

## 多重共线性 Multicollinearity

假设 Designed Matrix  $X$  "不再" 满秩：

$$\text{rank}(X) \approx p \quad X \in \mathbb{R}^{n \times p}$$

Assumption: ① 假设 Model 均满足 Gauss-Markov 条件

② 假设 已中心化，即

$\{y, x_1, \dots, x_p\}$  中的向量的 分量 均值为 0, 方差为 1

$$\frac{1}{n} \sum_{j=1}^n y_j = 0 \quad \frac{1}{n} \sum_{j=1}^n (y_j - 0)^2 = 1$$

$$\frac{1}{n} \sum_{j=1}^n x_{ij} = 0 \quad \frac{1}{n} \sum_{j=1}^n (x_{ij} - 0)^2 = 1$$

③  $\beta$  不含截距项  $\beta \in \mathbb{R}^p$

### 一、定义 Definition

#### 1. 完全多重共线性

若存在 非 0 向量  $v \in \mathbb{R}^p$ , 使设计矩阵  $X$  满足  $Xv = 0$ , 即  $X$  列向量 线性相关, 称为  $(x_1, \dots, x_p)$  完全多重共线性 Perfect Multicollinearity

#### 2. 多重共线性

若存在 非 0 向量  $v \in \mathbb{R}^p$ , 使设计矩阵  $X$  满足  $Xv \approx 0$ , 即  $X$  列向量 线性相关, 称为  $(x_1, \dots, x_p)$  多重共线性 Multicollinearity

## 二、影响 Effect

### 1. 完全多重共线性

当自变量  $\{\lambda_{(1)}, \lambda_{(2)}, \dots, \lambda_{(p)}\}$  存在完全多重共线性时，设计矩阵  $X$  不满秩， $\text{rank}(X) < p$  则  $\text{rank}(X^T X) = \text{rank}(X) < p$ ，即有  $X^T X$  不可逆，之后的 OLS 估计不可进行

### 2. 多重共线性

当自变量  $\{\lambda_{(1)}, \lambda_{(2)}, \dots, \lambda_{(p)}\}$  存在多重共线性时，质量差，因为方差  $\text{var}(v^T \hat{\beta})$  的上界趋向  $\infty$  (OLS: 无偏)  
 $\exists \lambda_{\min}\{(X^T X)\} \rightarrow 0$

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

$$\begin{aligned} \text{var}(v^T \hat{\beta}) &= \text{cov}(v^T \hat{\beta}, v^T \hat{\beta}) \\ &= v^T \text{cov}(\hat{\beta}) v \\ &= v^T \sigma^2 (X^T X)^{-1} v \\ &\leq \sigma^2 \cdot \|v\|^2 \cdot \lambda_{\max}\{(X^T X)^{-1}\} \\ &= \sigma^2 \cdot \|v\|^2 \cdot [\lambda_{\min}\{(X^T X)\}]^{-1} \end{aligned}$$

$X^T X$  存在接近 0 的特征根，故  $\frac{1}{\lambda_{\min}\{(X^T X)\}} \rightarrow \infty$

**Proof ②**  $\hat{\beta}$  各分量方差加和  $\text{tr}\{\sigma^2 (X^T X)^{-1}\} = \sigma^2 \cdot \sum_{j=1}^p \lambda_j^{-1} \rightarrow \infty$

### 三、诊断 Diagnoses

#### 1. 方差扩大因子法 Variance inflation factor ( VIF )

因为已经标准化，故  $\text{corr}(\hat{\beta}) = \sigma^2 (X^T X)^{-1} = (X^T X)^{-1}$ ，记

$$C = (X^T X)^{-1}$$

称  $C$  的第  $j$  个对角元素为：

$$VIF_j = C_{[j,j]}$$

称  $VIF_j$  为自变量  $X_{(j)}$  的 方差扩大因子，由于  $\text{cov}(\hat{\beta}) = \sigma^2 C$ ，故

$$\text{var}(\hat{\beta}_{[j]}) = \sigma^2 VIF_j$$

衡量标准：

①  $VIF_j$  衡量  $\hat{\beta}_j$  方差因多重共线性而扩大的程度

② 经验判断：  $VIF_j \geq 10$ ，认为  $X_{(j)}$  与其他有明显多重共线性

③ 检验整体：  $\frac{1}{P} \sum_{j=1}^P VIF_j >> 1$ ，认为整体有明显多重共线性

#### 2. 特征根判定法： Eigenvalue-based

三 多重共线性  $\Leftrightarrow \lambda_{\min}\{X^T X\}$  接近 0

经验判断：

考虑  $X^T X$  的特征根  $\lambda$ ：

①  $0 < \frac{\lambda_{\max}}{\lambda_{\min}} < 100$  不存在多重共线性

②  $100 < \frac{\lambda_{\max}}{\lambda_{\min}} < 10000$  存在较强多重共线性

③  $\frac{\lambda_{\max}}{\lambda_{\min}} > 10000$  存在严重多重共线性

### 3. 直观判断 Intuitive judgement

- 增加或剔除某个自变量，或者改变某个观测值时，回归系数的估计发生较大变化；
- 定性分析认定的一些重要自变量在拟合结果中没有通过显著性检验；
- 某些自变量回归系数估计值的正负号与定性分析结果不符；
- 自变量间的相关系数较大；
- 某些重要自变量回归系数估计值的标准误差较大。