

一	二	三	四	五	六	七	八	总分

(请在框内进行答题, 请标清题号)

$$-1. \begin{cases} (a-1)(b-1) = 2 \\ ab(m-1) = 6 \\ abm-1 = 11 \\ B \text{ 的自由度为 } 1 \end{cases} \Rightarrow \begin{cases} a = 3 \\ b = 2 \\ m = 2 \end{cases}$$

$$\text{故 } SSA = 0.0833 \times (3-1) = 0.1666$$

2. 由上可知, A 自由度为  $a-1 = 2$ .

3. 水平数为 2.

$$4. MSE = \frac{10}{6} = \frac{5}{3}$$

$$5. m = 2$$

$$6. \text{ 由于 } \frac{SSAB}{\sigma^2} \sim \chi^2((a-1)(b-1))$$

$$\text{故 } F_{AB} = \frac{SSAB / ((a-1)(b-1))}{SSE / (ab(m-1))} \sim F((a-1)(b-1), ab(m-1))$$

故 p 值为  $P(F > F_{AB})$  其中  $F \sim F((a-1)(b-1), ab(m-1))$

当 p 值小于 0.05, 我们拒绝原假设, 认为交互效应显著.

否则, 我们接受原假设, 认为交互效应不显著.

$$7. \text{ 证: } SS_T = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^m (y_{ijk} - \bar{y}_{...})^2$$

$$= \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^m ((y_{ijk} - \bar{y}_{ij.}) + (\bar{y}_{ij.} - \bar{y}_{i..}) + (\bar{y}_{i..} - \bar{y}_{...}) + (\bar{y}_{j..} - \bar{y}_{...}) + (\bar{y}_{...} - \bar{y}_{...}))^2$$

(请在框内进行答题, 请标清题号)

欲证  $SS_T = SS_A + SS_B + SS_{AB} + SS_E$ . 只需证上述平方和展开式交叉项为 0.

$$\text{令 } A = y_{ijk} - \bar{y}_{ij\cdot}, B = \bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot\cdot\cdot}, C = \bar{y}_{\cdot j\cdot} - \bar{y}_{\cdot\cdot\cdot}$$

$$D = y_{ij\cdot} - \bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot j\cdot} + \bar{y}_{\cdot\cdot\cdot}$$

$$\begin{aligned} "AB": & \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^m (y_{ijk} - \bar{y}_{ij\cdot})(\bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot\cdot\cdot}) \\ &= \sum_{i=1}^a (\bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot\cdot\cdot}) \sum_{j=1}^b \sum_{k=1}^m (y_{ijk} - \bar{y}_{ij\cdot}) \end{aligned}$$

由于  $\sum_{k=1}^m (y_{ijk} - \bar{y}_{ij\cdot}) = 0$  故该项为 0.

$$"AC": \text{同理: } \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^m (y_{ijk} - \bar{y}_{ij\cdot})(\bar{y}_{\cdot j\cdot} - \bar{y}_{\cdot\cdot\cdot}) = 0$$

$$\begin{aligned} "AD": & \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^m (y_{ijk} - \bar{y}_{ij\cdot})(y_{ij\cdot} - \bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot j\cdot} + \bar{y}_{\cdot\cdot\cdot}) \\ &= \sum_{i=1}^a \sum_{j=1}^b (y_{ij\cdot} - \bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot j\cdot} + \bar{y}_{\cdot\cdot\cdot}) \sum_{k=1}^m (y_{ijk} - \bar{y}_{ij\cdot}) \end{aligned}$$

由于  $\sum_{k=1}^m (y_{ijk} - \bar{y}_{ij\cdot}) = 0$  故该项为 0.

$$\begin{aligned} "BC": & \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^m (\bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot\cdot\cdot})(\bar{y}_{\cdot j\cdot} - \bar{y}_{\cdot\cdot\cdot}) \\ &= m \sum_{i=1}^a (\bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot\cdot\cdot}) \sum_{j=1}^b (\bar{y}_{\cdot j\cdot} - \bar{y}_{\cdot\cdot\cdot}) \end{aligned}$$

由于  $\sum_{j=1}^b (\bar{y}_{\cdot j\cdot} - \bar{y}_{\cdot\cdot\cdot}) = 0$  故该项为 0.

$$\begin{aligned} "BD": & \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^m (\bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot\cdot\cdot})(y_{ij\cdot} - \bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot j\cdot} + \bar{y}_{\cdot\cdot\cdot}) \\ &= m \sum_{i=1}^a (\bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot\cdot\cdot}) \sum_{j=1}^b (y_{ij\cdot} - \bar{y}_{i\cdot\cdot} - \bar{y}_{\cdot j\cdot} + \bar{y}_{\cdot\cdot\cdot}) \end{aligned}$$

(请在框内进行答题, 请标清题号)

由于  $\sum_{j=1}^a (y_{ij} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}_{...}) = 0$ , 故该项为 0.

同理,  $\sum_{i=1}^a \sum_{j=1}^m \sum_{k=1}^m (y_{ij} - \bar{y}_{i..}) (y_{ij} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}_{...}) = 0$

故所有交叉项均为 0. 故  $SS_T = SS_A + SS_B + SS_{AB} + SS_E$ .

二. / 1  $y = 0.6954 + 1.6034x_1 + \varepsilon$ ,  $\varepsilon \sim N(0, \sigma^2)$

由于 t 检验的 p 值为 0.000, 故我们认为该模型显著 (取  $\alpha = 0.05$ )

2.  $x_1 = 0.5$  时,  $\hat{y} = 0.6954 + 1.6034 \times 0.5 = 1.4971$

现在线性回归模型中,  $\hat{y}_0 \sim N(x_0' \beta, \sigma^2 x_0' (X'X)^{-1} x_0)$

而  $y \sim N(x_0' \beta, \sigma^2)$

故  $y - \hat{y}_0 \sim N(0, \sigma^2 (1 + x_0' (X'X)^{-1} x_0))$ .

又:  $\sigma^2$  未知, 故用估计值  $\hat{\sigma}^2 = \frac{SSE}{n-p-1}$  代替.

故  $t = \frac{y - \hat{y}_0}{\sqrt{(1 + x_0' (X'X)^{-1} x_0) SSE / (n-p-1)}} \sim t(n-p-1)$

用  $t(n-p-1)$  的分位数  $t_{1-\alpha/2}(n-p-1)$  给出预测区间

3. 这个结论普遍存在.

$$R^2 = \frac{SSE}{SST} = \frac{y'(I - H)y}{y'(I - H)y} = \frac{y'(I - H)y}{y'(I - H)y}$$

= ...



(请在框内进行答题, 请标清题号)

4. 由3中的结论: 放入模型的自变量越多则  $R^2$  越大, 则若以  $R^2$  作为模型选择标准, 当选模型正确, 该标准也会让我们选用全模型, 从而造成过拟合, 故不应以  $R^2$  作为标准。

改进方案: 用  $\sigma = \frac{SSE}{n-p-1}$  作为标准。

~~当自变量过多, 则  $\frac{SSE}{n-p-1}$  可作为罚项,~~

~~当自变量过多,  $SSE$  要小~~

三. 对  $X'X$  特征分解:  $X'X = U'\Lambda U$ , 其中  $U$  是对称正交阵,  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_{p+1})$ ,  $\lambda_1 \geq \dots \geq \lambda_{p+1}$  是  $X'X$  的特征值。

当数据集中出现多重共线性,  $|X'X| \approx 0$ 。

故又  $X'X\alpha = \lambda\alpha$ , 故  $X'X$  必然有比较接近于 0 的特征值。

取  $k_j = \sqrt{\frac{\lambda_1}{\lambda_j}}$ , 当  $\lambda_j$  接近于 0, 条件数  $k_j$  会很大, 当某个  $k_j$  大于某个特定的值, 就认为该数据集存在多重共线性。

四. 感知机损失函数为  $\sum_{x \in W} -y_i(w \cdot x_i + b)$

其中  $W$  是误分类点的集合, 只要所有点均正确分类, 损失函数即为 0。

线性 SVM 的损失函数为合页损失函数

$\sum_{i=1}^n \max(0, 1 - y_i(w \cdot x_i + b))$  其中  $[w_+] = \begin{cases} w, & w > 0 \\ 0, & w \leq 0. \end{cases} \sum_i \geq 0$ 。

(请在框内进行答题, 请标清题号)

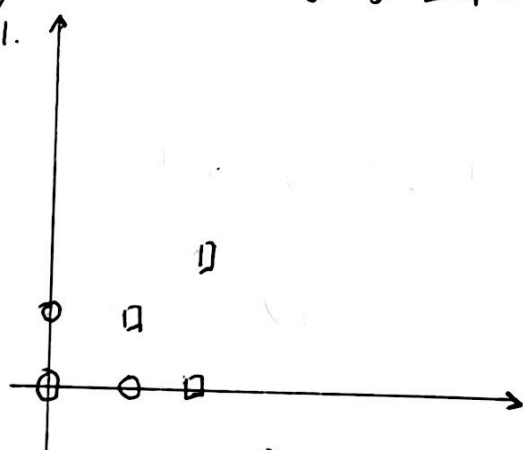
也就是说, 在线性 SVM 中, 不仅要完全分类正确, 所有点离分类超平面的距离要不小于 1 (即最大程度地把正负实例点分开), 损失函数的值为 0, 5/1. 生成式模型就是在给定的数据集在某种准则下逐步生成出来的模型;

判别式模型就是在给定的数据集学习出来, 并判断未知点所属类别的模型

不同在于判别式模型无显式的学习过程, 而生成式模型是通过多次的迭代逐渐生成出来的

2. 判别模型: kNN, 朴素贝叶斯, 逻辑斯蒂回归,  
生成模型: 支持向量机, 决策树.

六. 1.



□ 正例

○ 负例

显然, 数据集线性可分, 构造最优化问题

$$\min_w \frac{1}{2} \|w\|^2$$

$$\text{s.t. } y_i (w x_i + b) - 1 \geq 0, \quad i = 1, \dots, N$$

$$\text{即: } \min_{w_1, w_2} w_1^2 + w_2^2$$

$$\text{s.t. } w_1 + w_2 + b - 1 \geq 0$$

(请在框内进行答题，请标清题号)

$$2w_1 + 2w_2 + b - 1 \geq 0$$

$$2w_1 + b - 1 \geq 0$$

$$-b - 1 \geq 0$$

$$-w_1 - b - 1 \geq 0$$

$$-w_2 - b - 1 \geq 0$$

约束条件化简后得：

$$\begin{cases} w_1 + w_2 - 2 \geq 0 \\ w_1 + 2w_2 - 2 \geq 0 \\ 2w_1 - w_2 - 2 \geq 0 \\ w_2 \geq 2 \\ w_1 \geq 2 \\ 2w_1 + w_2 - 2 \geq 0 \end{cases}$$

解得： $w_1 = w_2 = 2$ .2.  $(1,1)$   $(2,0)$   $(0,1)$   $(1,0)$  都是支撑向量,

3. 构造对偶优化问题:

$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (x_i x_j) - \sum_{i=1}^n \alpha_i$$

$$\text{s.t.} \sum_{i=1}^n \alpha_i y_i = 0 \quad i=1, \dots, n.$$

$$\alpha_i \geq 0$$

$$\begin{aligned} \text{Eq} \min_{\alpha} & \frac{1}{2} (4\alpha_1\alpha_2 + 2\alpha_1\alpha_3 + 4\alpha_2\alpha_3 + 2\alpha_1^2 + 8\alpha_2^2 + 4\alpha_3^2 \\ & + \alpha_5^2 + \alpha_6^2 + 6\alpha_1\alpha_5 + 2\alpha_2\alpha_5 + 2\alpha_3\alpha_5 + \alpha_1\alpha_6 \\ & + 2\alpha_2\alpha_6) - \alpha_1 - \alpha_2 - \alpha_3 - \alpha_4 - \alpha_5 - \alpha_6 \\ \text{s.t.} & \sum_{i=1}^n \alpha_i + \alpha_2 + \alpha_3 - \alpha_4 - \alpha_5 - \alpha_6 = 0 \\ & \alpha_i \geq 0. \end{aligned}$$

(请在框内进行答题, 请标清题号)

令拉格朗日乘子为  $\lambda$ , 则有

$$L = \frac{1}{2} \alpha_1^2 + \alpha_2^2 + \alpha_3^2 + \alpha_4^2 + \alpha_5^2 + \alpha_6^2 + \lambda(-1)$$

$$\begin{cases} \frac{\partial L}{\partial \alpha_1} = 2\alpha_1 = 0 \\ \frac{\partial L}{\partial \alpha_2} = 2\alpha_2 = 0 \\ \frac{\partial L}{\partial \alpha_3} = 2\alpha_3 = 0 \\ \frac{\partial L}{\partial \alpha_4} = 2\alpha_4 = 0 \\ \frac{\partial L}{\partial \alpha_5} = 2\alpha_5 = 0 \\ \frac{\partial L}{\partial \alpha_6} = 2\alpha_6 = 0 \\ \frac{\partial L}{\partial \lambda} = -1 = 0 \end{cases}$$

$$\begin{cases} \alpha_1 = 0 \\ \alpha_2 = 0 \\ \alpha_3 = 0 \\ \alpha_4 = 0 \\ \alpha_5 = 0 \\ \alpha_6 = 0 \\ \lambda = -1 \end{cases}$$



课 程：《统计方法与机器学习》

页数：第(8)页 / 8页

学生姓名：温北和

学号：10205101432

专业：数据科学与大数据技术

(请在框内进行答题，请标清题号)