:

$$\begin{aligned} \|\mathbf{y} - \mathbf{1}\bar{y}\|^2 &= \|\hat{\mathbf{y}} - \mathbf{1}\bar{y}\|^2 + \|\mathbf{e}\|^2 \\ \text{SS}_{\text{总}} &= \text{SS}_{\text{回}} + \text{RSS} \end{aligned}$$

\bar{y}

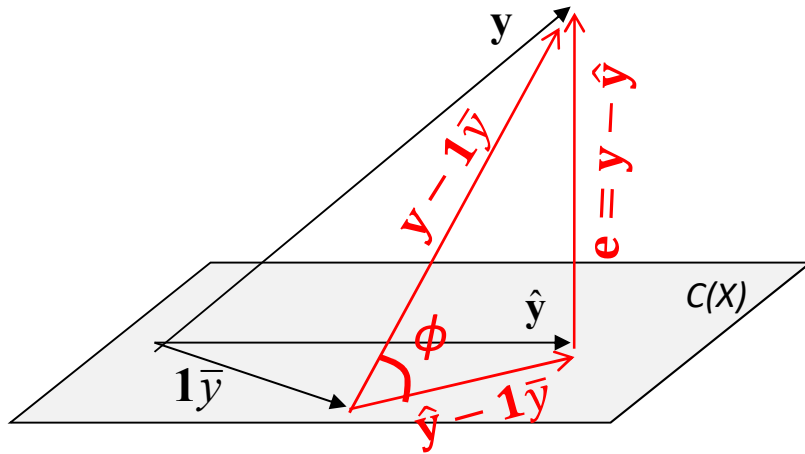
(样本版本的)
决定系数

$$R^2 = \frac{\text{SS}_{\text{回}}}{\text{SS}_{\text{总}}} = \frac{\text{VAR}(\hat{\mathbf{y}})}{\text{VAR}(\mathbf{y})} = \frac{\|\hat{\mathbf{y}} - \mathbf{1}\bar{y}\|^2}{\|\mathbf{y} - \mathbf{1}\bar{y}\|^2}, \quad 0 \leq R^2 \leq 1.$$

拟合值 $\hat{\mathbf{y}}$ 的样本方差占 \mathbf{y} 的样本方差的比例。

VAR 表示样本方差。

图示：总平方和分解 (红色)



中心化后的红色三角形：

$$\mathbf{y} - \mathbf{1}\bar{y} = (\hat{\mathbf{y}} - \mathbf{1}\bar{y}) + \mathbf{e}, \quad \mathbf{e} \perp \hat{\mathbf{y}} - \mathbf{1}\bar{y}$$

$$\|\mathbf{y} - \mathbf{1}\bar{y}\|^2 = \|\hat{\mathbf{y}} - \mathbf{1}\bar{y}\|^2 + \|\mathbf{e}\|^2$$

$$SS_{\text{总}} = SS_{\text{回}} + \text{RSS}$$

$$R^2 = \frac{\|\hat{\mathbf{y}} - \mathbf{1}\bar{y}\|^2}{\|\mathbf{y} - \mathbf{1}\bar{y}\|^2} = \cos^2(\phi) = (r_{y\hat{y}})^2$$

R^2 : y 与 \hat{y} 的相关系数的平方

命题1. (1) $R^2 = (r_{y\hat{y}})^2$, $r_{y\hat{y}}$ 为 $(y_i, \hat{y}_i), i=1, \dots, n$ 的样本相关系数。
 (2) $R^2 = \max_{\mathbf{u} \in L(X)} (r_{\mathbf{u}, y})^2$, 其中 $r_{\mathbf{u}, y}$ 为样本相关系数,
 且最大值在 $\mathbf{u} = \hat{\mathbf{y}} = X\hat{\boldsymbol{\beta}}$ 处达到, 其中 $\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y}$ 。

证明: (1) $[r_{y\hat{y}}]^2 = \frac{(\mathbf{y} - \mathbf{1}\bar{y}, \hat{\mathbf{y}} - \mathbf{1}\bar{y})^2}{\|\mathbf{y} - \mathbf{1}\bar{y}\|^2 \|\hat{\mathbf{y}} - \mathbf{1}\bar{y}\|^2}$.

因为 $\mathbf{e} \perp \hat{\mathbf{y}}, \mathbf{e} \perp \mathbf{1}$, 所以

$$(\mathbf{y} - \mathbf{1}\bar{y}, \hat{\mathbf{y}} - \mathbf{1}\bar{y}) = (\hat{\mathbf{y}} - \mathbf{1}\bar{y}, \hat{\mathbf{y}} - \mathbf{1}\bar{y}) = \|\hat{\mathbf{y}} - \mathbf{1}\bar{y}\|^2$$

$$\Rightarrow [r_{y\hat{y}}]^2 = \frac{\|\hat{\mathbf{y}} - \mathbf{1}\bar{y}\|^2}{\|\mathbf{y} - \mathbf{1}\bar{y}\|^2} = R^2$$

虚线框内容可忽略

(2) 记 $\mathbf{y}_c \triangleq \mathbf{y} - P_1 \mathbf{y} = \mathbf{y} - \mathbf{1}\bar{y}$ 。对任何 $\mathbf{u} = X\boldsymbol{\beta} = \mathbf{1}\beta_0 + Z\mathbf{b} \in C(X)$,

令 $\mathbf{u}_c \triangleq \mathbf{u} - P_1 \mathbf{u} = (\mathbf{1}\beta_0 + Z\mathbf{b}) - P_1(\mathbf{1}\beta_0 + Z\mathbf{b}) = Z\mathbf{b} - P_1 Z\mathbf{b} = Z_c \mathbf{b}$

$$\text{则 } (r_{\mathbf{u}, y})^2 = \frac{(\mathbf{u}_c^T \mathbf{y}_c)^2}{\|\mathbf{u}_c\|^2 \|\mathbf{y}_c\|^2} = \frac{(\mathbf{b}^T Z_c^T \mathbf{y}_c)^2}{\mathbf{b}^T Z_c^T Z_c \mathbf{b} \times \|\mathbf{y}_c\|^2}$$

$$\begin{aligned}
\text{令 } \mathbf{w} &= (Z_c^\top Z_c)^{1/2} \mathbf{b}, \text{ 则 } (r_{\mathbf{u}, \mathbf{y}})^2 = \frac{(\mathbf{w}^\top [(Z_c^\top Z_c)^{-1/2} Z_c^\top \mathbf{y}_c])^2}{\mathbf{w}^\top \mathbf{w} \cdot \|\mathbf{y}_c\|^2} \\
&\leq \frac{\mathbf{w}^\top \mathbf{w} \times \mathbf{y}_c^\top Z_c (Z_c^\top Z_c)^{-1} Z_c^\top \mathbf{y}_c}{\mathbf{w}^\top \mathbf{w} \times \|\mathbf{y}_c\|^2} \quad (\text{Cauchy-Schwartz不等式}) \\
&= \frac{\mathbf{w}^\top \mathbf{w} \times \|P_{Z_c} \mathbf{y}_c\|^2}{\mathbf{w}^\top \mathbf{w} \times \|\mathbf{y}_c\|^2} = \frac{\|P_{Z_c} \mathbf{y}_c\|^2}{\|\mathbf{y}_c\|^2} \quad \leftarrow P_{Z_c} \mathbf{y} = P_{Z_c} \mathbf{y}_c \\
&= \frac{\|\hat{\mathbf{y}} - \mathbf{1}\bar{y}\|^2}{\|\mathbf{y} - \mathbf{1}\bar{y}\|^2} = R^2 \quad \leftarrow \hat{\mathbf{y}} = \mathbf{1}\bar{y} + P_{Z_c} \mathbf{y}
\end{aligned}$$

当 $\mathbf{w} \propto (Z_c^\top Z_c)^{-1/2} Z_c^\top \mathbf{y}_c \Leftrightarrow \mathbf{b} = (Z_c^\top Z_c)^{-1/2} \mathbf{w} \propto (Z_c^\top Z_c)^{-1} Z_c^\top \mathbf{y}_c$ 时, 等号成立.

此时 $\mathbf{u}_c = Z_c \mathbf{b} \propto Z_c (Z_c^\top Z_c)^{-1} Z_c^\top \mathbf{y}_c = Z_c \hat{\mathbf{b}}$, 特别地在 $\mathbf{u} = Z_c \hat{\mathbf{b}} + \mathbf{1}\bar{y} = X\hat{\boldsymbol{\beta}}$ 达到极大。

LS估计的统计性质

即使没有概率模型假设，我们也能操作LS。

$\varepsilon_1, \dots, \varepsilon_n \text{ iid } \sim (0, \sigma^2)$ 的假设反映了我们对误差的认知：

$$\varepsilon_1, \dots, \varepsilon_n \text{ 大致相同, } \varepsilon_i \approx 0 \text{ 且 } |\varepsilon_i| \approx \sigma$$

该假设可用来评价LS估计是否准确，比如，LS估计 $\hat{\boldsymbol{\beta}}$ 在平均意义下是否等于或接近真正的 $\boldsymbol{\beta}$ ？LS方法的逼近效果如何，逼近误差即误差平方和RSS大概有多大？

例如， $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 中， $\text{RSS} = \|\mathbf{e}\|^2 = \mathbf{y}^\top (I_n - P_X)\mathbf{y} = (X\boldsymbol{\beta} + \boldsymbol{\varepsilon})^\top (I_n - P_X)(X\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \boldsymbol{\varepsilon}^\top (I_n - P_X)\boldsymbol{\varepsilon}$ ，平均来看（当 X 给定时）：

$$E(\text{RSS}) = E\boldsymbol{\varepsilon}^\top (I_n - P_X)\boldsymbol{\varepsilon} = \text{tr}((I_n - P_X)\text{var}(\boldsymbol{\varepsilon})) = (n - r)\sigma^2$$

其中 $r = \text{tr}(P_X) = \text{rank}(X)$.

第3讲命题1: 若 $\mathbf{x} \sim (\boldsymbol{\mu}, \Sigma)$, 则
 $E(\mathbf{x}^\top A\mathbf{x}) = \boldsymbol{\mu}^\top A\boldsymbol{\mu} + \text{tr}(A\Sigma)$.

LS估计的无偏性和方差

定理1. 假设线性模型 $\mathbf{y} = X_{n \times p} \boldsymbol{\beta} + \boldsymbol{\varepsilon}$, $E\boldsymbol{\varepsilon} = 0$, $\text{var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n$, 其中 $\boldsymbol{\varepsilon}$ 与 X 独立。假设 X 是列满秩的 ($n \geq p$, 且 $X^\top X$ 可逆), 则

(1) LS估计的无偏性: $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$, $E(\hat{\sigma}^2) = \sigma^2$.

(2) LS估计的方差: $\text{var}(\hat{\boldsymbol{\beta}} | X) = \sigma^2 (X^\top X)^{-1}$.

证明: 因为 $\boldsymbol{\varepsilon}$ 与 X 独立, 所以

- $E(\mathbf{y} | X) = E(X\boldsymbol{\beta} + \boldsymbol{\varepsilon} | X) = X\boldsymbol{\beta} + E(\boldsymbol{\varepsilon} | X) = X\boldsymbol{\beta} + E(\boldsymbol{\varepsilon}) = X\boldsymbol{\beta}$
- $\text{var}(\mathbf{y} | X) = \text{var}(X\boldsymbol{\beta} + \boldsymbol{\varepsilon} | X) = \text{var}(\boldsymbol{\varepsilon} | X) = \text{var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n$

$$(1) E(\hat{\boldsymbol{\beta}} | X) = E\left((X^\top X)^{-1} X^\top \mathbf{y} | X\right) = (X^\top X)^{-1} X^\top [E(\mathbf{y} | X)] = \boldsymbol{\beta} \Rightarrow E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}.$$

$$\begin{aligned} \text{上页我们已经证明了 } E(RSS | X) &= (n-p)\sigma^2 \Rightarrow E(RSS) = (n-p)\sigma^2 \\ \Rightarrow E(\hat{\sigma}^2) &= E(RSS / (n-p)) = \sigma^2. \end{aligned}$$

$$\begin{aligned} (2) \text{var}(\hat{\boldsymbol{\beta}} | X) &= \text{var}\left((X^\top X)^{-1} X^\top \mathbf{y} | X\right) \\ &= (X^\top X)^{-1} X^\top [\text{var}(\mathbf{y} | X)] X (X^\top X)^{-1} = \sigma^2 (X^\top X)^{-1} \end{aligned}$$

$\hat{\boldsymbol{\beta}}$ 的方差估计

在方差公式 $\text{var}(\hat{\boldsymbol{\beta}} | X) = \sigma^2 (X^T X)^{-1}$ 中代入 (Plug-in) σ^2 的估计, 即得到 $\hat{\boldsymbol{\beta}}$ 方差的估计:

$$\widehat{\text{var}(\hat{\boldsymbol{\beta}} | X)} = \hat{\sigma}^2 (X^T X)^{-1}$$

例1. 假设 (y_i, x_i) , $i = 1, 2, \dots, n$ 独立, 满足简单线性回归模型:

$$y_i = a + bx_i + \varepsilon_i,$$

第6讲命题1、2中我们已经求得LS估计及其方差, 这里我们以矩阵向量形式

再次计算如下: 记 $\boldsymbol{\beta} = \begin{pmatrix} a \\ b \end{pmatrix}$, $\mathbf{y} = (y_1, \dots, y_n)^T$, 设计阵 $X = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} = (\mathbf{1}, \mathbf{x})$, $\mathbf{x} \neq \mathbf{1c}$,

模型为 $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon} = \mathbf{1}a + \mathbf{x}b + \boldsymbol{\varepsilon}$, LS估计

$$\hat{\boldsymbol{\beta}} = \begin{pmatrix} \hat{a} \\ \hat{b} \end{pmatrix} = (X^T X)^{-1} X^T \mathbf{y} = \begin{pmatrix} n & n\bar{x} \\ n\bar{x} & \sum x_i^2 \end{pmatrix}^{-1} \begin{pmatrix} n\bar{y} \\ \sum x_i y_i \end{pmatrix}$$

$$= \frac{1}{n(\sum x_i^2 - n\bar{x}^2)} \begin{pmatrix} \sum x_i^2 & -n\bar{x} \\ -n\bar{x} & n \end{pmatrix} \begin{pmatrix} n\bar{y} \\ \sum x_i y_i \end{pmatrix} = \begin{pmatrix} \bar{y} - \bar{x}s_{xy}/s_{xx} \\ s_{xy}/s_{xx} \end{pmatrix}$$

$$\text{var}(\hat{\boldsymbol{\beta}} | X) = \sigma^2 (X^T X)^{-1} = \frac{\sigma^2}{ns_{xx}} \begin{pmatrix} \sum x_i^2 & -n\bar{x} \\ -n\bar{x} & n \end{pmatrix} = \sigma^2 \begin{pmatrix} 1/n + \bar{x}^2/s_{xx} & -\bar{x}/s_{xx} \\ -\bar{x}/s_{xx} & 1/s_{xx} \end{pmatrix}$$

LS估计的最优性

定义：对任何两个对称 $n \times n$ 矩阵 A, B ，若 $A - B \geq 0$ (非负定)，则称在Loewner偏序意义下 A 不小于 B ，记作 $A \geq B$ 。

性质：若 $A \geq B$ ，则对任何 $k \times n$ 矩阵 C ， $CAC^T \geq CBC^T$

Gauss-Markov 定理

定理2. 线性模型 $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 中，假设 $X_{n \times p}$ 列满秩，则 $\boldsymbol{\beta}$ 的最小二乘估计 $\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y}$ 是最优无偏线性估计 (BLUE, best linear unbiased estimate,)，即对任何 $\boldsymbol{\beta}$ 的线性无偏估计 $\tilde{\boldsymbol{\beta}} = C\mathbf{y}$ (其中 C 是仅与 X 有关的 $p \times n$ 常数矩阵)， $\text{var}(\tilde{\boldsymbol{\beta}} | X) \geq \text{var}(\hat{\boldsymbol{\beta}} | X)$ 。

推论：假设 $A_{k \times p}$ 是仅与 X 有关的常数矩阵，如果 A 行满秩，则 $A\hat{\boldsymbol{\beta}}$ 是 $A\boldsymbol{\beta}$ 的BLUE。

特别地， $\hat{\beta}_k$ 是 β_k 的BLUE， $\hat{\beta}_k - \hat{\beta}_j$ 是 $\beta_k - \beta_j$ 的BLUE，等等。

证明：给定 X ，设 $\tilde{\boldsymbol{\beta}} = C \mathbf{y}$ 是 $\boldsymbol{\beta}$ 的任一线性无偏估计，其中 C 是 $p \times n$ 常数矩阵(可能与 X 有关)，因为 $\text{var}(\tilde{\boldsymbol{\beta}} | X) = C \text{var}(\mathbf{y} | X) C^T = \sigma^2 C C^T$ ， $\text{var}(\hat{\boldsymbol{\beta}} | X) = \sigma^2 (X^T X)^{-1}$ 。我们需要证 $C C^T \geq (X^T X)^{-1}$ 。

因为 $\tilde{\boldsymbol{\beta}}$ 线性无偏，所以

$$\boldsymbol{\beta} = E(\tilde{\boldsymbol{\beta}} | X) = E(C\mathbf{y} | X) = C X \boldsymbol{\beta},$$

上式对任何 $\boldsymbol{\beta}$ 成立，故 $CX = I_p$ 。

因为 $CX = I_p$ ， $P_X \leq I_n$ ，所以

$$(X^T X)^{-1} = CX(X^T X)^{-1} X^T C^T = C P_X C^T \leq C C^T.$$

不确定性原理

例如，对于简单线性模型 $y_i = a + bx_i + \varepsilon_i$ 的任一无偏估计 \tilde{b} ，由GM定理

$$\text{var}(\tilde{b} | x) \geq \text{var}(\hat{b} | x) = \sigma^2 / s_{xx}$$

以 $\text{VAR}(x)$ 代表样本方差 $s_x^2 = s_{xx} / (n-1)$ ，上式等价于

$$\text{var}(\tilde{b} | x) \times \text{VAR}(x) \geq \sigma^2 / (n-1),$$

这说明左端两个方差不可能同时很小，即 x 和 \tilde{b} 不可能同时精确测量 (Uncertainty principle, 不确定性原理, 测不准原理)。

中心化

中心化是数据分析中常见的数据加工方法，它将每一个样本观测值减去样本均值，正等价于所有样本构成的向量或矩阵与向量 $\mathbf{1}$ 的正交化。

一元情形

样本: y_1, \dots, y_n ; $\mathbf{y} = (y_1, \dots, y_n)^\top$

样本均值 $\bar{y} = (y_1 + \dots + y_n) / n = \mathbf{1}^\top \mathbf{y} / n$, 其中 $\mathbf{1} = (1, \dots, 1)^\top$

中心化: $y_1 - \bar{y}, \dots, y_n - \bar{y}$ (样本均值为0)

$$\mathbf{y}_c = \begin{pmatrix} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{pmatrix} = \mathbf{y} - \mathbf{1}\bar{y} = \mathbf{y} - \mathbf{1}\mathbf{1}^\top \mathbf{y} / n = (\mathbf{I}_n - \mathbf{P}_1)\mathbf{y}$$

样本方差: $s^2 = \sum_{i=1}^n (y_i - \bar{y})^2 / (n-1) = \mathbf{y}_c^\top \mathbf{y}_c / (n-1) = \mathbf{y}^\top (\mathbf{I}_n - \mathbf{P}_1)\mathbf{y} / (n-1)$.

中心化基础上再除以标准差，称为标准化: $(y_1 - \bar{y}) / s, \dots, (y_n - \bar{y}) / s$,
标准化后样本均值为0，样本方差为1.

多元情形

样本: $\mathbf{x}_1, \dots, \mathbf{x}_n \in R^m$; 样本按行排列成 $n \times m$ 矩阵: $Z = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$

样本均值: $\bar{\mathbf{x}} = (\mathbf{x}_1 + \dots + \mathbf{x}_n) / n = Z^\top \mathbf{1} / n$.

样本协方差矩阵: $S = \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top / (n-1)$

中心化: $\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}$, 按行排列成中心化矩阵:

$$Z_c = \begin{pmatrix} (\mathbf{x}_1 - \bar{\mathbf{x}})^\top \\ \vdots \\ (\mathbf{x}_n - \bar{\mathbf{x}})^\top \end{pmatrix} = Z - \mathbf{1}\bar{\mathbf{x}}^\top = Z - \mathbf{1}\mathbf{1}^\top Z / n = (I_n - P_1)Z$$

(注意中心化矩阵 Z_c 的每列均值为 0, 即 $Z_c^\top \mathbf{1} = 0$)

所以样本协方差矩阵:

$$S = Z_c Z_c^\top / (n-1) = Z^\top (I_n - P_1) Z / (n-1).$$

如果中心化矩阵 Z_c 的每一列(变量)除以其标准差, 则称为 标准化

$$Z_s = Z_c D^{-1/2}, D = \text{diag}(S)$$

- 如果某个数据矩阵 A 已经中心化, 那么 $A^\top A / (n-1) = S$, 样本协方差矩阵。
- 如果某个数据矩阵 A 已经标准化, 那么 $A^\top A / (n-1) = R$, 样本相关系数矩阵。

线性模型的中心化

线性模型 $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 中我们把截距项单独列出来，为此划分

$$X_{n \times p} = \begin{pmatrix} 1 & \mathbf{x}_1^\top \\ \vdots & \vdots \\ 1 & \mathbf{x}_n^\top \end{pmatrix} = (\mathbf{1}_{n \times 1}, Z_{n \times (p-1)}), \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \mathbf{b} \end{pmatrix}.$$

- 模型: $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon} = \mathbf{1}\beta_0 + Z\mathbf{b} + \boldsymbol{\varepsilon}$
- Z中心化: 令中心化 $Z_c = Z - \mathbf{1}\bar{\mathbf{x}}^\top = (I_n - P_1)Z$, 模型等价地表示为

$$\mathbf{y} = \mathbf{1}\beta_0 + Z\mathbf{b} + \boldsymbol{\varepsilon} = \mathbf{1}\beta_0 + (Z_c + \mathbf{1}\bar{\mathbf{x}}^\top)\mathbf{b} + \boldsymbol{\varepsilon} = \mathbf{1}(\beta_0 + \bar{\mathbf{x}}^\top\mathbf{b}) + Z_c\mathbf{b} + \boldsymbol{\varepsilon}$$

$$\stackrel{\Delta}{=} \mathbf{1}\alpha + Z_c\mathbf{b} + \boldsymbol{\varepsilon} = X^*\boldsymbol{\beta}^* + \boldsymbol{\varepsilon}$$

其中 $X^* = (\mathbf{1}, Z_c)$, 回归系数 $\boldsymbol{\beta}^* = \begin{pmatrix} \alpha \\ \mathbf{b} \end{pmatrix}$, $\alpha = \beta_0 + \bar{\mathbf{x}}^\top\mathbf{b}$ 。

- 所以中心化不改变回归系数 \mathbf{b} , 但改变了截距项。

下面演示中心化对于计算过程的简化功能:

中心化模型 $\mathbf{y} = \mathbf{1}\alpha + Z_c\mathbf{b} + \boldsymbol{\varepsilon} = X^*\boldsymbol{\beta}^* + \boldsymbol{\varepsilon}$ 中 Z_c 的列与 $\mathbf{1}$ 正交, 这使得

$X^{*\top}X^* = \begin{pmatrix} n & 0 \\ 0 & Z_c^\top Z_c \end{pmatrix}$ 为分块对角, 对中心化模型应用最小二乘法

$$\hat{\boldsymbol{\beta}}^* = \begin{pmatrix} \hat{\alpha} \\ \hat{\mathbf{b}} \end{pmatrix} = (X^{*\top}X^*)^{-1}X^{*\top}\mathbf{y} = \begin{pmatrix} n & 0 \\ 0 & Z_c^\top Z_c \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{1}^\top \mathbf{y} \\ Z_c^\top \mathbf{y} \end{pmatrix} = \begin{pmatrix} \bar{y} \\ (Z_c^\top Z_c)^{-1} Z_c^\top \mathbf{y} \end{pmatrix}$$

所以LS估计 $\hat{\alpha} = \bar{y}$, $\hat{\mathbf{b}} = (Z_c^\top Z_c)^{-1} Z_c^\top \mathbf{y}$ 。由 $\alpha = \beta_0 + \bar{\mathbf{x}}^\top \mathbf{b}$, 得 $\hat{\beta}_0 = \bar{y} - \bar{\mathbf{x}}^\top \hat{\mathbf{b}}$ 。注意到

$$Z_c^\top Z_c = \sum (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top = (n-1)S_{\mathbf{xx}}, \quad Z_c^\top \mathbf{y} = (n-1) S_{\mathbf{xy}}$$

$$\Rightarrow \hat{\mathbf{b}} = (Z_c Z_c^\top)^{-1} Z_c^\top \mathbf{y} = S_{\mathbf{xx}}^{-1} S_{\mathbf{xy}}.$$

$$\text{var}(\hat{\mathbf{b}} | X) = \sigma^2 (Z_c^\top Z_c)^{-1} = \sigma^2 S_{\mathbf{xx}}^{-1} / (n-1).$$

另外, $\hat{\mathbf{y}} = P_1 \mathbf{y} + P_{Z_c} \mathbf{y} = \mathbf{1}\bar{y} + P_{Z_c} \mathbf{y}$, 所以 $\hat{\mathbf{y}} - \mathbf{1}\bar{y} = P_{Z_c} \mathbf{y} (= Z_c \hat{\mathbf{b}})$,

所以 $\text{SS}_{\text{回}} = \|\hat{\mathbf{y}} - \mathbf{1}\bar{y}\|^2 = \|P_{Z_c} \mathbf{y}\|^2 = \mathbf{y}^\top Z_c (Z_c^\top Z_c)^{-1} Z_c^\top \mathbf{y} = (n-1) S_{\mathbf{yx}} S_{\mathbf{xx}}^{-1} S_{\mathbf{xy}}$

而 $\text{SS}_{\text{总}} = \|\mathbf{y} - \mathbf{1}\bar{y}\|^2 = (n-1) S_{\mathbf{yy}}$

$$\Rightarrow \text{RSS} = \text{SS}_{\text{总}} - \text{SS}_{\text{回}} = (n-1) [S_{\mathbf{yy}} - S_{\mathbf{yx}} S_{\mathbf{xx}}^{-1} S_{\mathbf{xy}}] = (n-1) S_{\mathbf{yy} \cdot \mathbf{x}},$$

$$\text{所以 } R^2 = \frac{S_{\mathbf{yx}} S_{\mathbf{xx}}^{-1} S_{\mathbf{xy}}}{S_{\mathbf{yy}}}, \quad \hat{\sigma}^2 = \text{RSS} / (n-p) = \frac{n-1}{n-p} S_{\mathbf{yy} \cdot \mathbf{x}}.$$

我们有如下结论 (命题1)。

命题2. 假设模型 $y_i = \beta_0 + \mathbf{x}_i^\top \mathbf{b} + \varepsilon_i, i = 1, \dots, n$, 记 $X = (\mathbf{1}, Z)$ 列满秩,
 $Z = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$, $Z_c = Z - \mathbf{1}\bar{\mathbf{x}}^\top$, 样本方差 / 协方差矩阵

$$S_{\mathbf{xx}} = \frac{1}{n-1} \sum (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top, S_{\mathbf{xy}} = \frac{1}{n-1} \sum (\mathbf{x}_i - \bar{\mathbf{x}})(y_i - \bar{y})$$

则

(1) LS估计 $\hat{\mathbf{b}} = (Z_c Z_c^\top)^{-1} Z_c^\top \mathbf{y} = S_{\mathbf{xx}}^{-1} S_{\mathbf{xy}}$, $\hat{\beta}_0 = \bar{y} - \bar{\mathbf{x}}^\top \hat{\mathbf{b}}$, $\hat{\sigma}^2 = \frac{n-1}{n-p} S_{yy \cdot \mathbf{x}} \approx S_{yy \cdot \mathbf{x}}$.

(2) $\text{var}(\hat{\mathbf{b}} | X) = \sigma^2 (Z_c Z_c^\top)^{-1} = \sigma^2 S_{\mathbf{xx}}^{-1} / (n-1)$,

(3) $R^2 = \frac{S_{yx} S_{\mathbf{xx}}^{-1} S_{\mathbf{xy}}}{S_{yy}}$.

对比第5讲命题2(总体版本)

$$\mathbf{b} = \Sigma_{\mathbf{xx}}^{-1} \Sigma_{\mathbf{xy}}$$

$$\beta_0 = \mu_y - \mathbf{b}^\top \boldsymbol{\mu}_x$$

$$\sigma^2 = \Sigma_{yy \cdot \mathbf{x}}$$

$$R^2 = \Sigma_{yx} \Sigma_{\mathbf{xx}}^{-1} \Sigma_{\mathbf{xy}} / \Sigma_{yy}$$

命题2的第一部分 (\mathbf{b} 的LS估计) 可利用投影更简单地证明。

证1(投影): 令 $Z_c = Z - P_1 Z = Z - \mathbf{1}\bar{\mathbf{x}}^\top$, 则

$$\begin{aligned} \hat{\mathbf{y}} &= P_X \mathbf{y} = P_1 \mathbf{y} + P_{Z_c} \mathbf{y} = \mathbf{1}(\mathbf{1}^\top \mathbf{y} / n) + Z_c [(Z_c^\top Z_c)^{-1} Z_c^\top \mathbf{y}] \\ &= \mathbf{1}\bar{y} + (Z - \mathbf{1}\bar{\mathbf{x}}^\top) [S_{\mathbf{xx}}^{-1} S_{\mathbf{xy}}] = \mathbf{1}(\bar{y} - \bar{\mathbf{x}}^\top [S_{\mathbf{xx}}^{-1} S_{\mathbf{xy}}]) + Z [S_{\mathbf{xx}}^{-1} S_{\mathbf{xy}}] \end{aligned}$$

$\mathbf{1}$ 的系数为 $\hat{\beta}_0$, Z 的系数为 $\hat{\mathbf{b}}$, 所以 $\hat{\mathbf{b}} = S_{\mathbf{xx}}^{-1} S_{\mathbf{xy}}$, $\hat{\beta}_0 = \bar{y} - \bar{\mathbf{x}}^\top \hat{\mathbf{b}}$ 。

其它更直接 (但更复杂) 的证明:

证2. (应用分块矩阵的逆, 求 $\hat{\boldsymbol{\beta}}$ 的分量) $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon} = (\mathbf{1}, Z) \begin{pmatrix} \beta_0 \\ \mathbf{b} \end{pmatrix} + \boldsymbol{\varepsilon}$

由定理1, 已知LS估计 $\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y}$, 下面求分量 $\hat{\mathbf{b}}$ 。

$$\text{分块: } X^T X = \begin{pmatrix} \mathbf{1}^T \\ Z^T \end{pmatrix} (\mathbf{1}, Z) = \begin{pmatrix} \mathbf{1}^T \mathbf{1} & \mathbf{1}^T Z \\ Z^T \mathbf{1} & Z^T Z \end{pmatrix} = \begin{pmatrix} n & n\bar{\mathbf{x}}^T \\ n\bar{\mathbf{x}} & \sum \mathbf{x}_i \mathbf{x}_i^T \end{pmatrix}, X^T \mathbf{y} = \begin{pmatrix} \mathbf{1}^T \\ Z^T \end{pmatrix} \mathbf{y} = \begin{pmatrix} n\bar{y} \\ \sum \mathbf{x}_i y_i \end{pmatrix},$$

由分块矩阵求逆公式 (第3讲命题3)

$$(X^T X)^{-1} = \begin{pmatrix} a & -\bar{\mathbf{x}}^T A^{-1} \\ -A^{-1} \bar{\mathbf{x}} & A^{-1} \end{pmatrix}, \text{其中 } A = (n-1)S_{\mathbf{xx}}, \quad a = 1/n + \bar{\mathbf{x}}^T A^{-1} \bar{\mathbf{x}}$$

$$\Rightarrow \hat{\boldsymbol{\beta}} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\mathbf{b}} \end{pmatrix} = \begin{pmatrix} n & n\bar{\mathbf{x}}^T \\ n\bar{\mathbf{x}} & \sum \mathbf{x}_i \mathbf{x}_i^T \end{pmatrix}^{-1} \begin{pmatrix} n\bar{y} \\ \sum \mathbf{x}_i y_i \end{pmatrix} = \begin{pmatrix} a & -\bar{\mathbf{x}}^T A^{-1} \\ -A^{-1} \bar{\mathbf{x}} & A^{-1} \end{pmatrix} \begin{pmatrix} n\bar{y} \\ \sum \mathbf{x}_i y_i \end{pmatrix}$$

$$= \begin{pmatrix} na\bar{y} - \bar{\mathbf{x}}^T A^{-1} \sum \mathbf{x}_i y_i \\ A^{-1} \sum \mathbf{x}_i y_i - nA^{-1} \bar{\mathbf{x}} \bar{y} \end{pmatrix} = \begin{pmatrix} \bar{y} - \bar{\mathbf{x}}^T S_{\mathbf{xx}}^{-1} S_{\mathbf{xy}} \\ S_{\mathbf{xx}}^{-1} S_{\mathbf{xy}} \end{pmatrix}$$

所以 $\hat{\mathbf{b}} = S_{\mathbf{xx}}^{-1} S_{\mathbf{xy}}$, $\hat{\beta}_0 = \bar{y} - \bar{\mathbf{x}}^T \hat{\mathbf{b}}$.

第三讲命题3: $\Sigma > 0$,

$$\begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}^{-1} = \begin{pmatrix} \Sigma_{11\bullet 2}^{-1} & -\Sigma_{11\bullet 2}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \\ -\Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11\bullet 2}^{-1} & \Sigma_{22\bullet 1}^{-1} \end{pmatrix}$$

其中各个子块有不同的表达:

$$\Sigma_{11\bullet 2}^{-1} = \Sigma_{11}^{-1} + \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22\bullet 1}^{-1} \Sigma_{21} \Sigma_{11}^{-1},$$

$$\Sigma_{11\bullet 2}^{-1} \Sigma_{12} \Sigma_{22}^{-1} = \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22\bullet 2}^{-1} \text{ 等等}$$

证3 (误差平方和分别对 β_0 和 \mathbf{b} 求导) $\boldsymbol{\beta} = (\beta_0, \mathbf{b}^\top)^\top$

$$Q(\boldsymbol{\beta}) = \|\mathbf{y} - X\boldsymbol{\beta}\|^2 = \sum_{i=1}^n (y_i - \tilde{\mathbf{x}}_i^\top \boldsymbol{\beta})^2 = \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \mathbf{b})^2$$

$$\frac{\partial Q}{\partial \boldsymbol{\beta}} = 0 \Leftrightarrow \begin{cases} \frac{\partial Q}{\partial \beta_0} = \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \mathbf{b}) = 0 \\ \frac{\partial Q}{\partial \mathbf{b}} = \sum_{i=1}^n \mathbf{x}_i (y_i - \beta_0 - \mathbf{x}_i^\top \mathbf{b}) = 0 \end{cases} \Leftrightarrow \begin{cases} \sum y_i - n\beta_0 - (\sum \mathbf{x}_i)^\top \mathbf{b} = 0 \\ \sum \mathbf{x}_i y_i - (\sum \mathbf{x}_i) \beta_0 - (\sum \mathbf{x}_i \mathbf{x}_i^\top) \mathbf{b} = 0 \end{cases}$$

$$\Leftrightarrow \begin{cases} \bar{y} - \beta_0 - \bar{\mathbf{x}}^\top \mathbf{b} = 0 \\ \sum \mathbf{x}_i y_i - n\bar{\mathbf{x}} \beta_0 - (\sum \mathbf{x}_i \mathbf{x}_i^\top) \mathbf{b} = 0 \end{cases}$$

第一式左乘 $n\bar{\mathbf{x}}$ 减去第二式, 消去 β_0 , 得:

消去 β_0 的过程 \Leftrightarrow 中心化

$$(\sum \mathbf{x}_i \mathbf{x}_i^\top - n\bar{\mathbf{x}} \bar{\mathbf{x}}^\top) \mathbf{b} = \sum \mathbf{x}_i y_i - n\bar{\mathbf{x}} \bar{y}$$

$$\Rightarrow \hat{\mathbf{b}} = \left(\sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top \right)^{-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) y_i = S_{\mathbf{xx}}^{-1} S_{\mathbf{xy}}, \hat{\beta}_0 = \bar{y} - \bar{\mathbf{x}}^\top \hat{\mathbf{b}},$$

关于“不妨假设数据已经中心化”

我们已知：

- 如果 $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon} = \mathbf{1}\beta_0 + Z\mathbf{b} + \boldsymbol{\varepsilon}$ 中 Z 已经中心化，那么 $\hat{\mathbf{b}} = (Z^\top Z)^{-1} Z^\top \mathbf{y}$ ， $\hat{\beta}_0 = \bar{y}$
- 如果 \mathbf{y} ， X 都已经中心化，那么 $\hat{\mathbf{b}} = (Z^\top Z)^{-1} Z^\top \mathbf{y}$ ， $\hat{\beta}_0 = 0$ ，此时 β_0 无需估计。

在一些问题中，首先假设 Z 中心化，会简化后续论证过程。

- $\mathbf{y} = Z\mathbf{b} + \boldsymbol{\varepsilon}$ （不含截距项）， \mathbf{y} ， Z 都已经中心化

$$\text{LS估计 } \hat{\mathbf{b}} = (Z^\top Z)^{-1} Z^\top \mathbf{y} = S_{xx}^{-1} S_{xy}$$

- $\mathbf{y} = Z\mathbf{b} + \boldsymbol{\varepsilon}$ （不含截距项）， \mathbf{y} ， Z 都已经标准化

$$\text{LS估计 } \hat{\mathbf{b}} = (Z^\top Z)^{-1} Z^\top \mathbf{y} = R_{xx}^{-1} R_{xy}$$

Recap

我们已知，模型 $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 两端同时左乘 X^\top ，并令 $X^\top \boldsymbol{\varepsilon} = 0$

$$X^\top \mathbf{y} = X^\top X \boldsymbol{\beta}$$

即得 $\boldsymbol{\beta}$ 的LS估计 $\hat{\boldsymbol{\beta}} = (X^\top X)^{-1} X^\top \mathbf{y}$,

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \mathbf{b} \end{pmatrix}, \text{ LS估计 } \hat{\mathbf{b}} \text{ 的直观求解方法}$$

模型 $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon} = \mathbf{1}\beta_0 + Z\mathbf{b} + \boldsymbol{\varepsilon}$ 两边同时左乘 Z_c (Z_c 与 $\mathbf{1}$ 列正交, $Z_c^\top \mathbf{1} = 0$, 但与 Z 最接近), 并令 $Z_c^\top \boldsymbol{\varepsilon} = 0$:

$$Z_c^\top \mathbf{y} = Z_c^\top \mathbf{1}\beta_0 + Z_c^\top Z\mathbf{b} + Z_c^\top \boldsymbol{\varepsilon} = Z_c^\top Z\mathbf{b} = Z_c^\top Z_c \mathbf{b}$$

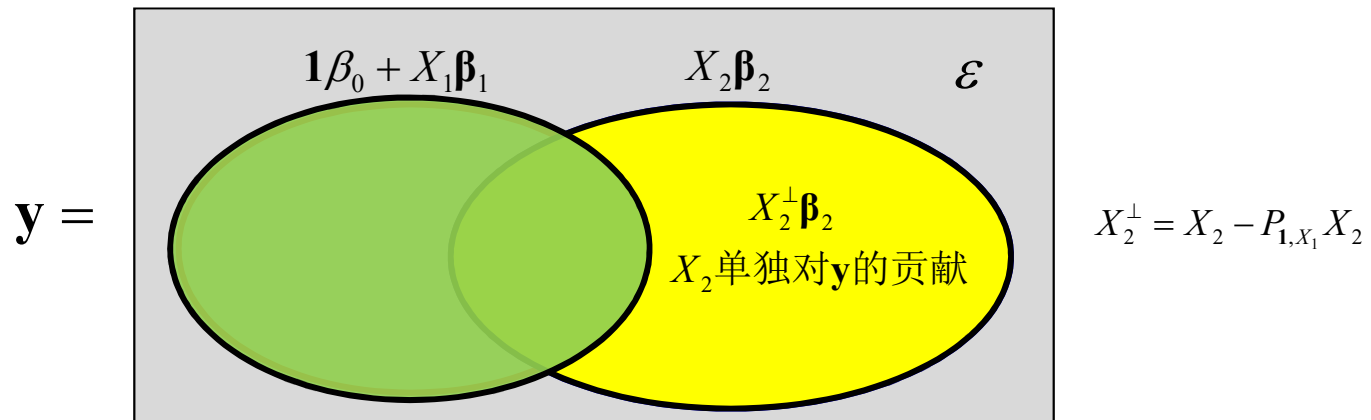
$$\Rightarrow \hat{\mathbf{b}} = (Z_c^\top Z_c)^{-1} Z_c^\top \mathbf{y}$$

LS估计的分量

假设设计阵 X 按列划分为 $X = (\mathbf{1}, X_1, X_2)$ ，模型为

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon} = \mathbf{1}\beta_0 + X_1\boldsymbol{\beta}_1 + X_2\boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$$

其中 X_1 的列是干扰因素， X_2 的列是感兴趣的变量。LS估计 $\hat{\boldsymbol{\beta}}_2$ 是否有效地消除了 X_1 的影响？



事实上LS估计 $\hat{\boldsymbol{\beta}}_2 = (X_2^{\perp\top} X_2^\perp)^{-1} X_2^{\perp\top} \mathbf{y} = (X_2^{\perp\top} X_2^\perp)^{-1} X_2^{\perp\top} \mathbf{y}^\perp$ ，
其中的 $X_2^\perp = X_2 - P_{1,X_1}X_2$ ， $\mathbf{y}^\perp = \mathbf{y} - P_{1,X_1}\mathbf{y}$ 都消除了 X_1 的影响。

记号：记所有数据 (\mathbf{y}, X_1, X_2) 的样本方差-协方差矩阵为

$$S = \begin{matrix} & \mathbf{y} & \mathbf{x}_1 & \mathbf{x}_2 \\ \mathbf{y} & S_{yy} & S_{y1} & S_{y2} \\ \mathbf{x}_1 & S_{1y} & S_{11} & S_{12} \\ \mathbf{x}_2 & S_{2y} & S_{21} & S_{22} \end{matrix}, \quad \begin{pmatrix} S_{yy \cdot 1} & S_{y2 \cdot 1} \\ S_{2y \cdot 1} & S_{22 \cdot 1} \end{pmatrix} = \begin{pmatrix} S_{yy} & S_{y2} \\ S_{2y} & S_{22} \end{pmatrix} - \begin{pmatrix} S_{y1} \\ S_{21} \end{pmatrix} S_{11}^{-1} (S_{1y}, S_{12})$$

LS估计的自相似性

命题3：假设线性模型： $\mathbf{y}_{n \times 1} = X_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1} = \mathbf{1} \beta_0 + X_1 \boldsymbol{\beta}_1 + X_2 \boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$

其中 $X = (\mathbf{1}, X_1, X_2)$, $\boldsymbol{\beta} = (\beta_0, \boldsymbol{\beta}_1^\top, \boldsymbol{\beta}_2^\top)^\top$, 记 $X_2^\perp = X_2 - P_{1, X_1} X_2$, 则

(1) $\hat{\boldsymbol{\beta}}_2 = (X_2^{\perp \top} X_2^\perp)^{-1} X_2^{\perp \top} \mathbf{y}$,

$$\text{var}(\hat{\boldsymbol{\beta}}_2 | X) = \sigma^2 (X_2^{\perp \top} X_2^\perp)^{-1}.$$

(注意： $\hat{\boldsymbol{\beta}}_2$ 也等于 $(X_2^{\perp \top} X_2^\perp)^{-1} X_2^{\perp \top} \mathbf{y}^\perp$, 其中 $\mathbf{y}^\perp = \mathbf{y} - P_{1, X_1} \mathbf{y}$)

(2) $\hat{\boldsymbol{\beta}}_2$ 也可表示为 $\hat{\boldsymbol{\beta}}_2 = S_{22 \cdot 1}^{-1} S_{2y \cdot 1}$,

$$\text{var}(\hat{\boldsymbol{\beta}}_2 | X) = \sigma^2 S_{22 \cdot 1}^{-1} / (n-1).$$

对比第5讲命题3(总体版本)

$$\boldsymbol{\beta}_2 = \Sigma_{22 \cdot 1}^{-1} \Sigma_{2y \cdot 1}$$

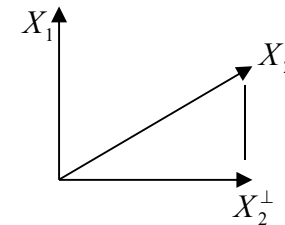
LS估计 $\hat{\beta}_2$ 的直观求解

模型 $\mathbf{y} = X_1\beta_1 + X_2\beta_2 + \varepsilon$ 两端同时左乘 $X_2^{\perp\top}$, 并令 $X_2^{\perp\top}\varepsilon = \mathbf{0}$ (矩估计方法), 得

$$X_2^{\perp\top}\mathbf{y} = X_2^{\perp\top}X_2\beta_2 \Rightarrow \hat{\beta}_2 = (X_2^{\perp\top}X_2^{\perp})^{-1}X_2^{\perp\top}\mathbf{y}.$$

为什么方程两端同时左乘 $X_2^{\perp\top}$?

- (1) X_2^{\perp} 与 X_1 列正交 (X_2^{\perp} : X_2 中消除了与 X_1 有关的部分)。
- (2) 在与 X_1 列正交的矩阵中, X_2^{\perp} 与 X_2 最为接近 (LS, 参见右图)。
- (3) 如此得到的解是 β_2 的 LS 估计 (命题3), 参见下页例子。



两步回归求 $\hat{\beta}_2$

$\hat{\beta}_2 = (X_2^{\perp\top}X_2^{\perp})^{-1}X_2^{\perp\top}\mathbf{y}^{\perp}$ 可以看作是两步回归的结果:

- $X_2 \sim X_1$ (多元回归) 得残差: $X_2^{\perp} = X_2 - P_{1,X_1}X_2$,
- $\mathbf{y} \sim X_1$ 得残差: $\mathbf{y}^{\perp} = \mathbf{y} - P_{1,X_1}\mathbf{y}$,
- $\mathbf{y}^{\perp} \sim X_2^{\perp}$: $\mathbf{y} = X_2^{\perp}\beta_2 + \delta$ 得 $\hat{\beta}_2 = (X_2^{\perp\top}X_2^{\perp})^{-1}X_2^{\perp\top}\mathbf{y}^{\perp}$.

理解为矩阵 $X_2 = (\mathbf{x}_1, \dots, \mathbf{x}_k)$ 的各列对 X_1 回归:
 $\mathbf{x}_i \sim X_1 \Rightarrow$ 残差 $\hat{\mathbf{e}}_i$,
 $X_2^{\perp} = (\hat{\mathbf{e}}_1, \dots, \hat{\mathbf{e}}_k)$

例2. 模型方程模型 $\mathbf{y} = X_1\boldsymbol{\beta}_1 + X_2\boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$, 设 X_2 为 $n \times q$ 矩阵, 设 W 为任何 $n \times q$ 矩阵, $W^\top X_1 = 0, W^\top X_2$ 可逆, 模型两端同时左乘 W^\top , 并令 $W^\top \boldsymbol{\varepsilon} / n = \mathbf{0}$:

$$W^\top \mathbf{y} = W^\top X_1 + W^\top X_2 \boldsymbol{\beta}_2 + W^\top \boldsymbol{\varepsilon} = W^\top X_2 \boldsymbol{\beta}_2$$

$\Rightarrow \tilde{\boldsymbol{\beta}}_2 = (W^\top X_2)^{-1} W^\top \mathbf{y}$, 则 $\tilde{\boldsymbol{\beta}}_2$ 是无偏估计, 但其方差大于 $\hat{\boldsymbol{\beta}}_2$ 的方差。

证明:

$$W^\top X_1 = 0 \Rightarrow E(\tilde{\boldsymbol{\beta}}_2 | X) = (W^\top X_2)^{-1} W^\top (X_1 \boldsymbol{\beta}_1 + X_2 \boldsymbol{\beta}_2) = (W^\top X_2)^{-1} W^\top X_2 \boldsymbol{\beta}_2 = \boldsymbol{\beta}_2.$$

$$\text{为证 } \text{var}(\tilde{\boldsymbol{\beta}}_2 | X) = \sigma^2 (W^\top X_2)^{-1} W^\top W (X_2^\top W)^{-1} \geq \sigma^2 (X_2^{\perp \top} X_2^\perp)^{-1} = \text{var}(\hat{\boldsymbol{\beta}}_2 | X),$$

$$\text{只需证 } W^\top W \geq W^\top X_2 (X_2^{\perp \top} X_2^\perp)^{-1} X_2^\top W.$$

因为 $W^\top X_1 = 0$, 所以 $W^\top X_2 = W^\top X_2^\perp$, 所以上述不等式等价于

$$W^\top W \geq W^\top X_2^\perp (X_2^{\perp \top} X_2^\perp)^{-1} X_2^{\perp \top} W$$

因为 $X_2^\perp (X_2^{\perp \top} X_2^\perp)^{-1} X_2^{\perp \top} \leq I_n$, 该不等式显然成立。

命题3的证明

证明1(投影): (1) 不妨设 X_1, X_2 已经中心化, 所以 $\mathbf{1}, X_1, X_2^\perp$ 三部分相互列正交, 所以

$$\begin{aligned}\hat{\mathbf{y}} &= P_X \mathbf{y} = P_{(\mathbf{1}, X_1, X_2)} \mathbf{y} = P_{(\mathbf{1}, X_1, X_2^\perp)} \mathbf{y} = P_{\mathbf{1}} \mathbf{y} + P_{X_1} \mathbf{y} + P_{X_2^\perp} \mathbf{y} \\ &= \mathbf{1}\bar{y} + X_1 \left((X_1^\top X_1)^{-1} X_1^\top \mathbf{y} \right) + X_2^\perp \left((X_2^{\perp\top} X_2^\perp)^{-1} X_2^{\perp\top} \mathbf{y} \right) \\ &\stackrel{\Delta}{=} \mathbf{1}\bar{y} + X_1 \mathbf{u} + X_2^\perp \mathbf{v} = \mathbf{1}\bar{y} + X_1 \mathbf{u} + \left(X_2 - X_1 (X_1^\top X_1)^{-1} X_1^\top X_2 \right) \mathbf{v} \\ &= \mathbf{1}\bar{y} + X_1 \left(\mathbf{u} - (X_1^\top X_1)^{-1} X_1^\top X_2 \mathbf{v} \right) + X_2 \mathbf{v},\end{aligned}$$

$$X_2 \text{ 的系数 } \mathbf{v} \text{ 即 } \hat{\boldsymbol{\beta}}_2, \quad \hat{\boldsymbol{\beta}}_2 = \mathbf{v} = (X_2^{\perp\top} X_2^\perp)^{-1} X_2^{\perp\top} \mathbf{y} = \left(X_2^{\perp\top} X_2^\perp \right)^{-1} X_2^{\perp\top} \mathbf{y}$$

$$\begin{aligned}(2) \quad X_2^\perp &= X_2 - P_{X_1} X_2 \Rightarrow X_2^{\perp\top} X_2^\perp = X_2^\top (I_n - P_{X_1}) X_2 = X_2^\top X_2 - X_2^\top X_1 (X_1^\top X_1)^{-1} X_1^\top X_2 \\ &= (n-1)S_{22} - (n-1)S_{21} S_{11}^{-1} S_{12} = (n-1)S_{22 \bullet 1}, \quad \text{同样 } X_2^{\perp\top} \mathbf{y} = (n-1)S_{2y \bullet 1}\end{aligned}$$

$$\text{所以 } \hat{\boldsymbol{\beta}}_2 = (X_2^{\perp\top} X_2^\perp)^{-1} X_2^{\perp\top} \mathbf{y} = S_{22 \bullet 1}^{-1} S_{2y \bullet 1},$$

$$\text{var}(\hat{\boldsymbol{\beta}}_2 | X) = \sigma^2 (X_2^{\perp\top} X_2^\perp)^{-1} = \sigma^2 S_{22 \bullet 1}^{-1} / (n-1).$$

证明2: 不妨假设 X_1, X_2 已经中心化, 这不影响 β_1, β_2 的LS估计.

$X_2^\perp = X_2 - P_{1, X_1} X_2 = X_2 - P_{X_1} X_2$, 改写模型:

$$\begin{aligned} \mathbf{y} &= \mathbf{1}\beta_0 + X_1\beta_1 + X_2\beta_2 + \boldsymbol{\varepsilon} = \mathbf{1}\beta_0 + X_1\beta_1 + (X_2^\perp + P_{X_1} X_2)\beta_2 + \boldsymbol{\varepsilon} \\ &= \mathbf{1}\beta_0 + X_1 \left[\beta_1 + (X_1^\top X_1)^{-1} X_1^\top X_2 \beta_2 \right] + X_2^\perp \beta_2 + \boldsymbol{\varepsilon} \\ &\triangleq \mathbf{1}\beta_0 + X_1 \beta_1^* + X_2^\perp \beta_2 + \boldsymbol{\varepsilon}, \text{ 其中 } \beta_1^* = \beta_1 + (X_1^\top X_1)^{-1} X_1^\top X_2 \beta_2 \end{aligned}$$

该模型得设计阵的三部分 $\mathbf{1}, X_1, X_2^\perp$ 相互正交, $X^\top X$ 分块对角, 所以由LS估计公式

$$\text{易得 } \hat{\boldsymbol{\beta}} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1^* \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{1}^\top \mathbf{1} & 0 & 0 \\ 0 & X_1^\top X_1 & 0 \\ 0 & 0 & X_2^{\perp \top} X_2^\perp \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{1}^\top \mathbf{y} \\ X_1^\top \mathbf{y} \\ X_2^{\perp \top} \mathbf{y} \end{pmatrix} = \begin{pmatrix} \bar{y} \\ (X_1^\top X_1)^{-1} X_1^\top \mathbf{y} \\ (X_2^{\perp \top} X_2^\perp)^{-1} X_2^{\perp \top} \mathbf{y} \end{pmatrix}$$

所以 $\hat{\beta}_2 = (X_2^{\perp \top} X_2^\perp)^{-1} X_2^{\perp \top} \mathbf{y}$.

证明3: (不假设 X_1, X_2 已经中心化)

$Q = \|\mathbf{y} - \mathbf{1}\beta_0 - X_1\boldsymbol{\beta}_1 - X_2\boldsymbol{\beta}_2\|^2$ 求导得正则方程 $X^T(\mathbf{y} - X\boldsymbol{\beta}) = \mathbf{0}$:

$$\begin{cases} \partial Q / \partial \beta_0 = -\mathbf{1}^T(\mathbf{y} - \mathbf{1}\beta_0 - X_1\boldsymbol{\beta}_1 - X_2\boldsymbol{\beta}_2) = 0 \\ \partial Q / \partial \boldsymbol{\beta}_1 = -X_1^T(\mathbf{y} - \mathbf{1}\beta_0 - X_1\boldsymbol{\beta}_1 - X_2\boldsymbol{\beta}_2) = \mathbf{0} \\ \partial Q / \partial \boldsymbol{\beta}_2 = -X_2^T(\mathbf{y} - \mathbf{1}\beta_0 - X_1\boldsymbol{\beta}_1 - X_2\boldsymbol{\beta}_2) = \mathbf{0} \end{cases}$$

下面应用初等消去法解正则方程(即求 $X^T X$ 的逆)

第一式整理得: $\beta_0 = \bar{y} - \bar{\mathbf{x}}_1^T \boldsymbol{\beta}_1 - \bar{\mathbf{x}}_2^T \boldsymbol{\beta}_2$, 代入二、三式

$$\begin{cases} X_1^T [\mathbf{y} - \mathbf{1}(\bar{y} - \bar{\mathbf{x}}_1^T \boldsymbol{\beta}_1 - \bar{\mathbf{x}}_2^T \boldsymbol{\beta}_2) - X_1\boldsymbol{\beta}_1 - X_2\boldsymbol{\beta}_2] = 0 \\ X_2^T [\mathbf{y} - \mathbf{1}(\bar{y} - \bar{\mathbf{x}}_1^T \boldsymbol{\beta}_1 - \bar{\mathbf{x}}_2^T \boldsymbol{\beta}_2) - X_1\boldsymbol{\beta}_1 - X_2\boldsymbol{\beta}_2] = 0 \end{cases}$$

整理得:

$$\begin{cases} X_1^T (X_1 - \mathbf{1}\bar{\mathbf{x}}_1^T) \boldsymbol{\beta}_1 + X_1^T (X_2 - \mathbf{1}\bar{\mathbf{x}}_2^T) \boldsymbol{\beta}_2 = X_1^T (\mathbf{y} - \mathbf{1}\bar{y}) = (X_1 - \mathbf{1}\bar{\mathbf{x}}_1^T)^T \mathbf{y} \\ X_2^T (X_1 - \mathbf{1}\bar{\mathbf{x}}_1^T) \boldsymbol{\beta}_1 + X_2^T (X_2 - \mathbf{1}\bar{\mathbf{x}}_2^T) \boldsymbol{\beta}_2 = X_2^T (\mathbf{y} - \mathbf{1}\bar{y}) \end{cases}$$

记（中心化） $X_{1c} = X_1 - \mathbf{1}\bar{x}_1^\top$, $X_{2c} = X_2 - \mathbf{1}\bar{x}_2^\top$, $\mathbf{y}_c = \mathbf{y} - \mathbf{1}\bar{y}$,

前述方程组简写为

$$\begin{cases} X_{1c}^\top X_{1c} \boldsymbol{\beta}_1 + X_{1c}^\top X_{2c} \boldsymbol{\beta}_2 = X_{1c}^\top \mathbf{y} \\ X_{2c}^\top X_{1c} \boldsymbol{\beta}_1 + X_{2c}^\top X_{2c} \boldsymbol{\beta}_2 = X_{2c}^\top \mathbf{y} \end{cases}$$

第一式左乘 $X_{2c}^\top X_{1c} (X_{1c}^\top X_{1c})^{-1}$ 得:

$$X_{2c}^\top X_{1c} \boldsymbol{\beta}_1 + X_{2c}^\top X_{1c} (X_{1c}^\top X_{1c})^{-1} X_{1c}^\top X_{2c} \boldsymbol{\beta}_2 = X_{2c}^\top X_{1c} (X_{1c}^\top X_{1c})^{-1} X_{1c}^\top \mathbf{y}$$

即 $X_{2c}^\top X_{1c} \boldsymbol{\beta}_1 + X_{2c}^\top P_{X_{1c}} X_{2c} \boldsymbol{\beta}_2 = X_{2c}^\top P_{X_{1c}} \mathbf{y}$, 与第二式相减(消去 $\boldsymbol{\beta}_1$)得:

$$X_{2c}^\top (I - P_{X_{1c}}) X_{2c} \boldsymbol{\beta}_2 = X_{2c}^\top (I - P_{X_{1c}}) \mathbf{y},$$

$$\Leftrightarrow X_2^{\perp\top} X_2^{\perp} \boldsymbol{\beta}_2 = X_2^{\perp\top} \mathbf{y} \Rightarrow \hat{\boldsymbol{\beta}}_2 = \left(X_2^{\perp\top} X_2^{\perp} \right)^{-1} X_2^{\perp\top} \mathbf{y}$$

证明4: 应用分块矩阵的逆 (略)