

线性回归 :  $X \in \mathbb{R}^{n \times p}$   $n$  个样本,  $p$  个特征

low dimension :  $\begin{cases} \text{univariate} : p=1 & \text{单变量} \\ \text{multivariate} : 1 < p < n & \text{多变量} \end{cases}$  } 低维

high dimension :  $p \geq n$  高维

## 1. Form of Multivariate Linear Models

① 假设有  $n$  组样本 :  $\{(y_i, x_i) \mid i=1, 2, \dots, n\}$  其中标量  $y_i \in \mathbb{R}$  为 response 向量  $x_i = [x_{i1}, x_{i2}, \dots, x_{ip}]^T \in \mathbb{R}^p$  包含  $p$  个 variate,  $y_i$  与  $x_i$  满足如下  
的 Linear Model

向量形式  $y_i = x_i^T \beta + \varepsilon_i = \beta_1 + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \varepsilon_i \quad i=1, 2, \dots, n \quad (1)$

其中  $\beta = [\beta_1, \beta_2, \dots, \beta_p]^T \in \mathbb{R}^p$  是 regression parameters (包含截距 intercept 因为  $x_{i1} = 1$ ),  $\varepsilon_i$  为第  $i$  组观测样本的 random error.

称 (1) 为 vector form of the linear model

② 记 design matrix  $X = [x_1^T, x_2^T, \dots, x_n^T]^T \in \mathbb{R}^{n \times p}$

response vector  $y = [y_1, y_2, \dots, y_n]^T \in \mathbb{R}^n$

random error vector  $\varepsilon = [\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n]^T \in \mathbb{R}^n$

regression parameters  $\beta = [\beta_1, \beta_2, \dots, \beta_p]^T \in \mathbb{R}^p$

则有 Matrix form :

矩阵形式  $y = X \cdot \beta + \varepsilon \quad (2)$

## 2. Basic assumptions of multivariate linear models 基本假设

① Assumption 1: fixed designed: designed matrix  $X$  is fixed / nonrandom

固定设计: 设计矩阵  $X$  是固定非随机的

② Assumption 2: full rankness of the design matrix:  $\text{rank}(X) = p < n$

设计矩阵满秩:  $\text{rank}(X) = p < n$  即  $X$  的列向量线性独立

Tips: 这一假设是为了保证可识别性 (identifiability), 例如: 若  $x_{i2} = k x_{i3}$

则只能 estimate/learn 到  $(\beta_2 + \beta_3)$ , 具体的  $\beta_2, \beta_3$  无法得到

Tips: 但只对模型的选择、估计、推断是必要的, 对预测不必要

③ Assumption 3: Gauss - Markov condition of the random errors

随机误差满足 Gauss - Markov 条件:  $E(\varepsilon) = 0$   $\text{cov}(\varepsilon) = \sigma^2 I$

$E(\varepsilon) = 0 \Rightarrow y$  的期望是  $X$  的线性函数

$\text{cov}(\varepsilon) = \sigma^2 I \Rightarrow$   $n$  个样本是不相关的 (未必独立), 且等方差

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

(3)

## 3. Interpretation of Regression Parameters 参数的含义

In formula (1),

$$\frac{\partial E(y_i)}{\partial x_{ij}} = \beta_j$$

表示在其他自变量保持不变时, 自变量  $x_{ij}$  每增加 1 单位,

因变量  $y_i$  的平均增幅

## 4. Estimation of regression parameters 估计

### 4.1. Ordinary Least Squares Estimation (OLSE) 最小二乘

① 拟合值: fitted values

$$\{\hat{y}_i = x_i^T \hat{\beta} \mid i=1, 2, \dots, n\} \subset \mathbb{R}$$

② 残差: residuals

$$\{e_i = y_i - \hat{y}_i \mid i=1, 2, \dots, n\} \subset \mathbb{R}$$

③ 估计参数  $\hat{\beta}$

$$\hat{\beta} = \arg \min_{b \in \mathbb{R}^p} \sum_{i=1}^n (y_i - x_i^T b)^2 = \arg \min_{b \in \mathbb{R}^p} \|y - Xb\|^2 \quad (4)$$

④ 求导即可:



$$\sum x_i e_i = 0 \text{ 而 } x_{i1}=1 \Rightarrow \sum e_i = 0$$

$$\mathbb{R}^p \ni 0 = \frac{\partial Q}{\partial b} \Big|_{b=\hat{\beta}} = \sum_{i=1}^n x_i e_i = \sum_{i=1}^n x_i (y_i - x_i^T \hat{\beta}) = X^T (y - X\hat{\beta}) \quad (5)$$

$$e_i \in \mathbb{R}$$

$$x_i \in \mathbb{R}^p$$

$$\hat{\beta} = (X^T X)^{-1} \cdot X^T \cdot y$$

(6)

$$X \in \mathbb{R}^{n \times p} \quad y \in \mathbb{R}^n \quad \hat{\beta} \in \mathbb{R}^p$$

Proof:  $X \in \mathbb{R}^{n \times p}$  当  $\text{rank}(X) = p$  有  $X^T X$  可逆

Lemma 对  $\forall$  矩阵  $A \in \mathbb{R}^{m \times n}$  有  $A^T A x = 0 \Leftrightarrow Ax = 0$

①  $\Leftarrow$  显然,

②  $\Rightarrow A^T A x = 0$  有  $x^T A^T A x = 0$

$$\text{即 } (Ax)^T (Ax) = 0 \text{ 即 } \|Ax\|_2 = 0$$

$$\text{则 } Ax = 0 \quad \square$$

回到原命题:

$$\text{记 } N(X) = \{v \mid Xv = 0\} \quad N(X^T X) = \{v \mid X^T X v = 0\}$$

由 Lemma, 二者相同, 即有:  $\dim(N(X)) = \dim(N(X^T X))$

$$\text{又} \because p = \text{rank}(X) + \dim(N(X)) \quad \text{rank}(X) = p$$

$$\therefore \dim(N(X^T X)) = \dim(N(X)) = 0$$

故  $\text{rank}(X^T X) = p$  加之  $X^T X \in \mathbb{R}^{p \times p}$

则  $X^T X$  满秩 (可逆)  $\square$

Supplementary:  $\frac{\partial \|y - Xb\|^2}{\partial b} = \frac{\partial}{\partial b} \{(y - Xb)^T (y - Xb)\}$

$$= -X^T \cdot (y - Xb) - (y - Xb)^T X = -X^T y + X^T X b - (y^T X)^T + X^T X b$$

$$= -2X^T (y - Xb)$$

## 4.2 Maximum Likelihood Estimation (MLE) 极大似然,

④ Assumption 4: random error vector follows a multivariate normal

$$\varepsilon \sim N(0, \sigma^2 I_n)$$

or  $y \sim N(X\beta, \sigma^2 I_n)$

(7)

①  $\{y_i \mid i=1, 2, \dots, n\}$  相互独立

② 对数似然函数

$$\log L(\beta, \sigma^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \|y - X\beta\|^2$$

(8)

③ 关于  $\beta$  最大化, 等价于 (4), 即

在  $\varepsilon \sim N(0, \sigma^2 I)$  的假设下,  $\beta$  的 MLE 与 OLSE 相同

$$\hat{\beta} = (X^T X)^{-1} \cdot X^T \cdot y$$

(6)

④ 关于  $\sigma^2$  最大化

$$\hat{\sigma}_{MLE}^2 = n^{-1} \|y - X\hat{\beta}\|^2 = n^{-1} \|e\|^2$$

## 5. Properties of estimates and residuals 性质

5.1  $\{e_i = y_i - x_i^T \hat{\beta}\}$  残差的性质

fitted value of  $y$ :

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

记  $H = X(X^T X)^{-1} X^T$ , 则有  $\hat{y} = H \cdot y$

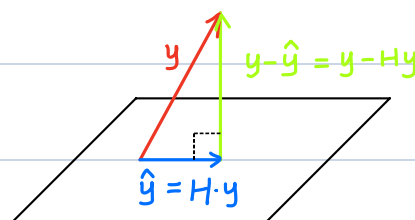
实际上,  $H$  是一个投影矩阵, 将  $y$  投影到设计矩阵  $X$  的列空间

①  $H$  的性质 (投影矩阵)

(1)  $H = X(X^T X)^{-1} X^T \in \mathbb{R}^{n \times n}$

(2)  $\sum_{i=1}^n H_{[i,i]} = P \leftarrow \text{tr}(AB) = \text{tr}(BA)$  易证

✓ (3)  $H^2 = H \quad H^T = H$



✓ (4)  $(I-H)^2 = I-H \Rightarrow I-H$  特征值为 0 或 1

✓ (5)  $(I-H)X = 0$   $\text{rank}(I-H) = \sum \lambda_i = \text{tr}(I-H) = n-p$

$\text{tr}(I-H) = n-p$

## ② 残差 $e_i$ 的性质

(1)  $e_i = y_i - \hat{y}_i$

(2)  $E(e_i) = 0$

$\text{cov}(e) = \text{cov}(y - \hat{y}) = \text{cov}((I-H)y)$

(3)  $E(e_i^2) = \text{var}(e_i) = (1 - H_{[i,i]}) \cdot \sigma^2$   $\begin{aligned} &= (I-H) \sigma^2 I (I-H)^T \\ &= \sigma^2 (I-H) \end{aligned}$

(4)  $E(\sum_{i=1}^n e_i^2) = (n - \sum_{i=1}^n H_{[i,i]}) \sigma^2 = (n-p) \sigma^2$

表明:  $\hat{\sigma}^2 = (n-p)^{-1} \cdot \sum_{i=1}^n e_i^2$  是  $\sigma^2$  的无偏估计 ★

注: 模型建立时, 截距项作为  $x_{i1}=1$  的参数, 故与其他书本略有差别

## 5.2 回归参数的性质

(a). 线性 Linearity:  $\hat{\beta}$  是  $y$  的线性函数

$$\hat{\beta} = (X^T X)^{-1} X^T \cdot y$$

$E(y) = X\beta$

(b) 无偏性 Unbiasedness:  $\hat{\beta}$  是  $\beta$  的无偏估计

$$E(\hat{\beta}) = \beta$$

$\Rightarrow \hat{y} = X\hat{\beta}$  是  $E(y)$  的无偏估计

$$E(\hat{y}) = E(y)$$

## (c) 协方差矩阵 Covariance matrices

$$\text{cov}(\hat{\beta}) = (X^T X)^{-1} X^T \cdot \text{cov}(y) \cdot X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}$$

$$\text{cov}(e) = \sigma^2 (I - H)$$

☆

见上

(d) Gauss-Markov 定理:

在 Assumption 3 条件下, 向量  $\beta$  的任意线性组合  $v^T \beta$  的 最小方差 线性无偏估计 (Best Linear Unbiased Estimator, BLUE) 是  $v^T \hat{\beta}$ , 其中  $v \in \mathbb{R}^p$

$\hat{\beta}$  是  $\beta$  的 OLSE (见作业)

在 Assumption 4 正态假设也成立时,  $v^T \hat{\beta}$  作为 MLE 是  $v^T \beta$  的 最小方差 无偏估计 (Best Unbiased Estimator, BUE), 即  $v^T \beta$  所有 (线性或非线性) 无偏估计中方差最小的

(e).  $\sigma^2$  的无偏估计 Unbiased estimate for  $\sigma^2$

$\frac{e^T e}{n-p}$  是  $\sigma^2$  的无偏估计

见前

$$E\left(\frac{e^T e}{n-p}\right) = \sigma^2$$

Lemma:  $\text{tr}(H) = p$

(f). 与残差的不相关性:

$$e = y - \hat{y} = y - X(X^T X)^{-1} X^T y = (I - H) y$$

☆  
↓

$$\text{cov}(\hat{\beta}, e) = 0$$

$$e = (I - H) y$$

Note: 当正态假设 Assumption 4 成立时,  $\hat{\beta} \perp e$

19) 分布: 当正态假设 Assumption 4 成立时

$$\hat{\beta} \sim N(\beta, \sigma^2 (X^T X)^{-1})$$

$$\frac{e^T e}{\sigma^2} \sim \chi^2(n-p)$$

$$y \sim N(X\beta, \sigma^2 I)$$

$$e^T e = y^T (I-H) y$$

$$(X\beta)^T (I-H) (X\beta) = 0 \quad (E-H)^2 = I-H$$

$$X \sim N(\mu, I) \quad \text{rank}(A)=r \quad \mu^T A \mu = 0 \quad A^2 = A \quad A \text{ 是实对称}$$

$$X^T A X \sim \chi^2(r)$$

$$\text{rank}(I-H) = n-p$$

在 Assumption 3 条件下, 向量  $\beta$  的任意线性组合  $v^T \beta$  的 最小方差 线性

无偏估计 (Best Linear Unbiased Estimator, BLUE) 是  $v^T \hat{\beta}$ , 其中  $v \in \mathbb{R}^p$

$\hat{\beta}$  是  $\beta$  的 OLSE

$v^T \hat{\beta}$  的方差最小

4. 证明课件第五章第 5 节的性质(d).

设  $\tilde{\beta}$  为  $\beta$  的无偏估计量, 且存在  $D \in \mathbb{R}^{p \times n}$  以及  $B = (X^T X)^{-1} X^T + D$ , 使得  $\tilde{\beta} = By$ . 由于  $\tilde{\beta}$  的无偏性, 对任意  $\beta \in \mathbb{R}^p$ , 都有

$$\beta = E(\tilde{\beta}) = E(By) = (X^T X)^{-1} X^T X \beta + DX \beta = \beta + DX \beta.$$

( $\because v\beta$  成立)

$\therefore$  这表明  $DX = 0$ . 我们进而有

$$\begin{aligned} \text{cov}(\tilde{\beta}) &= B \text{cov}(y) B^T = \sigma^2 \{ (X^T X)^{-1} X^T + D \} \{ X (X^T X)^{-1} + D^T \} \\ &= \sigma^2 \{ (X^T X)^{-1} + \underline{DX} (X^T X)^{-1} + (X^T X)^{-1} (\underline{DX})^T + DD^T \} \\ &= \sigma^2 \{ (X^T X)^{-1} + DD^T \}. \end{aligned}$$

$$\text{故 } \text{cov}(\tilde{\beta}) - \text{cov}(\hat{\beta}) = \sigma^2 DD^T, \text{ 而 } \text{var}(v^T \tilde{\beta}) - \text{var}(v^T \hat{\beta}) = \sigma^2 v^T DD^T v = \sigma^2 \|D^T v\|^2 \geq 0.$$

$$\begin{aligned} E(e^T e) &= E\{ y^T (I-H)^T (I-H) y \} = E\{ y^T (I-H) y \} \\ &= E(y)^T (I-H) E(y) + \text{tr}\{ (I-H) \text{cov}(y) \} \end{aligned}$$





$$= \beta^T X^T (I - H) X \beta + \sigma^2 \text{tr}(I - H)$$

法2

$$= \sigma^2 \text{tr}(I - H)$$

$$= \sigma^2 (\text{tr}(I) - \text{tr}(H))$$

$$= \sigma^2 (n - p) \quad \square$$

or  $E(e^T e) = E(\sum e_i^2)$

$$= \sum [E(e_i^2) - (E(e_i))^2] \quad (\because E(e_i) = 0)$$

法3

$$= \sum \text{var}(e_i)$$

$$= \text{tr}[\text{cov}(e, e)]$$

$$= \text{tr}[\sigma^2(I - H)]$$

$$= \sigma^2 (n - p) \quad \square$$