

温冰如. 10205501432. 统计方法与机器学习作业 2.

1. 4.1: 设转换后的系数为 $\tilde{\beta}_0, \tilde{\beta}_1$.

$$\text{则 } Q(\tilde{\beta}_0, \tilde{\beta}_1) = \sum_{i=1}^n \left(\frac{y_i - c_1}{d_1} - \tilde{\beta}_0 - \tilde{\beta}_1 \frac{x_i - c_2}{d_2} \right)^2.$$

$$\begin{cases} \frac{\partial Q}{\partial \tilde{\beta}_0} = -2 \sum_{i=1}^n \left(\frac{y_i - c_1}{d_1} - \tilde{\beta}_0 - \tilde{\beta}_1 \frac{x_i - c_2}{d_2} \right) \\ \frac{\partial Q}{\partial \tilde{\beta}_1} = -2 \sum_{i=1}^n \left(\frac{y_i - c_1}{d_1} - \tilde{\beta}_0 - \tilde{\beta}_1 \frac{x_i - c_2}{d_2} \right) \frac{x_i - c_2}{d_2} \end{cases}$$

$$\text{令 } \begin{cases} \frac{\partial Q}{\partial \tilde{\beta}_0} = 0 \\ \frac{\partial Q}{\partial \tilde{\beta}_1} = 0 \end{cases} \text{ 解得: } \begin{cases} \tilde{\beta}_0 = \frac{\bar{y} - c_1}{d_1} - \frac{\bar{x} - c_2}{d_1} \frac{l_{xy}}{l_{xx}} \\ \tilde{\beta}_1 = \frac{d_2}{d_1} \frac{l_{xy}}{l_{xx}} \end{cases}$$

$$SS_T = \sum_{i=1}^n \left(\frac{y_i - c_1}{d_1} - \frac{\bar{y} - c_1}{d_1} \right)^2$$

$$= \frac{1}{d_1^2} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{d_1^2} SS_T$$

$$SS_R = \sum_{i=1}^n \left(\tilde{\beta}_0 + \tilde{\beta}_1 \frac{x_i - c_2}{d_2} - \tilde{\beta}_0 - \tilde{\beta}_1 \frac{\bar{x} - c_2}{d_2} \right)^2$$

$$= \sum_{i=1}^n \left(\tilde{\beta}_1 \left(\frac{x_i - c_2}{d_2} - \frac{\bar{x} - c_2}{d_2} \right) \right)^2$$

$$= \frac{1}{d_1^2} \tilde{\beta}_1^2 l_{xx} = \frac{1}{d_1^2} SS_R$$

$$\text{又: } \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2$$

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

$$\text{故 } SS_T = SS_R + SS_E$$

$$\text{故 } SS_E = \frac{1}{d_1^2} SS_E$$

由于调整后的 SS_T, SS_R, SS_E 均乘以调整前 $\frac{1}{d_1^2}$ 倍,

故调整后的检验统计量 $F_0 = SS_R(n-2)/SS_E$

直接把转换后的值代入
课件中的计算过程



$$= \frac{SSR}{SSE / (n-2)} = F_0 \text{ 没有变化,}$$

故转换前后统计量的值保持不变.

2. 证: 已知, 在线性回归模型 $\hat{y} = a + bx$ 中,

$$\hat{b} = \frac{l_{xy}}{l_{xx}}, \quad \hat{a} = \bar{y} - \hat{b}\bar{x} = \bar{y} - \frac{l_{xy}}{l_{xx}}\bar{x}$$

则在线性回归模型 $\hat{x} = c + dy$ 中

$$\hat{d} = \frac{l_{xy}}{l_{yy}}, \quad \hat{c} = \bar{x} - \hat{d}\bar{y} = \bar{x} - \frac{l_{xy}}{l_{yy}}\bar{y}$$

$$\text{故 } \hat{x} = \bar{x} - \frac{l_{xy}}{l_{yy}}\bar{y} + \frac{l_{xy}}{l_{yy}}y$$

将其与 $y = \bar{y} - \frac{l_{xy}}{l_{xx}}\bar{x} + \frac{l_{xy}}{l_{xx}}x$ 联立, 可得

$$\begin{cases} \hat{x} = \bar{x} \\ y = \bar{y} \end{cases} \text{ 故这两条直线有交点 } (\bar{x}, \bar{y}).$$

直接代入最简估计.

$$3. \text{证: } (I-H)^T = [I - X(X^T X)^{-1}X^T]^T$$

$$= I - X^T [X(X^T X)^{-1}]^T X$$

$$= I - X^T [X^T X]^{-1} X^T$$

$$= I - X(X^T X)^{-1}X^T = I - H$$

故对称性成立.

$$(I-H)^2 = (I-H)(I-H)$$

$$= I - 2H + H^2$$

$$= I - 2X(X^T X)^{-1}X^T + X(X^T X)^{-1}X^T X(X^T X)^{-1}X^T$$

$$= I - 2X(X^T X)^{-1}X^T + X(X^T X)^{-1}X^T$$

$$= I - X(X^T X)^{-1}X^T = I - H$$

故 ~~界等性~~ $I-H$ 是对称界等矩阵.

特征值 0.1



扫描全能王创建

证明数学作业2/Q2

由于 $\text{tr}(H) = p+1$, 故 $\text{tr}(I-H) = n-p-1$.

下证: 对称幂等矩阵特征值只可能是0或1.

考虑对称幂等矩阵 P . $P^2 = P$.

设 $P\alpha = \lambda\alpha$. 则 $P^2\alpha = \lambda P\alpha = \lambda^2\alpha = \lambda\alpha$

故 $\lambda^2 = \lambda$. $\lambda = 0$ 或 1 .

故 $I-H$ 的特征值只能是0或1.

又: $\text{tr}(I-H) = n-p-1$.

故 $(I-H) \sim \text{diag}(\underbrace{1, \dots, 1}_{n-p-1 \text{ 个}}, \underbrace{0, \dots, 0}_{p+1 \text{ 个}})$

由于相似矩阵有相同的秩. $r(I-H) = n-p-1$.

4. 证: $\sum_{i=1}^n (y_i - \hat{y}_i)$

$$= \sum_{i=1}^n [x_i' (\beta - \hat{\beta}) + \varepsilon_i]$$

$$= \sum_{i=1}^n [x_i' (\beta - (X'X)^{-1} X'Y) + \varepsilon_i]$$

$$\frac{X'X\beta = X'Y}{\sum_{i=1}^n} \sum_{i=1}^n [x_i' (\beta - \boxed{(X'X)^{-1} X'X} \beta) + \varepsilon_i]$$

$$= \sum_{i=1}^n \varepsilon_i = 0.$$

问题: 未做2又偏差之和为0

5. 证: 令 $L = \text{diag} \left\{ \frac{1}{\sqrt{L_{11}}}, \dots, \frac{1}{\sqrt{L_{pp}}} \right\}$

则 X_c 为中心化后的 X .

设 $X_s = X_c L$.

$$\begin{aligned} \beta_{s-c, \text{slope}} &= (X_s' X_s)^{-1} X_s' y^* \\ &= (L' X_c' X_c L)^{-1} L' X_c' y^* \end{aligned}$$



$$= L^{-1} (X_c' X_c)^{-1} X_c' y^*$$

$$\text{而 } \hat{\beta}_{c, \text{slope}} = \hat{\beta}_{\text{slope}} = (X_c' X_c)^{-1} X_c' y^*$$

$$\text{故 } \hat{\beta}_{s.c., \text{slope}} = L^{-1} \hat{\beta}_{\text{slope}}$$

$$\hat{\beta} = \begin{pmatrix} 0 \\ 1 \\ L^{-1} \hat{\beta}_{\text{slope}} \end{pmatrix}$$

$$E(\hat{\beta}) = \begin{pmatrix} 0 \\ 1 \\ L^{-1} E(\hat{\beta}_{\text{slope}}) \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ L^{-1} \beta_{\text{slope}} \end{pmatrix}$$

$$= (0, \sqrt{L} \beta_1, \dots, \sqrt{L} \beta_p)$$

6. 证: 在单因素方差分析模型中, $\sum_{i=1}^a \alpha_i = 0$.

$$y_{ijk} = \mu + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_{a-1} x_{a-1} + \epsilon_{ij}$$

当 y 在第 a 个水平下, 则 $x_1 = \dots = x_{a-1} = -1$

其余情况下, 取到哪个水平, 哪个水平对应的 x 取 1, 其余 x 取 0.

$$y = \begin{pmatrix} y_{1j} \\ y_{2j} \\ \vdots \\ y_{aj} \end{pmatrix} \quad \beta = \begin{pmatrix} \mu \\ \alpha_1 \\ \vdots \\ \alpha_{a-1} \end{pmatrix} \quad X = \begin{pmatrix} 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & -1 & -1 & \dots & -1 \end{pmatrix}$$

$$\epsilon = \begin{pmatrix} \epsilon_{1j} \\ \epsilon_{2j} \\ \vdots \\ \epsilon_{aj} \end{pmatrix}$$

相当于用来自 a 个不同水平的响应变量值及其对应水平下的 α 的取值来估计 μ 和诸 α_i .
 $\epsilon_{ij} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$

下面用最小二乘法求 μ 和诸 α_i 的估计:



$$L = \sum_i \sum_j (y_{ij} - \mu - \alpha_1 x_1 - \dots - \alpha_{a-1} x_{a-1})^2$$

$$\begin{cases} \frac{\partial L}{\partial \mu} = -2 \sum_i \sum_j (y_{ij} - \mu - \alpha_1 x_1 - \dots - \alpha_{a-1} x_{a-1}) = 0 \\ \frac{\partial L}{\partial \alpha_i} = -2 \sum_i \sum_j (y_{ij} - \mu - \alpha_1 x_1 - \dots - \alpha_{a-1} x_{a-1}) x_i = 0 \end{cases}$$

$$\text{解得: } \begin{cases} \hat{\mu} = \frac{\sum_{i=1}^a \bar{y}_i}{a} \\ \hat{\alpha}_i = \bar{y}_i - \frac{\sum_{i=1}^a \bar{y}_i}{a} \end{cases}$$

下面对该线性回归模型进行显著性检验。

$$SST = \sum_i \sum_j (y_{ij} - \bar{y}_i)^2$$

$$\hat{y} = \sum_{i=1}^a \frac{1}{a} \mu + \alpha_i = \bar{y}_i$$

$$SSR = \sum_i \sum_j (\bar{y}_i - \bar{y})^2$$

$$SSE = \sum_i \sum_j (y_{ij} - \bar{y}_i)^2$$

样本量 $n = am$. 变量数 $p = a - 1$

$$F_0 = \frac{SSR/p}{SSE/(n-p-1)} \sim F(p, n-p-1)$$

而 $F(p, n-p-1)$ 相当于 $F(a-1, n-a)$

与单因素方差分析中检验 $\alpha_1 = \dots = \alpha_a = 0$ 检验统计量是完全相同。

综上：单因素方差分析模型可看作多元线性回归模型。

显著分析： $\alpha_i = 0?$ \Rightarrow 因子是否显著

相当于回归分析： α_i 作系数

看回归方程是否显著



其实 y, β, X, ε 应该写成这样:

$$y = \begin{pmatrix} y_{11} \\ \vdots \\ y_{im} \\ y_{a1} \\ \vdots \\ y_{a1} \\ \vdots \\ y_{am} \end{pmatrix} \quad X = \begin{pmatrix} 1 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & -1 & \dots & \dots & -1 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & -1 & \dots & \dots & -1 \end{pmatrix}$$

$$\beta = \begin{pmatrix} \mu \\ \alpha_1 \\ \vdots \\ \alpha_{a-1} \end{pmatrix} \quad \varepsilon = \begin{pmatrix} \varepsilon_{11} \\ \vdots \\ \varepsilon_{im} \\ \vdots \\ \varepsilon_{a1} \\ \vdots \\ \varepsilon_{am} \end{pmatrix}$$

这样可以使所有的 y_{ij} 都参与到 μ 和诸 α_i 的估计中。

