



Step 2:  $\frac{\partial F}{\partial \beta} = 0 \Rightarrow \hat{\beta} = (X^T X)^{-1} X^T Y = A^{-1} Y \Rightarrow \hat{\beta} = (A^T A)^{-1} A^T Y = (A^T A)^{-1} A^T (A(A^T A)^{-1} A^T) Y = (A^T A)^{-1} A^T Y$

定理: 对上述模型, 满足  $A\hat{\beta} = \hat{y}$  且OLS估计为  $\hat{\beta} = (X^T X)^{-1} X^T Y = A^{-1} Y$  —— 计算时不要直接代入公式, 用特殊矩阵所算

**Example (1.3.1)**

【例3.3.1】在天文测量中, 对天空中三个星位点构成的三角形ABC的三个内角  $\theta_1, \theta_2, \theta_3$  进行测量, 得到的测量值分别为  $y_1, y_2, y_3$ . 由于存在测量误差, 所以需要同时对  $\theta_1, \theta_2, \theta_3$  进行估计, 我们用线性模型表示有关的量:

$$\begin{cases} y_1 = \theta_1 + e_1 \\ y_2 = \theta_2 + e_2 \\ y_3 = \theta_3 + e_3 \\ \theta_1 + \theta_2 + \theta_3 = \pi \end{cases}$$

其中  $e_i, i=1, 2, 3$  表示测量误差. 求  $\theta_1, \theta_2, \theta_3$  的最小二乘估计  $\hat{\theta}_1, \hat{\theta}_2, \hat{\theta}_3$ , 并求出当观测值  $y_1 = \frac{3}{2}\pi, y_2 = \frac{1}{2}\pi, y_3 = \frac{1}{2}\pi$  时  $\hat{\theta}_1, \hat{\theta}_2, \hat{\theta}_3$  的具体数值.

• 问题等价于在约束条件  $1^T \theta = \pi$  下求解如下模型参数的约束OLS估计:

$$\hat{\theta} = \hat{\theta} + \hat{e}$$

• 构造辅助函数

$$F(\hat{\theta}, \lambda) = \|\hat{y} - \hat{\theta}\|^2 + 2\lambda(1^T \hat{\theta} - \pi)$$

• 求F的临界点, 求解方程

$$\frac{\partial F}{\partial \hat{\theta}} = 2(\hat{y} - \hat{\theta}) = 0$$

得  $\hat{\theta} = \hat{y} - \lambda 1$  (5)

• 代入(5), 由定理1.3.1即可得

$$\hat{\theta}_{con} = \hat{y} - (\frac{\pi}{1^T \hat{y}}) 1$$
 (6)

• 将观测值  $y_1 = \frac{3}{2}\pi, y_2 = \frac{1}{2}\pi, y_3 = \frac{1}{2}\pi$  代入(6), 可得

$$\hat{\theta}_{con} = \begin{pmatrix} \frac{3}{2}\pi \\ \frac{1}{2}\pi \\ \frac{1}{2}\pi \end{pmatrix} - \left( \frac{\pi}{\frac{3}{2}\pi + \frac{1}{2}\pi + \frac{1}{2}\pi} \right) \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

解答完毕.

• 上述等式左右两边同时左乘  $1^T$  并利用约束条件  $1^T \theta = \pi$ , 可求得  $\lambda = (\frac{\pi}{1^T \hat{y}})$ .

**线性分析** 定义: 对线性回归模型  $\hat{y} = X\hat{\beta} + \hat{e}$ , 有  $\hat{\beta} = (X^T X)^{-1} X^T Y$ . 有  $Y$  的拟合值  $\hat{y} = X\hat{\beta} = X(X^T X)^{-1} X^T Y = H\hat{y}$ . 称  $H = X(X^T X)^{-1} X^T$  为投影矩阵 (恒为  $\hat{y} = H\hat{y}$  的充分必要条件).

性质: ①  $H^T = H, H^2 = H$ . 即  $H$  为对称幂等阵. ②  $\text{rank } H = \text{rank } X = p$ . ③  $\sum_{i=1}^n h_{ii} = p$  ④  $\sum_{i=1}^n h_{ij} = \sum_{j=1}^n h_{ij} = 1$  即  $H 1_n = 1_n$  ⑤  $HX = X$ .

定义: 对单元  $h_{ii}$  称为第  $i$  个案例的杠杆 (杠杆).  $h_{ii}$  越大, 第  $i$  个案例的作用越大. 有  $h_{ii} = \frac{1}{n} + X_i^T (X_c^T X_c)^{-1} X_i$   $X_i = (1, x_{i1}, \dots, x_{ip})$  为案例点.

定理: 对残差向量  $\hat{e} = \hat{y} - \hat{y} = (I_n - H)\hat{y} = (I_n - H)(X\hat{\beta} + \hat{e}) = (I_n - H)\hat{e}$ . 有 ①  $E[\hat{e}] = 0, \text{Cov}(\hat{e}) = \sigma^2(I_n - H)$   $h_{ii} = \frac{1}{n} + X_i^T (X_c^T X_c)^{-1} X_i$

第  $i$  个案例有  $\text{Var}(\hat{e}_i) = [\text{Cov}(\hat{e})]_{ii} = \sigma^2(1 - h_{ii})$

②  $\text{Cov}(\hat{e}, \hat{y}) = 0, \text{Cov}(\hat{y}) = \sigma^2 H$

由此可见,  $\hat{y}$  的方差比  $\hat{e}$  大, 这并不奇怪. 考虑标准化 ③ 若  $\hat{e} \sim N_n(0, \sigma^2(I_n - H))$ , 则  $\hat{y} \sim N_n(X\hat{\beta}, \sigma^2 H)$

$$r_i = \frac{\hat{e}_i}{\sigma \sqrt{1 - h_{ii}}}$$

定义: 第  $i$  个案例的标准化残差  $r_i = \frac{\hat{e}_i}{\sigma \sqrt{1 - h_{ii}}}$

定义: 若一个案例不遵从某模型且其余数据均遵从此模型, 称其为异常值.

方法: 考虑均压漂移 (线性) 模型. 设第  $i$  个案例可称为异常值.  $\begin{cases} y_i = (1, x_{i1}^T)^T \hat{\beta} + e_i \\ y_i = (1, x_{i1}^T)^T \hat{\beta} + e_i \end{cases} \hat{e} \sim N(0, \sigma^2 I_n)$  在  $2\sigma$  以下检验  $H_0: \eta = 0$  若不成立, 则第  $i$  个案例为异常值.

有结论: 若  $H_0: \eta = 0$  成立, 则  $|t_i| = r_i \sqrt{\frac{n-p-1}{n-p-r_i}} \sim t_{n-p-r_i}$ . 当  $|t_i| \geq t_{n-p-r_i}(\alpha)$  时判定为异常值.

$$t_i = r_i \sqrt{\frac{n-p-1}{n-p-r_i}} \sim t_{n-p-r_i}$$

**Example (2.1.1)**

在Forbes数据中, 检验案例12 ( $x_{12} = 204.6, y_{12} = 142.44$ ) 是否是异常值, 检验水平  $\alpha = 0.05$ . 备注:  $n = 17, \bar{x} = 202.95, \bar{y} = 0.379, SXX = 530.78$ , 拟合回归方程  $\hat{Y} = -42.131 + 0.895X$ .

$\hat{e}_{12} = y_{12} - \hat{y}_{12} = 1.36$  ← 一元回归中  $SXX = (X_{12} - \bar{x}_{12})^2 = X_c^T X_c$

$h_{12,12} = \frac{1}{n} + (x_{12} - \bar{x})^T (SXX)^{-1} (x_{12} - \bar{x}) = 0.0639$  ←  $15x$ , 共2个参数.

$r_{12} = \frac{\hat{e}_{12}}{\sigma \sqrt{1 - h_{12,12}}} = 3.7078 \Rightarrow t_{12} = r_{12} \sqrt{\frac{n-p-1}{n-p-r_{12}}} = 12.4 > t_{16}(0.025) = 2.1448$ . 是异常值.

**强影响分析**

定义: 在某些数据条件下, 若一个案例被删除后, 拟合统计量有重要改变, 则称其为强影响案例.

定义: Cook 距离  $D_i = \frac{(\hat{\beta} - \hat{\beta}_{(i)})^T X^T X (\hat{\beta} - \hat{\beta}_{(i)})}{p \hat{\sigma}^2}$   $p$  为参数个数;  $\hat{\beta}_{(i)}$  表示剔除第  $i$  个案例后 OLS 估计.

定理:  $D_i = \frac{1}{p} \left( \frac{h_{ii}}{1 - h_{ii}} \right) r_i^2$   $i = 1, \dots, n$ .  $h_{ii}$  为中度杠杆阵  $H = X(X^T X)^{-1} X^T$  的对角元,  $r_i$  为第  $i$  个案例的标准化残差.

tip: 若  $p$  固定,  $D_i$  由  $r_i$  与  $h_{ii}$  决定.  $r_i$  反映在第  $i$  个模型中其偏离程度,  $h_{ii}$  反映第  $i$  个案例与  $\bar{x}$  接近程度.

也可用 Cook 距离判定某案例是否为强影响案例?

$$D_i = \frac{1}{p} \frac{h_{ii}}{1 - h_{ii}} \cdot r_i^2 - F_{p, n-p}(\alpha)$$

Lemma:  $X \sim N_n(0, I)$  则  $X^T X \sim \chi_p^2$  类比  $X \sim N(0, \sigma^2)$ .  $\frac{\chi_p^2}{p} \sim \chi_p^2$

由上述 Lemma 得到: 由于  $\hat{\beta} - \hat{\beta}_{(i)} \sim N(0, \sigma^2(X^T X)^{-1})$ , 则  $\frac{(\hat{\beta} - \hat{\beta}_{(i)})^T X^T X (\hat{\beta} - \hat{\beta}_{(i)})}{p \hat{\sigma}^2} \sim F_{p, n-p}$

对于给定  $\alpha$ ,  $\frac{(\hat{\beta} - \hat{\beta}_{(i)})^T X^T X (\hat{\beta} - \hat{\beta}_{(i)})}{p \hat{\sigma}^2} \leq F_{p, n-p}(\alpha)$  的概率为  $1 - \alpha$ , 以  $\beta$  为变量, 其在  $p$  维空间中以  $\hat{\beta}$  为中心分布.

$\Rightarrow$  若  $D_i > F_{p, n-p}(\alpha)$ , 表示剔除第  $i$  个案例后, 落在了置信椭圆上,  $D_i$  越大表明第  $i$  个案例影响越大.

总结: 对线性回归模型  $\hat{y} = X\hat{\beta} + \hat{e}, \hat{e} \sim N(0, \sigma^2 I_n)$  时  $D_i \leq F_{p, n-p}(\alpha)$  则认为  $\hat{\beta}$  与两个 OLS 估计  $\hat{\beta}, \hat{\beta}_{(i)}$  间无重要差别, 个案例不是强影响案例.

**Example (2.1.1)**

在Forbes数据中, 检验案例12 ( $x_{12} = 204.6, y_{12} = 142.44$ ) 是否是异常值, 检验水平  $\alpha = 0.05$ . 备注:  $n = 17, \bar{x} = 202.95, \bar{y} = 0.379, SXX = 530.78$ , 拟合回归方程  $\hat{Y} = -42.131 + 0.895X$ .

**Example (2.1)**

在Forbes数据中, 由本章例1.1得, 案例12其  $h_{12,12} = 0.0639$ , 学生化残差  $r_{12} = 3.7078$ . 判断其是否是强影响数据, 检验水平  $\alpha = 0.05$ .

有  $h_{12,12} = 0.0639, r_{12} = 3.7078 \Rightarrow \text{Cook 距离 } D_{12} = \frac{1}{p} \left( \frac{h_{12,12}}{1 - h_{12,12}} \right) r_{12}^2 = 0.9 < F_{2, 15}(0.05) = 3.68$  于是不是强影响数据.

**Box-Cox 变换**

"好的残差图: 残差方差稳定在图中分布均匀"  $y$  与  $x$  线性. "坏的残差图: 残差方差不稳定在图中分布不均匀"

定理 (方差稳定变换):  $E[Y] = \mu, \text{Var}(Y) = \sigma^2$ , 若  $\sigma = \sigma(\mu)$ , 则  $Y \rightarrow \tilde{Y} = g(Y) = \int_0^Y \frac{1}{\sigma(\mu)} d\mu$  方差基本为常数 (为了满足 Gauss-Markov 假设)

例: 常用变换: ①  $\sqrt{\cdot}$  若  $\text{Var}(Y_i) \propto E[Y_i]^2$  ②  $\frac{1}{\cdot}$  若  $\text{Var}(Y_i) \propto E[Y_i]^4$   
③  $\ln Y$  若  $\text{Var}(Y_i) \propto E[Y_i]^2$  ④  $\sin^{-1}(\sqrt{\cdot})$  若  $\text{Var}(Y_i) \propto E[Y_i]^2 (1 - E[Y_i]^2)$   $0 \leq Y_i \leq 1$

例: 设  $Y$  为各生身重数的倒数,  $X$  为年龄.  $Y \sim P(\lambda)$ , 研究  $Y$  与  $X$  关系, 方差稳定. 若为线性变量  $Y = \beta_0 + \beta_1 X + e$ , 则  $\beta_0 + \beta_1 X = E[Y] = \lambda$ , 又  $\text{Var}(e) = \text{Var}(Y) + \text{Var}(\beta_0 + \beta_1 X) - 2\text{Cov}(Y, \beta_0 + \beta_1 X) = \text{Var}(Y)$ . 不满足 Gauss-Markov 假设中方差性.

于是作变换:  $Y \rightarrow \tilde{Y} = g(Y) = \int_0^Y \frac{1}{\lambda} d\lambda = 2\sqrt{Y}$ . 考虑  $\tilde{Y} = \beta_0 + \beta_1 X + e$ , 则  $\text{Var}(\tilde{Y}) \approx (\frac{1}{2\sqrt{Y}})^2 \lambda = \frac{1}{4} (\text{Var}(g(Y)) \approx (g'(\lambda))^2 \text{Var}(Y))$

定义: 假设  $Y > 0$ , 则 Box-Cox 变换为  $Y \rightarrow Y^{(\lambda)} = \begin{cases} \frac{Y^\lambda - 1}{\lambda} & \lambda \neq 0 \\ \ln Y & \lambda = 0 \end{cases}$  tip: 变换后  $Y^{(\lambda)} \sim N_n(X\hat{\beta}, \sigma^2 I_n)$

也可不用代入变换而对  $Y^{(\lambda)} = X\hat{\beta} + e, e \sim N(0, \sigma^2 I_n)$  拟合 OLS 估计?

Step 1:  $Y^{(\lambda)}$  的联合 p.d.f. 为  $L(\hat{\beta}, \sigma^2) = \frac{1}{(\sqrt{2\pi})^n \sigma^n} \exp\{-\frac{1}{2\sigma^2} (Y^{(\lambda)} - X\hat{\beta})^T (Y^{(\lambda)} - X\hat{\beta})\}$  则  $\hat{\beta}$  的 p.d.f. 为  $L(\lambda, \hat{\beta}, \sigma^2) = \frac{1}{(\sqrt{2\pi})^n \sigma^n} \exp\{-\frac{1}{2\sigma^2} (Y^{(\lambda)} - X\hat{\beta})^T (Y^{(\lambda)} - X\hat{\beta})\}$  其中  $J$  为 Box-Cox 变换的 Jacobian,  $J = \prod_{i=1}^n \left| \frac{dY_i^{(\lambda)}}{dY_i} \right| = \prod_{i=1}^n \lambda^{n-1}$

Step 2: 若  $\lambda$  已知, 有 MLE:  $\hat{\beta}(\lambda) = (X^{(\lambda)T} X^{(\lambda)})^{-1} X^{(\lambda)T} Y^{(\lambda)}, \hat{\sigma}^2(\lambda) = \frac{1}{n} Y^{(\lambda)T} (I_n - H) Y^{(\lambda)} = \frac{RSS(\lambda)}{n}$

Step 3: 将 MLE 代入, 取对数, 有  $\ln L_{max}(\lambda) = -\frac{n}{2} \ln \frac{RSS(\lambda)}{n} - \frac{n}{2} \ln \left( \frac{2\pi}{n} \right) - \frac{n}{2} \ln \frac{RSS(\lambda)}{n} + \text{Const}$

Step 4: 对  $\lambda$  求导为 0 得到  $\lambda$  的估计, 令  $l(\lambda) = \ln L_{max}(\lambda) = \ln \left( \frac{RSS(\lambda)}{n} \right)^{-\frac{n}{2}} \left( \frac{2\pi}{n} \right)^{-\frac{n}{2}} \left( \frac{RSS(\lambda)}{n} \right)^{-\frac{n}{2}} = -\frac{n}{2} \ln \frac{RSS(\lambda)}{n} + \text{Const}$ . 找到  $\lambda$  使  $l(\lambda)$  最大, 即  $\lambda = \hat{\lambda}$ . 找到  $\lambda$  使  $l(\lambda)$  最大, 即  $\lambda = \hat{\lambda}$ .

**OLS 估计**

对线性回归模型  $\hat{y} = X\hat{\beta} + \hat{e}, E[\hat{e}] = 0, \text{Cov}(\hat{e}) = \sigma^2 I, I$  为正定阵 (不设为  $I_n$ )

$$RSS = \hat{y}^T (I_n - H) \hat{y} = e^T (I_n - H) e$$

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

① 已知 线性变换:  $T: \mathbb{R}^n \rightarrow \mathbb{R}^n$   $T \rightarrow T^{-1}$   $T^{-1}y = X^{-1}y$

于是有  $\hat{\beta}_a = (U^T U)^{-1} U^T z = (X^T X)^{-1} X^T z$  称为  $\beta$  的 OLS 估计 (GLS 估计)

$$(X^T X)^{-1} X^T z$$

定理: ①  $E(\hat{\beta}_a) = \beta$  ②  $Cov(\hat{\beta}_a) = \sigma^2 (X^T X)^{-1}$  ③ 对  $v \in \mathbb{R}^p$  线性函数  $C^T \hat{\beta}_a$  为  $C^T \beta$  的唯一最小方差线性估计 (对一般线性回归模型, 无论与最小方差意义下 GLS 相对于 OLS 估计)

例: 若  $\Sigma = Cov(\hat{\epsilon}) = \text{diag}(\sigma_1^2, \dots, \sigma_k^2)$ , 则  $\hat{\beta}$  为 GLS 估计.

② 未知

设  $X$  与  $n$  个向量为  $x_1, \dots, x_n$ , 则由公式  $\hat{\beta}_a = \left( \sum_{i=1}^n \frac{1}{\sigma_i^2} x_i x_i^T \right)^{-1} \sum_{i=1}^n \frac{1}{\sigma_i^2} x_i y_i$

$$\sum x_i x_i = X^T X = X^T e^T X$$

例: Example (4.2)

【例 3.6.1】在线性回归模型(1)中, 若观测向量  $y$  是  $n$  件样品的某项指标, 它们是在  $k$  台仪器上测试得到的. 不妨设前  $n_1$  个是在第一台仪器上测得的, 试验和测量误差为  $\sigma_1^2$ , 接下来的  $n_2$  个是在第二台仪器上测得的, 测量误差为  $\sigma_2^2$ , 依此类推. 并且假设仪器之间的试验和测量误差互不相关. 这里  $n_1 \geq p$ ,  $\sigma_i^2$  未知,  $i = 1, \dots, k$ . 求模型(1)中回归系数  $\beta$  的 GLS 估计.

模型可写为  $y_i = x_i \beta + e_i$ ,  $E(e_i) = 0$ ,  $Cov(e_i) = \sigma_i^2 I_{n_i}$

对  $\sigma_1^2$ : 有  $\sigma_1^2$  的 OLS 估计:  $\hat{\sigma}_1^2 = \frac{||y_1 - X_1 \hat{\beta}_1||^2}{n_1 - p}$

若  $\sigma_1^2, \dots, \sigma_k^2$  已知, 有  $\beta$  的 GLS 估计  $\hat{\beta}_a = \left( \sum_{i=1}^k \frac{X_i^T X_i}{\sigma_i^2} \right)^{-1} \sum_{i=1}^k \frac{X_i^T y_i}{\sigma_i^2}$ , 称为  $\beta$  的 GLS 估计.

$$SS_{reg} = \frac{S_{xy}^2}{S_{xx}} \quad R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

tip: 无截距时  $SYT + R_{SS} + ESS$

第一套:

期中部分总结

$$1. \hat{y} = X\hat{\beta} + \hat{e}, \hat{e} \sim N(0, \sigma^2 I_n) \quad E[E(\hat{e})] = 0, Cov(\hat{e}) = \sigma^2 I_n$$

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

$$2. \hat{\beta} = (X^T X)^{-1} X^T y, \text{ 无截距时 } \hat{\beta} = \frac{S_{xy}}{S_{xx} + n\bar{x}^2}$$

$$RSS = SYT + n\bar{y}^2 - \frac{(S_{xy} + n\bar{x}\bar{y})^2}{S_{xx} + n\bar{x}^2} \quad R^2 = \frac{(S_{xy} + n\bar{x}\bar{y})^2}{(S_{xx} + n\bar{x}^2)(SYT + n\bar{y}^2)}$$

合截距

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

$$\hat{\beta} = \frac{S_{xy}}{S_{xx}} = (X^T X_c)^{-1} X_c^T y$$

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = SYT - \frac{S_{xy}^2}{S_{xx}} \quad R^2 = \frac{ESS}{SYT} = 1 - \frac{RSS}{SYT} = \frac{(S_{xy})^2}{S_{xx} \cdot SYT}$$

$$3. \hat{y} = \hat{\alpha} + X_c \hat{\beta} + \hat{e} \quad \text{中心化时 } \hat{\alpha} = \bar{y}, \hat{\beta}_0 = \hat{\alpha} - \bar{\beta} \bar{x}$$

$$4. E(\hat{\beta}) = \beta, Cov(\hat{\beta}) = \sigma^2 (X^T X)^{-1}, Cov(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

$Cov(\hat{\beta}) = \sigma^2 (X^T X_c)^{-1}$  其中  $X_c$  是不含截距  $\beta_0$  只有斜率部分  $E, Cov$

特别地,  $n=1$  时,  $\hat{\beta} = \frac{\sigma^2}{S_{xx}} (S_{xx} = X_c^T X_c)$

5. 约束 OLS 估计先构造辅助函数, 一般后面用  $+2(A\hat{\beta} - b)^T$

$$\frac{1}{n} + (y_i - \bar{y}) (X_c^T X_c)^{-1} (y_i - \bar{y}), \hat{e} = \frac{1}{\sqrt{1 - h_{ii}}} \hat{e}_i \quad D_i = \frac{1}{1 - h_{ii}} \cdot \frac{h_{ii}^2}{1 - h_{ii}} \cdot t_i^2$$

第二套:

$$6. \text{ 中间矩阵 } H = X(X^T X)^{-1} X^T \quad HX = X, h_{11} = x_1^T (X^T X)^{-1} x_1, h_{ii} = \frac{1}{n} + x_i^T (X_c^T X_c)^{-1} x_i \quad \text{中心化的性质: } h_{SS} = \frac{1}{n} + (x_S - \bar{x})^T (X_c^T X_c)^{-1} (x_S - \bar{x})$$

$$7. \text{ 学生化残差 } r_i = \frac{\hat{e}_i}{\sigma \sqrt{1 - h_{ii}}} \Rightarrow \text{ 判断异常值: 若 } |r_i| \sqrt{\frac{n-p-1}{n-p-2}} \geq t_{n-p-1}(\frac{\alpha}{2}) \text{ 为异常值. } \textcircled{1} E(\hat{e}) = 0, Cov(\hat{e}) = \sigma^2 (I_n - H)$$

$$A = \begin{bmatrix} X_c^T X_c & X_c^T \bar{y} \\ \bar{y}^T X_c & S_{yy} \end{bmatrix} \quad ESS = (X_c^T y)^T (X_c^T X_c)^{-1} (X_c^T y)$$

$$\left( \frac{S_{xy}^2}{S_{xx}} \right)$$

$$t_i = |r_i| \sqrt{\frac{n-p-1}{n-p-2}}$$

- ②  $Cov(\hat{e}, \hat{y}) = 0, Cov(\hat{y}) = \sigma^2 H$
- ③ 若  $\hat{e} \sim N_n(0, \sigma^2 I_n)$ , 则  $\hat{e} \sim N_n(0, \sigma^2 (I_n - H)), \hat{y} \sim N_n(X\beta, \sigma^2 H)$

Cook 距离:  $D_i = \frac{1}{p} \left( \frac{h_{ii}}{1 - h_{ii}} \right) \cdot r_i^2 \Rightarrow \text{ 判断异常值: 若 } D_i \geq F_{p, n-p}(\alpha)$  为异常影响点

8. Box-Cox 选择

$$GLS \text{ 估计 } \hat{\beta} = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} y$$

关于 Gauss-Markov 假设的若干问题

- 问题 1: 如果随机误差项不满足 Gauss-Markov 假设, 上述性质是否仍成立?
- 问题 2: 观测值向量  $y$  与其它随机向量的相关关系?  
 $Cov(\hat{y}, \hat{e}) = ? \quad Cov(\hat{y}, \hat{e}) = ? \quad Cov(\hat{y}, \hat{\beta}) = ? \quad Cov(\hat{y}, \hat{y}) = ?$   
 $= Cov(y, (I_n - H)y) = Cov(y, (I_n - H)y) = \sigma^2 (I_n - H)$
- 问题 3: 拟合值向量  $\hat{y}$  与其它随机向量的相关关系?  
 $Cov(\hat{y}, \hat{e}) = ? \quad Cov(\hat{y}, \hat{e}) = ? \quad Cov(\hat{y}, \hat{\beta}) = ?$   
 $= \sigma^2 H(I_n - H) \quad = X(X^T X)^{-1} X^T \sigma^2 (I_n - H)$

- 问题 4: 考虑简单线性回归模型  $y_i = \beta_0 + \beta_1 x_i + e_i, i = 1, 2, \dots, n, e_i$  满足 Gauss-Markov 假设. 若  $x_i, y_i$  均随机, 此时是否仍成立  $E(y_i | x_i) = \beta_0 + \beta_1 x_i$ ?  
假设  $Var(x_i) = \tau^2, i = 1, \dots, n$ , 且  $x_i$  与  $e_i$  独立, 计算  $Var(y_i)$  和  $Var(y_i | x_i) = Var(e_i) = \sigma^2$   
 $y_i = \beta_0 + \beta_1 x_i + e_i, Var(y_i) = \beta_1^2 \tau^2 + \sigma^2$

- 问题 1: 性质(3)中第二个等式表明残差向量与拟合值向量相互正交, 那么残差向量与解释向量  $(x_1, \dots, x_{p-1})$  呢?  
答案: 相互正交.
- 问题 2: 对于无截距的回归模型,  
 $y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_{p-1} x_{i,p-1} + e_i, i = 1, 2, \dots, n,$   
 $e_i, i = 1, \dots, n$  满足 Gauss-Markov 假设. 性质(3)中两个等式是否仍然成立? 性质(1)和(2)是否也仍然满足?  
答案: 除了性质(3)中第二个等式仍然成立, 其余性质(1)-(3)均不一定成立.

关于 OLS 估计是 BLUE 的若干问题

- 问题 1: 对于线性回归模型(1),  $\hat{\beta}_0$  是否是  $\beta_0$  的 BLUE? 也即,  $\hat{\beta}_0$  是否是  $\beta_0$  在线性无偏估计类  $\{ \sum_{i=1}^n a_i y_i : a_i \in \mathbb{R}, i = 1, \dots, n \}$  中的 UMVUE?
- 问题 2: 对于中心化的线性回归模型  $\hat{y} = a \cdot \bar{1}_n + X_c \hat{\beta}_c + \hat{e}, E(\hat{e}) = 0, Cov(\hat{e}) = \sigma^2 I_n, rank(X_c) = p - 1,$  其 OLS 估计  $\hat{a}$  是否是  $a$  的 BLUE?
- 问题 3: 若  $rank(X) < p$ , 则本章第一节(5)的解  $\hat{\beta} = (X^T X)^{-1} X^T y$  是否成立 Gauss-Markov 定理?  
答案: 不一定, 参考【1】定理 4.1.2.

结论: 是  
解释:  
1. 线性性: 截距 OLS 估计量  $\hat{a} = \hat{\beta}_0 = \hat{\beta}_0$ , 是  $y$  的线性函数, 均  $y$  的线性组合, 满足线性条件.  
2. 无偏性: 在 Gauss-Markov 假设下,  $E(\hat{a}) = a$ , 无偏性成立.  
3. 最佳性: 由 Gauss-Markov 定理, OLS 估计量是所有线性无偏估计中的方差最小者, 因此  $\hat{a}$  是  $a$  的 BLUE.  
4. 补充: 若进一步假设误差服从正态分布  $e \sim N(0, \sigma^2 I_n)$ , 则  $\hat{a}$  同时是一致最小方差无偏估计 (UMVUE), 即所有无偏估计 (不限于线性) 中方差最小.

结论: 是  
解释与证明:  
中心化的线性回归模型:  
 $\hat{y} = a \cdot \bar{1}_n + X_c \hat{\beta}_c + \hat{e}, E(\hat{e}) = 0, Cov(\hat{e}) = \sigma^2 I_n, rank(X_c) = p - 1$   
其中  $\bar{1}_n$  是全 1 向量,  $X_c$  是中心化后的解释变量矩阵 (均值为 0, 因此  $\bar{1}_n^T X_c = 0$ ), 第一个方程简化为  $na = \sum_{i=1}^n y_i$ .  
1. OLS 估计量: 正规方程:  
 $\begin{cases} \bar{1}_n^T (a \cdot \bar{1}_n + X_c \hat{\beta}_c - \hat{y}) = 0 \\ X_c^T (a \cdot \bar{1}_n + X_c \hat{\beta}_c - \hat{y}) = 0 \end{cases}$   
由  $\bar{1}_n^T X_c = 0$ , 第一个方程简化为  $na = \sum_{i=1}^n y_i$ .  
2. BLUE 证明:  
线性性:  $\hat{a} = \frac{1}{n} \sum_{i=1}^n y_i$ , 是  $y$  的线性函数.  
- 无偏性:  $E(\hat{a}) = E(\frac{1}{n} \sum_{i=1}^n y_i) = \frac{1}{n} \sum_{i=1}^n E(y_i) = a$ , 因此  $\hat{a}$  是无偏的.  
- 最佳性: 由 Gauss-Markov 定理, 中心化的 OLS 估计量  $(\hat{a}, \hat{\beta}_c)^T$  是 BLUE, 因此  $\hat{a}$  是  $a$  的 BLUE.  
3. 解释: 证明: 中心化的 OLS 估计量  $\hat{a}$  是最佳线性无偏估计 (BLUE) 估计量, 符合 BLUE 的定义.

结论: 不一定成立  
解释:  
1. 最小方差性: Gauss-Markov 定理的必要前提是  $X$  列满秩 ( $rank(X) = p$ ), 此时  $X^T X$  可逆, OLS 估计量  $\hat{\beta} = (X^T X)^{-1} X^T y$  唯一可识别. 当  $rank(X) < p$  时,  $X^T X$  不可逆, 只能推广使用  $(X^T X)^+$  求解, 此时:  
- 估计量不唯一: 不同广义逆会得到不同的  $\hat{\beta}$ ;  
- 参数不可识别: 无法唯一区分不同解释变量的系数 (多重共线性完全);  
- 无偏性仅对可估计函数成立, 单个系数  $\beta_j$  不再具有唯一的无偏估计.  
2. 广义逆的性质:  $\hat{\beta} = (X^T X)^+ X^T y$  是正规方程的解, 但:  
- 仅保证拟合值  $\hat{y} = X \hat{\beta}$  唯一 (由广义逆选择决定);  
- 单个系数  $\beta_j$  不唯一, 因此不满足 Gauss-Markov 定理对“唯一线性无偏最小方差估计”的要求.  
3. 补充说明: 仅当估计的可估计函数 (如  $c^T \beta$ , 其中  $c \in Col(X)$ ,  $C$  为列空间) 时,  $c^T \hat{\beta}$  是唯一的 BLUE; 但单个系数  $\beta_j$  不可作为估计函数, 不存在 BLUE.

$$S_{xx} = X_c^T X_c$$

- ④ 约束最小二乘估计  $\hat{\beta}_{con}$  是否是  $\beta$  的无偏估计?  
答案: 在约束条件下, 是, 无约束条件下, 是否无偏?
- ⑤ 对于回归参数有线性约束条件的回归模型(1), 在任一关于  $\hat{\beta}$  的线性函数  $c^T \hat{\beta}$  的所有线性无偏估计中, OLS 估计  $c^T \hat{\beta}$  是否仍是唯一具有最小方差的估计?  
答案: 不是. 约束 OLS 估计  $c^T \hat{\beta}_{con}$  是无偏的且方差可以更小, 参考无约束条件下 OLS 估计的 Gauss-Markov 定理证明.  
⑥ 定理 1.3.1 结果对于无截距模型是否成立?  
答案: 仍成立.  
⑦ 中心化模型下, 假设回归参数有约束条件  $A\hat{\beta}_c = b$ , 求此时参数的约束最小二乘估计?  
答案:  $\hat{\beta}_{con} = \hat{\beta} - (X^T X)^{-1} A^T (A (X^T X)^{-1} A^T)^{-1} (A \hat{\beta} - b)$

Problem

与最小二乘 (OLS) 估计  $\hat{\beta} = (X^T X)^{-1} X^T y$  相比, 参数  $\eta = c^T \hat{\beta}$  的另一估计  $\hat{\eta}_c = c^T \hat{\beta}_c$  方差更小还是更大?

答案: 当随机误差项  $\hat{e}$  满足  $E(\hat{e}) = 0, Cov(\hat{e}) = \sigma^2 I$  时,  $c^T \hat{\beta}$  方差最小; 而若  $Cov(\hat{e}) = \sigma^2 \Sigma \neq \sigma^2 I$  时, 由定理 4.1 (c) 可知, 对于任意  $\hat{c} \in \mathbb{R}^p$ ,

$$Var\left(c^T \hat{\beta}_c\right) \leq Var\left(c^T \hat{\beta}\right)$$

- 即对于一般线性回归模型(1), 以线性无偏和最小方差意义而言, 广义最小二乘 (GLS) 估计优于最小二乘 (OLS) 估计.

# 方差分析表 1. 单因素

问题:  $\beta \sim N(0, \sigma^2 I)$ , 两模型为  $\gamma = \beta_0 + \beta_1 X_1 + \dots + \beta_{p-1} X_{p-1} + e$   
 $\gamma = \beta_0 + e$ . 是否有显著区别? 检验水平为  $\alpha$ .

答: 检验  $H_0: \beta_1 = \dots = \beta_{p-1} = 0 \leftrightarrow H_1: \beta_1, \dots, \beta_{p-1}$  中至少有 1 个非 0.

若  $\sigma^2$  已知,  $\frac{ESS}{\sigma^2} \sim \chi_{p-1}^2$ . 当  $\frac{ESS}{\sigma^2} > \chi_{p-1}^2(\alpha)$  时认为两模型有显著区别

若  $\sigma^2$  未知, 有  $\hat{\sigma}^2 = \frac{RSS}{n-p}$ .  $\frac{ESS}{\hat{\sigma}^2} \sim F_{p-1, n-p}$ . 当  $\frac{ESS}{\hat{\sigma}^2} > F_{p-1, n-p}(\alpha)$  时认为两模型有显著区别

记回归均方为  $MS_{reg} = SS_{reg} / (p-1)$ , 则上述讨论可归结为下表

方差分析表 (ANOVA)

来源	d.f.	SS	MS	F
X 的回归	p-1	SS <sub>reg</sub>	SS <sub>reg</sub> / (p-1)	MS <sub>reg</sub> / $\hat{\sigma}^2$
残差	n-p	RSS	RSS / (n-p)	
总的	n-1	SS <sub>Y</sub>		

定义  $MSS \equiv \frac{ESS}{p-1}$

## 2. 单个因变量的显著性检验

问题:  $\beta \sim N(0, \sigma^2 I)$ , 两模型为  $\gamma = \beta_0 + \beta_1 X_1 + \dots + \beta_{p-1} X_{p-1} + e$   
 $\gamma = \beta_0 + \beta_2 X_2 + \dots + \beta_{p-1} X_{p-1} + e$ . 是否有显著区别? 检验水平为  $\alpha$ .

答: 检验  $H_0: \beta_1 = 0 \leftrightarrow H_1: \beta_1 \neq 0$

一般地, 对  $H_0: \beta_j = 0 \leftrightarrow H_1: \beta_j \neq 0, |j| = 1, \dots, p-1$  时均适用.

定义  $(C_j)_{j=1, \dots, p-1} = (X_2^T X_2, \dots, X_{p-1}^T X_{p-1})^T$ .  $C_{00} = \frac{1}{n} + \bar{m}^T (X_2^T X_2)^{-1} \bar{m}$ .  $\hat{\beta}_j \sim N_{1, n-p}$ . 则当  $|t_{j1}| > t_{n-p}(\alpha)$  时拒绝  $H_0$ , 认为两模型有显著区别

或: 完整模型下残差平方和为  $RSS$ . 定义  $SS_{reg}^{(j)} = RSS - RSS_{(j)}$ ,  $\hat{MS}_{reg}^{(j)} = \frac{SS_{reg}^{(j)}}{1}$ . 有  $F_j = \frac{\hat{MS}_{reg}^{(j)}}{\hat{\sigma}^2} \sim F_{1, n-p}$ .  $F_j > F_{1, n-p}(\alpha)$  时认为两模型有显著区别

去除  $X_1$  模型下残差平方和为  $RSS_{(1)}$

## 3. 两个因变量的显著性检验

问题:  $\beta \sim N(0, \sigma^2 I)$ , 两模型为  $\gamma = \beta_0 + \beta_1 X_1 + \dots + \beta_{p-1} X_{p-1} + e$   
 $\gamma = \beta_0 + \beta_2 X_2 + \dots + \beta_{p-1} X_{p-1} + e$ . 是否有显著区别? 检验水平为  $\alpha$ .

或: 完整模型下残差平方和为  $RSS$ . 定义  $SS_{reg}^{(1,2)} = RSS - RSS_{(1,2)}$ ,  $\hat{MS}_{reg}^{(1,2)} = \frac{SS_{reg}^{(1,2)}}{2}$ . 有  $F_2 = \frac{\hat{MS}_{reg}^{(1,2)}}{\hat{\sigma}^2} \sim F_{2, n-p}$ .  $F_2 > F_{2, n-p}(\alpha)$  时认为两模型有显著区别

去除  $X_1, X_2$  模型下残差平方和为  $RSS_{(1,2)}$

一般地,  $\hat{MS}_{reg}^{(1, \dots, k)} = \frac{SS_{reg}^{(1, \dots, k)}}{k}$  有  $F_k = \frac{\hat{MS}_{reg}^{(1, \dots, k)}}{\hat{\sigma}^2} \sim F_{k, n-p}$ .  $F_k > F_{k, n-p}(\alpha)$  时认为两模型有显著区别

### 一般线性假设的应用

思考: 考虑  $y = X\beta + e, e \sim N(0, \sigma^2 I)$ . 检验  $\beta$  是否满足假设  $H_0: A\beta = b$ .

若不考虑限制有  $RSS = y^T (I - X(X^T X)^{-1} X^T) y$

附加线性假设, 有  $RSS_H = (y - X\hat{\beta}_H)^T (y - X\hat{\beta}_H)$ ,  $\hat{\beta}_H = (X^T X)^{-1} A^T (A(X^T X)^{-1} A^T)^{-1} (A\hat{\beta} - b)$

若参数服从满足  $A\beta = b$ , 对数据拟合程度一样,  $RSS_H - RSS$  较小. 反之较大, 于是我们可在  $RSS_H - RSS$  较大时拒绝  $H_0$ .

定理: ①  $\frac{RSS_H - RSS}{\sigma^2} \sim \chi_{m}^2$  ② 若  $A\beta = b$  成立,  $\frac{RSS_H - RSS}{\sigma^2} \sim \chi_{m}^2$  ③  $RSS$  与  $RSS_H - RSS$  相互独立. ④  $A\beta = b$  时  $\frac{(RSS_H - RSS)/m}{RSS/(n-p)} \sim F_{m, n-p}$ . ( $m = \text{rank } A$ )

tip: 由 ④ 知  $\frac{(RSS_H - RSS)/m}{RSS/(n-p)} > F_{m, n-p}(\alpha)$  时拒绝  $H_0$ .  $RSS_H - RSS = (A\beta - b)^T (A(X^T X)^{-1} A^T)^{-1} (A\hat{\beta} - b)$

## 回归方程的显著性检验

思考:  $H: \beta_1 = \dots = \beta_{p-1} = 0$ . 若拒绝  $H_0$ ,  $\gamma$  线性依赖于至少一个  $X_j$ ; 若接受  $H_0$ , 可认为相对于误差而言, 所有自变量与  $\gamma$  均不显著.

明显地, 这是上一步中  $A = (0, I_{p-1})$ ,  $b = 0$  的特殊情况, 于是  $m = p-1$

结论: 由上述分析,  $\frac{(RSS_H - RSS)/(p-1)}{RSS/(n-p)} > F_{p-1, n-p}(\alpha)$  时拒绝  $H_0$ .

## 因变量的预测

考虑线性回归模型  $y_0 = (1, x_0^T) \beta + e$ ;  $e_i \in N(0, \sigma^2), 1 \leq i \leq n$ . 给定  $x_0 = (x_{01}, \dots, x_{0, p-1})^T, y_0$  未知. 若也满足  $y_0 = (1, x_0^T) \beta + e_0, e_0 \in N(0, \sigma^2)$

结论:  $y_0$  的点预测:  $\hat{y}_0 = (1, x_0^T) \hat{\beta}$ .

性质: ①  $\hat{y}_0$  为  $y_0$  的无偏预测 ② 在  $y_0$  的线性无偏预测中,  $\hat{y}_0$  具有最小方差. tip: 此处“无偏”指二者同维空间.

定义: 使预测值落在其中概率达到预定置信区间的称为预测区间.

记  $d_0 = \frac{1}{n} + (x_0^T - \bar{m}^T) (X^T X)^{-1} (x_0^T - \bar{m})$

$y_0$  的点预测区间:  $\sigma^2$  已知时  $\frac{y_0 - \hat{y}_0}{\sigma \sqrt{1+d_0}} \sim N(0, 1)$ , 置信度为  $1-\alpha$  的  $y_0$  点预测区间为  $[\hat{y}_0 - \sigma \sqrt{1+d_0} u_{\alpha/2}, \hat{y}_0 + \sigma \sqrt{1+d_0} u_{\alpha/2}]$

$(\frac{n-p}{\sigma^2}) \hat{\sigma}^2 \sim \chi_{n-p}^2$

$\sigma^2$  未知时  $\frac{y_0 - \hat{y}_0}{\hat{\sigma} \sqrt{1+d_0}} \sim t_{n-p}$ . 置信度为  $1-\alpha$  的  $y_0$  点预测区间为  $[\hat{y}_0 - \hat{\sigma} \sqrt{1+d_0} t_{n-p}(\alpha/2), \hat{y}_0 + \hat{\sigma} \sqrt{1+d_0} t_{n-p}(\alpha/2)]$

$y_0$  的区间预测:  $\sigma^2$  已知时  $\frac{E(y_0) - \hat{y}_0}{\sigma \sqrt{d_0}} \sim N(0, 1)$ , 置信度为  $1-\alpha$  的  $y_0$  区间预测为  $[\hat{y}_0 - \sigma \sqrt{d_0} u_{\alpha/2}, \hat{y}_0 + \sigma \sqrt{d_0} u_{\alpha/2}]$

$(\frac{n-p}{\sigma^2}) \hat{\sigma}^2 \sim \chi_{n-p}^2$

$\sigma^2$  未知时  $\frac{E(y_0) - \hat{y}_0}{\hat{\sigma} \sqrt{d_0}} \sim t_{n-p}$ . 置信度为  $1-\alpha$  的  $y_0$  区间预测为  $[\hat{y}_0 - \hat{\sigma} \sqrt{d_0} t_{n-p}(\alpha/2), \hat{y}_0 + \hat{\sigma} \sqrt{d_0} t_{n-p}(\alpha/2)]$

## 共线性相关

定义: 设  $\hat{\beta}$  为未知参数向量  $\hat{\beta}$  由  $\hat{\beta} = (X^T X)^{-1} X^T y$  估计.  $\hat{\beta}$  的协方差定义为  $MSE(\hat{\beta}) = E[\|\hat{\beta} - \beta\|^2]$ ,  $\|\cdot\|$  为  $R^p$  上的欧氏范数

定理:  $MSE(\hat{\beta}) = \text{tr } Cov(\hat{\beta}) + \|E[\hat{\beta}] - \beta\|^2$

考虑  $\hat{\beta} = \alpha_2 + X_3 \beta_3 + e$ .  $X_3$  为中心化矩阵.  $X_3$  的协方差矩阵  $X_3 X_3^T = \frac{X_{ij} - m_j}{s_j} s_j = \sqrt{S} X_3 X_3^T$ .  $\hat{\alpha} = \hat{\beta}$ .  $\hat{\beta}_3 = (X_3^T X_3)^{-1} X_3^T y$  为  $\beta_3$  的无偏估计

$MSE(\hat{\beta}_3) = \text{tr } Cov(\hat{\beta}_3) = \sigma^2 \text{tr } (X_3^T X_3)^{-1} = \sigma^2 \sum_{i=1}^p \frac{1}{\lambda_i}$ .  $\lambda_i$  为  $X_3^T X_3$  的特征值

定义:  $X_3^T X_3$  的特征值为  $\lambda_1 \geq \dots \geq \lambda_{p-1} > 0$ . 条件数  $\kappa = \frac{\lambda_1}{\lambda_{p-1}}$  为最大与最小特征值之比

$\kappa < 100$  认为共线性程度低  
 $100 \leq \kappa \leq 1000$  认为共线性程度较高  
 $\kappa > 1000$  认为共线性程度严重

问题: 认为预测变量  $X_1, \dots, X_{p-1}$  之间存在较强/多重共线性时, 怎样找共线性关系?

回答: 对正定阵  $X_3^T X_3$ , 存在交阵  $S$ , s.t.  $X_3^T X_3 = S \text{diag}(\lambda_1, \dots, \lambda_{p-1}) S^T$ . 其中  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{p-1} > 0$ .

若  $\lambda_j \approx 0, \hat{y}_j = (\hat{y}_{j1}, \dots, \hat{y}_{j, p-1})^T, X_3 = (x_{31}, x_{32}, \dots, x_{3, p-1})$ . 有  $\hat{y}_{j1} x_{31} + \dots + \hat{y}_{j, p-1} x_{3, p-1} \approx 0 \Rightarrow \hat{y}_{j1} \frac{x_{31}}{s_1} + \dots + \hat{y}_{j, p-1} \frac{x_{3, p-1}}{s_{p-1}} \approx 0$

Example (4.2.1)

【例3.7.1】考虑一个有六个回归自变量的线性回归模型，  

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_6 X_6 + e,$$
 原始数据见[0]表3.7.1，共12组数据。  
 对于如上模型，记中心化和标准化设计阵为 $X_s$ ，则自变量 $X_1, \dots, X_6$ 的样本相关系数矩阵 $X_s^T X_s$ 的上三角阵为

1.000	0.052	-0.343	-0.498	0.417	-0.192
1.000	-0.432	-0.371	0.485	-0.317	
	1.000	-0.355	-0.505	0.494	
		1.000	-0.215	-0.087	
			1.000	-0.123	
				1.000	

注 从非对角元的绝对值看，任两个回归自变量之间似乎不存在较严重的线性依赖关系，但事实是，多个自变量之间存在严重的共线性。

• 计算矩阵 $X_s^T X_s$ 的六个特征根，分别为  
 $\lambda_1 = 2.24879, \lambda_2 = 1.54615, \lambda_3 = 0.92208,$   
 $\lambda_4 = 0.79399, \lambda_5 = 0.30789, \lambda_6 = 0.00111.$

• 样本均值和样本标准方差  
 $m_1 = 2.5, m_2 = 2, m_3 = 2.5, m_4 = 3.083,$   
 $s_1 = 11.358, s_2 = 10.1, s_3 = 13.077, s_4 = 13.89.$

•  $\lambda_6 \approx 0$ ，因此存在共线性关系；同时共线性的度量：条件数  

$$\kappa = \frac{\lambda_1}{\lambda_6} = \frac{2.24879}{0.00111} = 2025.94 > 1000.$$

• 根据我们本节约定的标准，模型的设计阵（或自变量之间）存在严重的（复）共线性。  
 • 利用 $\lambda_6 = 0.00111 \approx 0$ ，我们进一步算出 $X_s^T X_s$ 对应于 $\lambda_6$ 的特征向量为  

$$\vec{\varphi} = (-0.45, -0.42, -0.54, -0.57, -0.006, -0.002)^T,$$

• 因而回归自变量之间有如下（复）共线关系  

$$0.45 \frac{X_1 - m_1}{s_1} + 0.42 \frac{X_2 - m_2}{s_2} + 0.54 \frac{X_3 - m_3}{s_3} + 0.57 \frac{X_4 - m_4}{s_4} \approx 0. \quad (4)$$
 其中由于 $\varphi_5, \varphi_6 \approx 0$ ，已删除相关变量。  
 • 代入样本均值和样本标准方差，可得  

$$0.0396X_1 + 0.0416X_2 + 0.0413X_3 + 0.0410X_4 \approx 0.4120$$

岭估计

考虑回归模型  $\hat{y} = \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_k x_k + \epsilon, \epsilon \sim N(0, \sigma^2 I_n)$   
 定义：对上述模型， $\hat{\beta}_s$  的岭估计为  $\hat{\beta}_s(k) = (X_s^T X_s + k I_{p+1})^{-1} X_s^T y$   
 这里  $k > 0$  称岭系数或偏差参数  
 $\hat{\beta}_s(k)$  是一组估计量， $\hat{\beta}_s(0)$  即为常规的 OLS 估计  
 定理：在 MSE 意义下，岭估计优于 OLS 估计。即  $\exists k > 0, \text{MSE}(\hat{\beta}_s(k)) < \text{MSE}(\hat{\beta}_s)$

岭系数的选择

Hoerl-Kennard 公式： $\hat{\sigma}^2 = \frac{\hat{\sigma}^2}{\max \alpha_i^2}, \hat{\alpha} = \hat{\alpha}^T \hat{\beta}_s$  称为正则化因子

岭估计：原理：找岭系数  $k$ ，使得岭估计估计值趋于稳定，且  $k$  不宜过于过大

定义：记  $\hat{\beta}_{s,j}(k)$  为  $\hat{\beta}_s(k)$  的第  $j$  分量，为  $k$ -元函数， $k \in [0, \infty)$  变化时， $\hat{\beta}_{s,j}(k)$  的图形称为岭迹。

- 岭迹法选择  $k$  的具体操作：
  - 将  $\hat{\beta}_{s,1}(k), \dots, \hat{\beta}_{s,p+1}(k)$  的岭迹画在同一图上；
  - 根据岭迹的变化趋势选择  $k$  值，使得各个回归系数的估计大体上稳定，并且各个回归系数岭估计值的符号比较合理；
  - 同时选取的  $k$  应尽量接近零，使得残差平方和不要上升太多。

一些公式： $\hat{\alpha}_j = \hat{\beta}_j \cdot \frac{s_j}{\sigma_j}, \hat{\beta}_{s,j}(k) = \hat{\alpha}_j(k) \cdot \frac{\sigma_j}{s_j}, \hat{\beta}_s(k) = \hat{\alpha} - \sum_{j=1}^p \hat{\beta}_s(k) \hat{\alpha}_j$

主成分估计

考虑线性回归模型  $y = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_p x_p + e, e \sim N(0, \sigma^2 I_n), \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p > \lambda_{p+1} = 0$  为  $X_s^T X_s$  的特征值， $\varphi_1, \dots, \varphi_p$  为标准正交特征向量， $\Phi = (\varphi_1, \dots, \varphi_p)$  为  $p \times p$  正交阵。

记  $Z = X_s \Phi, \alpha = \alpha^T \hat{\beta}_s$ ，则上述模型为  $y = \alpha_0 + Z\alpha + e, \alpha$  为正则化回归系数。

定义： $Z = (\bar{z}_1, \bar{z}_2, \dots, \bar{z}_p)$ ，称  $\bar{z}_k$  为第  $k$  主成分

解法：将对应特征值  $\approx 0$  的  $p+1-r$  主成分从模型中剔除，再对剩下  $r$  个主成分的 OLS 估计，得到主成分回归方程。此过程称为主成分回归。

定理：当设计阵存在共线性时，适当选择保留的主成分个数可使主成分估计的 MSE 比 OLS 估计更小，即  $\text{MSE}(\hat{\beta}_{s,p}) < \text{MSE}(\hat{\beta}_s)$

$\frac{\hat{\beta}_j}{\hat{\sigma}^2 \alpha_i}$

$C_{eo} = \frac{1}{n} m^T (X_s^T X_s)^{-1} m$

$\frac{1}{n} + \frac{\sigma^2}{s^2 + \lambda}$

1. 单个检验:  $H_0: \beta_j = 0$ . 则取  $t_j = \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{C_{jj}}} \sim t_{n-p}(\frac{\alpha}{2})$   $\hat{\sigma} = \sqrt{MSE} = \sqrt{\frac{RSS}{n-p}}$   $C_{00} = \frac{1}{n} + \bar{m}^T (X_0^T X_0)^{-1} \bar{m}$   $C_{ii}$  为  $(X_0^T X_0)^{-1}$  的对角元 ( $-C_{ii}$  为  $\frac{1}{S_{XX}}$ )

2.  $STT = SS_{reg} + RSS \Rightarrow$  残差均方  $MSE = \frac{RSS}{df_e}$   $(-C_{ii}$  为  $\frac{1}{n} + \frac{\bar{x}^2}{S_{XX}})$   $(S_{XX} = X_0^T X_0)$

回归均方  $MS_{reg} = \frac{SS_{reg}}{df_R}$   $\rightarrow$  一定是所有自变量系数均为0. 全回归方程不显著  
 对检验  $H_0: \beta_1 = \dots = \beta_k = 0$ . 则  $df_R = k$ .  $df_e = n - k - 1$ . F统计量:  $F = \frac{MS_{reg}}{MSE} \sim F_{k, n-k-1}(\alpha)$

表1 回归平方和  $SS_{reg}$

来源	d.f.	SS	MS	F
关于所有X的回归	$\delta_1$ k	$\delta_2$ $STT - RSS$	$\delta_3$ $MS_{reg}$	$\delta_4 = \frac{MS_{reg}}{MSE} = \frac{SS_{reg}/(p-1)}{RSS/(n-p)}$
残差	$\delta_5$ n-p	$\delta_6$ (n-p) MSE	4816	
总的	$\delta_7$ n-1	216571		

记回归均方为  $MS_{reg} = SS_{reg}/(p-1)$ , 则上述讨论可归结为如下表

方差分析表 (ANOVA)

来源	d.f.	SS	MS	F
X的回归	p-1	$SS_{reg}$	$\frac{SS_{reg}}{p-1}$	$\frac{MS_{reg}}{MS_{reg}/\hat{\sigma}^2}$
残差	n-p	RSS	$\frac{RSS}{n-p}$	
总的	n-1	SYT		

$F = \frac{(RSS_H - RSS) / m}{RSS / (n-p)} \rightarrow$  原检验统计量

$\frac{(RSS_H - RSS) / m}{RSS / (n-p)}$

除去原检验部分, p为总自变量数,  $\alpha = \beta_0 + \beta_1 x + e$   
 $\gamma = \beta_0 + e$  中  $p=2$

如:  $\gamma = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + e$  总为  $\gamma = \beta_0 + \dots + \beta_r x_r + e$   
 $\gamma = \beta_0 + \beta_1 x_1 + e$

$(A\hat{\beta} - b)^T (A(X^T X)^{-1} A^T)^{-1} (A\hat{\beta} - b) / m$   
 $\sim F_{m, n-p}$

例合子为  $X_i$  (新增变量) 的  $R^2$ ,  $m=1$  为新自变量数

分母为  $X_1, X_2, \dots, X_r$  的  $R^2$ . ( $R^2_{full} = STT - SS_{reg}(z, x, y)$ )  $n-p = n-3$   $p=3$  为检验中总变量数 ( $\beta_0, \beta_1, \beta_2$ )

3. 线性约束检验统计量:

对  $\gamma = X\beta + e$ ,  $e \sim N(0, \sigma^2 I)$ ,  $H_0: A\beta = b$  的 F 统计量为  $F = \frac{(A\hat{\beta} - b)^T (A(X^T X)^{-1} A^T)^{-1} (A\hat{\beta} - b) / m}{\hat{\sigma}^2} \sim F_{m, n-p}$   $m$  为约束个数  $= r - k$   $\hat{\sigma}^2 = \frac{RSS}{n-p}$

tip: 若  $A$  为对称矩阵,  $A = (a_{11}, 0, \dots, 0)^T$ , 则  $A(X^T X)^{-1} A^T = \frac{1}{n} + m^T (X_0^T X_0)^{-1} m$   $m = (m_{11}, \dots, m_{p1})$

$d_0 = \frac{1}{n} (\sum_{i=1}^n x_i - \bar{x})^T (X_0^T X_0)^{-1} (\sum_{i=1}^n x_i - \bar{x})$

不过若  $A$  不对称矩阵, 一般用  $X_0^T X_0$  计算, 会直接得出, 不存在此情况.  $A\hat{\beta} - b \rightarrow A\hat{\beta}_0 - b$ .  $\hat{\beta}_0$  为去掉  $\beta_0$  的  $(\hat{\beta}_1, \dots)$   
 $(X^T X)^{-1} \rightarrow (X_0^T X_0)^{-1}$

4.  $d_0 = \frac{1}{n} (\sum_{i=1}^n x_i - \bar{x})^T (X_0^T X_0)^{-1} (\sum_{i=1}^n x_i - \bar{x})$   $d_0 = \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{XX}}$

$\frac{1}{n} (\sum_{i=1}^n x_i - \bar{x})^T (X_0^T X_0)^{-1} (\sum_{i=1}^n x_i - \bar{x})$

$\sigma^2$  未知, 如预测:  $\hat{y}_0 \pm \hat{\sigma} \sqrt{1 + d_0} \cdot t_{n-p}(\frac{\alpha}{2})$   $E(y_0) = \hat{y}_0 \pm \hat{\sigma} \sqrt{d_0} \cdot t_{n-p}(\frac{\alpha}{2})$   
 $\hat{y}_0 \pm \hat{\sigma} \sqrt{1 + d_0} \cdot t_{n-p}(\frac{\alpha}{2})$   $\hat{\beta}_0 \pm \hat{\sigma} \sqrt{C_{jj}} \cdot t_{n-p}(\frac{\alpha}{2})$

5. 共线性: 若  $\lambda_1 \approx 0$ , 可判断有共线性.  $V_i = \frac{1}{\lambda_i} > 1000$  则有强共线性. 记标准自变量为  $z_i = \frac{x_i - m_i}{s_i}$

对具体共线性关系: 若  $\hat{\alpha} = (a_1, \dots, a_m)$ , 则为  $a_1 z_1 + \dots + a_m z_m = 0$  (可去除  $< 0.01$  的  $a_i$  部分)

6. 方差表中 t 检验 =  $\frac{\text{系数}}{\text{标准误}}$

7. 岭参数与岭回归方程下计算:

若有中心化的 OLS 系数 (原系数)  $\hat{\beta}_0$ , 则有  $\hat{\beta}_k = \hat{\beta}_0 (s_1, s_2, \dots)^T$ . 有主成分坐标下  $\hat{\alpha} = \sum \hat{\beta}_k$ .  $\Phi$  为特征向量矩阵

由 Hoerl-Kennard 公式:  $\hat{k} = \frac{\hat{\sigma}^2}{\max_i \hat{\alpha}_i^2}$  可得  $\hat{k}$ . 则岭估计在主成分坐标下为  $\hat{\alpha}(k) = (1 + k_2)^{-1} \hat{\alpha}$  ( $1$  为  $\text{diag}(\lambda_1, \dots)$  特征值对角阵)

岭参数  $\sum_{i=1}^k \frac{\lambda_i}{\lambda_i + k}$

对  $\bar{y} = \bar{\beta}_0 + \bar{m}^T \hat{\beta}_0$ , 有岭回归方程:  $\bar{y} = \bar{\beta}_0(k) + \bar{m}^T \hat{\beta}_0(k)$   $(s_1, \dots, s_p)^T$

对  $\bar{y} = \bar{\beta}_0 + \bar{m}^T \hat{\beta}_0$ , 有岭回归方程:  $\bar{y} = \bar{\beta}_0(k) + \bar{m}^T \hat{\beta}_0(k)$

岭回归误差的 RSS 的期望:  $OLS: RSS(0)$   $岭: RSS(k)$   $RSS(k) - RSS(0) = k^2 \sum_{i=1}^p \frac{\lambda_i}{(\lambda_i + k)^2} \hat{\alpha}_i^2$  附  $\lambda$  即可

8. 主成分回归方程: 先找解释几乎全部方差的主成分, 前  $r$  个成分  $\hat{\beta}_{r,p} = \tilde{\beta}_{r,p} \Phi (s_1, \dots, s_r)^T$

则  $\bar{y} = \bar{\beta}_0 + \bar{m}^T \hat{\beta}_{r,p}$

若含变量  $m$  个主成分

9.  $SS_{reg} = (X_0^T X_0)^{-1} X_0^T Y$   $\Rightarrow$  对于子集  $S$ , 有  $G_S$  为  $X_0^T X_0$  对应在  $S$  上主成分矩阵,  $V_S$  为  $X_0^T Y$  对应在  $S$  上

$R^2$  损失:  $S^T T = SS_{reg} + RSS$ ,  $R^2 \hat{\beta}_k^2 = \frac{RSS_p - RSS}{S^T T}$   $RSS_p - RSS = \sum \lambda_m \hat{\alpha}_m^2$

系数估计向量长度平方:  $OLS: \|\hat{\beta}\|^2$ , 主:  $\|\hat{\beta}_{r,p}\|^2 \Rightarrow \sqrt{\|\hat{\beta}_0\|^2 - \|\hat{\beta}_{r,p}\|^2}$

模型公式:  $\hat{y}$  为自变量  $x_i$  (学成地)  $(\hat{\alpha}^2 = \frac{RSS}{n-p}$  对应在全模型) (LCP 不是子集 SCP 是全模型 CP)

$RMS_e = \frac{RSS}{n-p}$   $C_p = \frac{RSS}{\hat{\sigma}^2} - n + 2p$ ,  $AIC = n \ln RSS_p + 2p$ ,  $BIC = n \ln RSS_p + p \ln n$

引进线性回归模型 (1) 的典则形式:  
 $\bar{y} = \alpha_0 \bar{1}_n + Z \bar{\alpha} + \bar{e}$ ,  $E(\bar{e}) = 0$ ,  $Cov(\bar{e}) = \sigma^2 I$  (3)  
 其中设计阵  $Z = X_0 \Phi$ , 而  $\bar{\alpha} = \Phi^T \hat{\beta}_0$  称为典则回归系数. 这里  $\Phi = (\phi_1, \dots, \phi_{p-1})$ ,  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_{p-1})$ , 其中  $\lambda_1 \geq \dots \geq \lambda_{p-1}$  为  $X_0^T X_0$  的特征根,  $\phi_1, \dots, \phi_{p-1}$  为对应的标准正交特征向量.

$\hat{\beta}_0 \rightarrow \hat{\beta}_k \rightarrow \hat{\alpha} = \sum \hat{\beta}_k \rightarrow \hat{F} = \frac{\hat{\sigma}^2}{\max_i \hat{\alpha}_i}$   
 $\hat{\alpha}(k) = (k_2 + \Lambda)^{-1} \hat{\alpha}$   
 $\hat{\beta}_k(k) = \sum \hat{\alpha}(k)$   
 $\hat{\beta}_0(k) = \hat{\beta}_0(k) \Phi (s_1, \dots)$   
 $\bar{y} = \hat{\beta}_0 + \bar{m}^T \hat{\beta}_0(k)$   
 $\hat{y} = \bar{y} + \bar{m}^T \hat{\beta}_0(k) (x_i - m_i) + \dots$

$\sum \lambda_m \frac{k^2}{(\lambda_m + k)^2} \hat{\alpha}_m^2$  岭

10. 岭回归法

$RSS_{\beta} = \frac{RSS_{OLS}}{n-q}$ 
 $RSS_{\beta} = n + 2q$ 
 $n - RSS_{\beta} + 2q$ 
 $\dots + 2q/n$

方差分析表			
方差来源	自由度	平方和	均方
回归	3	73.506	24.502
残差	66	282.873	4.286

  

系数表			
变量	系数	标准误差	t检验
常数项	-0.070	0.251	-0.28
FAM	1.101	1.411	0.78
PEER	2.322	1.481	1.57
SCHOOL	-2.281	2.220	-1.03

Table: 4.1.1

问题 考察学校资源设施(SCHOOL)这个变量在模型(1)中的作用, 检验其是否有存在的必要, 检验水平 $\alpha = 0.05$ .

- 问题等价于检验 $H_0: \beta_3 = 0$ , 检验统计量

$$|t_3| = \frac{|\hat{\beta}_3|}{\sqrt{\widehat{var}(\hat{\beta}_3)}} = \frac{-2.281}{\sqrt{2.220}} = 1.03$$

$$< t_{66}(0.05/2) = 1.997$$

因此接受 $H_0$ , 即SCHOOL这一变量在模型(1)中不存在存在的必要。  
注 这里 $\widehat{var}$ 指代样本方差, 从而 $\widehat{var}(\hat{\beta}_3) = c_{33}\hat{\sigma}^2$ , 其中 $c_{33}$ 是校正叉积矩阵的逆 $(X'X)^{-1}$ 的第三个对角元。

• 从而

$$MSE(\hat{\beta}_s) = \sigma^2 \text{tr}((X_s'X_s)^{-1}) = \sigma^2 \sum_{j=1}^{p-1} \frac{1}{\lambda_j}$$

• 由此可得

- 如果 $X_s'X_s$ 至少有一个特征根非常小, 即非常接近于零, 那么 $MSE(\hat{\beta}_s)$ 就会很大。从均方误差的标准来看, 这时的最小二乘估计 $\hat{\beta}$ 就不是一个好的估计。

注意 这一点和Gauss-Markov定理并无抵触, 因为Gauss-Markov定理仅仅保证了最小二乘(OLS)估计在线性无偏估计类中的方差最小性。但在 $X_s'X_s$ 至少有一个特征根很小时, 这个最小的方差值本身却很大, 因而导致了很大的均方误差。

- 对最小二乘估计根长的影响: 由

$$E\|\hat{\beta}_s\|^2 = \|\hat{\beta}_s\|^2 + MSE(\hat{\beta}_s) = \|\hat{\beta}_s\|^2 + \sigma^2 \sum_{j=1}^{p-1} \frac{1}{\lambda_j}$$

可知, 当 $X_s'X_s$ 至少有一个特征根很小时, OLS估计 $\hat{\beta}_s$ 的长度平均说来要比真正的未知向量 $\beta_s$ 的长度大得多。这就导致了 $\hat{\beta}$ 的某些分量的绝对值太大。

总结 当 $X_s'X_s$ 至少有一个特征值很小时, 最小二乘(OLS)估计不再是一个好的估计。

Definition (4.3.2)

记 $\hat{\beta}_{s,j}(k)$ 为 $\hat{\beta}_s(k)$ 的第 $j$ 个分量, 它是 $k$ 的一元函数, 当 $k$ 在 $(0, \infty)$ 上变化时,  $\hat{\beta}_{s,j}(k)$ 的图形称为岭迹 (ridge trace)。

- 岭迹法选择 $k$ 的具体操作:
  - 将 $\hat{\beta}_{s,1}(k), \dots, \hat{\beta}_{s,p-1}(k)$ 的岭迹画在同一图上;
  - 根据岭迹的变化趋势选择 $k$ 值, 使得各个回归系数的岭迹大体上稳定, 并且各个回归系数岭估计值的符号比较合理;
  - 同时选取的 $k$ 尽量接近零, 使得残差平方和不要上升太多。

总结 岭回归方法可对线性回归方程通常的最小二乘(OLS)估计的稳定性质提供一个判断的工具:

- 在严重复共线性的情况下, 数据的一个微小的变动(扰动)会造成回归系数估计的很大的变动, 而岭回归能揭示这种现象。
- 相对于OLS估计来说, 岭回归估计具有抵抗数据扰动的稳健性。
- 岭回归的稳健性定义: 它不受数据中微小的变动的影响。
- 相对于OLS估计, 岭估计具有更小的均方误差, 因此:
  - 回归系数的岭估计更倾向于接近回归系数的真实值;
  - 对于不在数据集中的预测变量的值, 预测其相应的响应的变量的值时, 用岭回归方法来预测将获得更精确地预报结果。

指标	OLS ( $k=0$ )	最优 $k$ 的岭回归	$k$ 持续增大
训练RSS	全局最小值	变大	持续单调增大
参数偏差	0 (无偏)	轻微 $>0$	持续单调飙升
参数方差	最大、易爆炸	大幅缩小	持续单调减小
参数总MSE	基准很高	显著变小	先降后升, 过大 $k$ 后反弹

一句话终极总结  
岭回归本质是主动牺牲一点点训练拟合精度 (RSS变大)、换取参数稳定性, 解决OLS多重共线性方差爆炸的问题, 最终实现整体估计误差MSE下降、泛化预测能力变强。

$C_p$ 准则选择的标准之一:  $C_p$ 与 $q$ 的靠近程度。

图形法进行变量子集的选择 ( $C_p$ 图):

- 对于每一个变量子集, 在二维图上画一个点 $(q, C_p)$ ,
- 在图上画一条直线 $C_p(q) = q$ , 即第一象限平分线;
- 选点 $(q, C_p)$ , 最接近第一象限平分线且 $C_p$ 值最小的的那个点所对应的自变量子集就是建立回归模型的最优子集。

- 问题: 当自变量个数很多时, 计算所有可能子集回归的计算量很大。
- 解决: 不需要计算所有可能子集回归的变量选择算法——逐步回归算法。
- 基本思想:

- 将变量一个一个引入, 引入变量的条件是其偏回归平方和经检验是显著的。
- 同时, 每引入一个新的变量, 对已入选变量的老变量逐个进行检验, 将检验认为不显著的变量剔除, 以保证所得自变量子集中的每一个变量都是显著的。
- 此过程经若干步直到不能再引入新变量, 也没有已入选变量的老变量需要剔除为止。

注 经逐步回归后, 最终, 回归方程中所有自变量对因变量 $Y$ 都是显著的, 而在回归方程中的变量对 $Y$ 经检验均不显著。

逐步回归法选择变量的两个基本步骤:

- 引入新变量到回归方程中;
- 从回归方程中剔除经检验不显著的变量。

一 利用主成分回归处理预测变量之间的非正交性, 这种非正交性导致模型的共线性。

- 利用主成分的方法可以减缓数据中的共线性现象。
- 其方法就是用部分主成分去解释预测变量的变动。
- 若主成分回归中使用了全部主成分作为解释变量, 那么主成分回归就等价于通常的最小二乘法。

二 主成分回归与回归系数 $\hat{\beta}_s$ 之间的约束条件有关。

岭回归:

$$RSS_R(k) = RSS_{OLS} + \sum_{j=1}^{p-1} \frac{k^2 \lambda_j}{(\lambda_j + k)^2} \hat{\alpha}_j^2$$

主成分回归:

$$RSS_P(r) = RSS_{OLS} + \sum_{j=r+1}^{p-1} \lambda_j \hat{\alpha}_j^2$$

岭回归是对所有主成分方向做连续收缩; 主成分回归是直接舍去后面的主成分方向。因此二者的RSS增量分别来自:

岭回归:  $\sum_{j=1}^{p-1} \frac{k^2 \lambda_j}{(\lambda_j + k)^2} \hat{\alpha}_j^2$

主成分回归:  $\sum_{j=r+1}^{p-1} \lambda_j \hat{\alpha}_j^2$

二者差别在于:

1. 岭回归不删除任何主成分, 而是把第 $j$ 个方向按

$$\frac{\lambda_j}{\lambda_j + k}$$

进行连续收缩;

- 主成分回归直接保留前 $r$ 个主成分, 并把后面的主成分系数估计为0;
- 岭回归的偏差来自所有方向的收缩;
- 主成分回归的偏差只来自被舍去的主成分方向。