

<p>均值: 平均水平, 稳健性</p> <p>中位数: 中心位置, 稳健性</p> <p>三均值 = $\frac{1}{4}(\text{上四分位数} + \text{下四分位数}) + \frac{1}{2}\text{中位数}$</p> <p>差: 有量纲</p> <p>极差, 四分位极差 $R_i = Q_3 - Q_1$</p> <p>分散性</p>	<p>经验分布与位置</p> <p>$X_{[np]+1}$, NP 非整数.</p> <p>$\frac{1}{2}(X_{(np)} + X_{(np+1)})$ 整数.</p> <p>变异系数: 相对分散性</p> <p>无量纲</p> <p>$CV = \frac{S}{\bar{x}}$ 越大波动越大</p> <p>上下截断法 (异常值)</p> <p>$Q_1 - 1.5R, Q_3 + 1.5R$</p> <p>下四分位数 上四分位数 四分位极差</p> <p>形状</p> <p>偏度: < 0 左偏, > 0 右偏 右侧更分散</p> <p>$G_1 = \frac{\mu_3}{\sigma^3} - 3$</p> <p>$\mu_3$: 三阶中心矩</p>	<p>高维</p> <p>$(x_1, \dots, x_p)^T$ 高维变量</p> <p>数据为 $(x_1, \dots, x_p)^T$ $(x_1, \dots, x_p)^T$ $(x_1, \dots, x_p)^T$</p> <p>$x = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix}$</p> <p>$S = \begin{bmatrix} S_{11} & S_{12} & \cdots & S_{1p} \\ S_{21} & S_{22} & \cdots & S_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ S_{p1} & S_{p2} & \cdots & S_{pp} \end{bmatrix} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$</p> <p>$s_{ijk} = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)$</p> <p>残差: 列残差矩阵</p> <p>$r_{ijk} = \frac{s_{ijk}}{\sqrt{S_{jj}} \sqrt{S_{kk}}}$</p> <p>Pearson 相关矩阵</p> <p>$R = \{r_{ijk}\}_{n \times p}$, $R = D^{-1} S D^{-1}, D = \text{Diag}(\sqrt{S_{11}}, \dots, \sqrt{S_{pp}})$.</p> <p>斯皮尔曼相关</p> <p>$Q = \{q_{ijk}\}$</p> <p>标准化后</p> <p>$R = \frac{1}{n-1} \sum_{i=1}^n x_i x_i^T = \frac{1}{n-1} (X^*)^T X^{**}$</p> <p>单体的数学特征</p> <p>$\Sigma = \text{COV}(X) = E((X - \mu)(X - \mu)^T)$</p> <p>$\widehat{\Sigma}_{jk} = \text{COV}(X_j, X_k) = E(X_j - \mu_j)(X_k - \mu_k)^T$</p> <p>$\rho_{jk} = \frac{\widehat{\Sigma}_{jk}}{\widehat{\Sigma}_{jj} \widehat{\Sigma}_{kk}}$ 相关系数</p> <p>Y=f(X_1, \dots, X_p)</p> <p>Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p</p> <p>Y = X\beta + \varepsilon</p> <p>估计</p> <p>$\hat{\beta} = (X^T X)^{-1} X^T Y$</p> <p>全: $(I - H)Y$</p> <p>SSE: $\varepsilon^T (I - H)\varepsilon$</p> <p>E(SSE): $\sigma^2(n-p)$</p> <p>$\hat{\sigma}^2$: $\frac{SSE}{n-p} = \frac{1}{n-p} Y^T (I - H)Y$</p> <p>惯性:</p> <p>① 全, B 无偏.</p> <p>② $\hat{\beta} \sim N(\beta, \sigma^2(X^T X))$</p> <p>$\hat{\beta} = (X^T X)^{-1} X^T (X\beta + \varepsilon)$</p> <p>$\text{cov}(\hat{\beta}) = E[(X^T X)^{-1}]^T \Sigma \varepsilon$</p> <p>③ $\frac{SSE}{\sigma^2} \sim \frac{n-p}{\sigma^2}$ 分布</p> <p>$SSE = \varepsilon^T (I - H)\varepsilon$</p> <p>$P(\varepsilon \perp \varepsilon) \sim N(0, \sigma^2)$</p> <p>④ $\hat{\sigma}^2$ 与 $\hat{\beta}$ 独立</p> <p>对称性.</p> <p>\Leftrightarrow 验证 $(X^T X)^{-1} X^T$</p>
--	---	--

註記: HYS I-HYR

直方圖: 組距 \sim 歷史 累積頻數. $\frac{1}{n} \sum_{i=1}^k I(x_i \leq x)$

QQ圖: (X, Y) 相同名稱.

正态QQ (理論路徑). $Y = \sigma X + \mu$.

線性關係 \Leftrightarrow 直線.

χ^2 檢驗 \rightarrow 單尾檢驗

$$\chi^2 = \sum_{i=1}^k \frac{(m_i - np_i)^2}{np_i} \sim \chi^2(k-h-1), k \text{為分布中待估參數}$$

Reject for large $H_0: F_X = F_0$.

m_i : 観測值. np_i : 理論頻數

經驗分布擬合指標

正态性 V 檢驗

x_1, \dots, x_m 觀測
 $\underbrace{x_1, \dots, x_{(k)}}_{d_1} \xrightarrow{\text{分佈}} \underbrace{x_{(k+1)}, \dots, x_m}_{d_2}$ 次序
 $W = \left(\frac{k}{\sum_{i=1}^k d_i} \alpha_i d_i \right)^2$

$W = \frac{\sum_{i=1}^k (x_{(i)} - \bar{x})^2}{\sum_{i=1}^k (x_{(i)} - \bar{x})^2}$

$W_k = \frac{1}{\prod_{i=1}^k \Gamma(k_i)}$

隨機向量性質.

$E(\bar{X}) = A E X = A \mu$.

$Cov(AX) = ACov(X)A^T = A \Sigma A^T$

$E(c^T X) = c^T \mu. \quad \text{Var}(c^T X) = c^T \Sigma c$

$\text{Cov}(AX, BY) = A \text{Cov}(XY) B^T$

$\text{Cov}(XX) = E[(X - E[X])(X - E[X])^T]$

$\text{Cov}(c^T X, d^T Y) = c^T \text{Cov}(X, Y) d$

註: $c^T \mu = c^T$
 $\text{Var}(a^T y) = \sigma^2 \|a\|^2$
 $E(c^T \theta) = c^T \theta$
 $\text{Cov}(c^T \theta) = \sigma^2 c^T (X' X)^{-1} c$
Gauss-Markov 定理
 註: $a^T X = c^T$
 $\text{Var}(a^T y) = \sigma^2 \|a\|^2$
 $\|a\|^2 = \|a - X(X^T X)^{-1} c + X(X^T X)^{-1} c\|^2$
 $= \|a - X(X^T X)^{-1} c\|^2 +$
 $+ 2c^T (X(X^T X)^{-1} X^T c)$
 $= \|a - X(X^T X)^{-1} c\|^2$
 故 $\text{Var}(a^T y) \geq \text{Var}(a^T y)$
 當且僅當 $a = X(X^T X)^{-1} X^T c$
 $a^T y =$

$$S = \begin{bmatrix} S_{xx} & S_{xy} \\ S_{yx} & S_{yy} \end{bmatrix} \quad S_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$S_{xy}^2 \leq S_{xx} S_{yy}$$

1) $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p + \varepsilon$
 $\varepsilon \sim N(0, \sigma^2 I)$

$\hat{Y} = f(X(X^T X)^{-1} X^T) Y$

$SST = \sum_{i=1}^n (y_i - \bar{y})^2 = Y^T (I - \frac{1}{n} J) Y$ $J = 1 \mathbf{1}^T$
 $SSE = \sum_{i=1}^n (\hat{y}_i - y_i)^2 = Y^T (I - H) Y$
 $SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = Y^T (H - \frac{1}{n} J) Y$ $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$

$SST = SSE + SSR$
 $R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$

全局显著性
 $F = \frac{SSR/(p-1)}{SSE/(n-p)} \sim F(p-1, n-p)$.
 Reject For large F .

残差分析

	DF	SS	MS	F ₀	P值
R	p-1	$\sum (\hat{y}_i - \bar{y})^2$	$SSR/p-1$	MSR	P_0
E	n-p	$\sum (y_i - \hat{y}_i)^2$	$SSE/n-p$		
T	n-1	$\sum (y_i - \bar{y})^2$			

$\varepsilon = Y^T (I - H) Y = \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 = \hat{\sigma}^2$
 p). 证明 取加权
 $H(Y)$

$\hat{\beta}_k = \hat{\beta}_0 + \hat{\beta}_1 x_k$ $(X^T X)^{-1}$ 矩阵上乘以 x_k
 $t_k = \frac{\hat{\beta}_k - \beta_0}{\sqrt{\hat{\sigma}^2}} \sim t(n-p)$
 $\beta_0 = 2Pf + t(n-p) > t(n-p)$ 有偏
 $\hat{\beta}_k \pm \hat{\sigma} \sqrt{t_k} t_{n-p}^{-1/2}$

置信区间
预测 对 y 预测
 $\hat{y}_0 \sim N(\hat{\beta}_0, \sigma^2 \hat{X}^T (X^T X)^{-1} \hat{X})$
 $\hat{y}_0 \sim N(\hat{\beta}_0, \sigma^2 \hat{X}^T (X^T X)^{-1} \hat{X})$, studentized
 $\text{Plug-in } \hat{y}_0$

$\mu_1 = \mu_2 = \mu_3$ $MS_{M1} = MS_{M2}$
 $MS_{M1} > MS_{M2}$
因素: 对 y 影响的效应量, 不同水平
因子效应:
交互作用

模型
 $\eta_{ij} = \mu + \delta_i + \varepsilon_{ij}$, $i = 1, 2, \dots, n_1$, $j = 1, 2, \dots, n_2$
 $\varepsilon_{ij} \sim N(0, \sigma^2)$ 假设
 $\sum_{j=1}^n \eta_{ij} = 0$.

$\bar{\eta}_i = \frac{1}{n_2} \sum_{j=1}^n \eta_{ij}$,
 $\bar{\varepsilon}_i = \frac{1}{n_2} \sum_{j=1}^n \varepsilon_{ij}$,
 $\bar{\mu} = \frac{1}{n} \sum_{i=1}^n \bar{\eta}_i$,
 $\bar{\eta}_i = \frac{1}{n_2} \sum_{j=1}^n \eta_{ij}$,
 $\bar{\varepsilon}_i = \frac{1}{n_2} \sum_{j=1}^n \varepsilon_{ij}$,
 $\bar{y} = \mu + \bar{\varepsilon}$

$SST = \sum_{i=1}^n \sum_{j=1}^{n_i} (\eta_{ij} - \bar{y})^2$ 总平方和
 $SSE = \sum_{i=1}^n \sum_{j=1}^{n_i} (\hat{\eta}_{ij} - \bar{y})^2$ 残差平方和
 $SSA = \sum_{i=1}^n n_i (\bar{\eta}_i - \bar{y})^2$ 因素 A

$SST = SSE + SSA$
 $E(SSE) = \sum_{i=1}^n E(\eta_{i-1} \varepsilon_i^2) = (n-1) \sigma^2$
 $E(SSA) = E\left(\sum_{i=1}^n n_i (\bar{\eta}_i - \bar{y})^2\right) = (n-1) \sigma^2 +$

LOGISTIC	$\pi_k = \frac{\exp(\beta_0 + \beta_1 x_1 + \dots + \beta_m x_m)}{1 + \exp(\beta_0 + \beta_1 x_1 + \dots + \beta_m x_m)}$	$\mu_{ij} = \mu - (\mu_i - \mu_j)$
因 (x_0, \dots, x_{m-1}) 為 m 項取值 $x_k = (x_{k,0}, \dots, x_{k,m-1})^T$	$\pi_k = \frac{\exp(\beta_0 + \beta_1 x_{k,1} + \dots + \beta_m x_{k,m})}{1 + \exp(\beta_0 + \beta_1 x_{k,1} + \dots + \beta_m x_{k,m})}$	$= (\mu_{ij} - \mu) - (\mu_i + \beta_j)$
對 A (Bernoulli 分佈 Y) 進行 凡得之迴歸		
$\begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} \sim \text{Bin}(n_k, \pi_k)$	$\begin{pmatrix} 1 & x_{1,1} & x_{1,2} & \dots & x_{1,m-1} \\ 1 & x_{2,1} & x_{2,2} & \dots & x_{2,m-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & x_{n,2} & \dots & x_{n,m-1} \end{pmatrix} \sim \tilde{x}_k^T$	表:
	$N = \begin{bmatrix} n_1 \pi_1(x_1) \\ n_2 \pi_2(x_2) \\ \vdots \\ n_m \pi_m(x_m) \end{bmatrix}$	
$\tilde{x}_k^T \beta$ 為量 $\tilde{x}_k^T \beta$, $Y=1$ 的概率		系源 自由度 平方和
$\ln\left(\frac{\pi(x_k)}{1-\pi(x_k)}\right) = \tilde{x}_k^T \beta$, $k=1, \dots, m$		A $a-1$ SS_A B $b-1$ SS_B
$\pi(x_k) = \frac{\exp(\tilde{x}_k^T \beta)}{1 + \exp(\tilde{x}_k^T \beta)}$, $k=1, 2, \dots, m$		之五 $(a-1)(b-1)$ SS_{AB}

$$\begin{aligned} \text{误差.} & ab(c-1) & S_{\text{E}} \\ X^T N = X^T Y & abc-1 & S_T \\ W = \text{Diag} [m\pi(x_0)\left(1-\pi(x_0)\right), \dots, m\pi(x_m)\left(1-\pi(x_m)\right)] & \text{总和} & \\ I(\theta) = X^T W X & S_{\text{A}} = b \sum_{i=1}^a (\bar{y}_{i..} - \bar{y})^2 \end{aligned}$$

$$\begin{aligned} \hat{\beta}^{(t+1)} &= \hat{\beta}^{(t)} - \left[I(\hat{\beta}) \right]^{-1} \left(X^T (Y - N^{(t)}) \right) \\ \text{推斷} \quad H_0: \beta_1 = \dots = \beta_p = 0. \\ \textcircled{1} \quad \text{至模型} \quad X_1, \dots, X_{p-1} \quad \text{與 SS 檢驗} \\ K^2 &= 2 \ln \frac{L(\hat{\beta}; Y_1, \dots, Y_m)}{L(\hat{\beta}_{H_0}, Y_1, \dots, Y_m)} \sim \chi^2(r) \\ \text{reject for large } K^2 \\ \textcircled{2} \quad \text{Wald 檢驗} \\ \hat{\beta}_j = 0, \quad \rightarrow \text{當對系數 } \hat{\beta}_j \text{ 是 } S(\hat{\beta}) \\ \hat{\beta} - \hat{\beta}_j \sim N(0, I(\hat{\beta})) \\ w_j = \left(\frac{\hat{\beta}_j}{S(\hat{\beta}_j)} \right)^2 \sim \chi^2(1) \\ SS_B &= a \sum_{j=1}^k (\bar{y}_{\cdot j} - \bar{y})^2 \\ SS_{AB} &= c \sum_{i=1}^m \sum_{j=1}^k (\bar{y}_{ij} - \bar{y}_{\cdot i} - \bar{y}_{\cdot j} + \bar{y})^2 \\ SS_E &= \sum_{i=1}^m \sum_{j=1}^k \sum_{h=1}^c (y_{ijh} - \bar{y}_{ij\cdot})^2 \end{aligned}$$

主成分分析

$$\text{cov}(X) = \Sigma, \quad Y = a^T X.$$

$$\begin{array}{l} \text{max}_{a_i} \quad \text{Var}(Y_i) = a_i^T \Sigma a_i \\ \text{s.t.} \quad a_i^T a_i = 1 \end{array} \quad \Rightarrow \quad Y_i = a_i^T X$$

$$\begin{array}{l} \text{max}_{a_k} \quad a_k^T \Sigma a_k \\ \text{s.t.} \quad a_k^T a_k = 1, \quad a_k^T \Sigma a_i = 0 \quad (i=1, 2, \dots, k-1) \end{array} \quad \Rightarrow \quad Y_k = a_k^T X$$

式 $\sum \lambda_i$ 特征值与相关系数化特征向量,

$$\lambda_1 > \lambda_2 > \dots > \lambda_p > 0.$$

$$e_1, e_2, \dots, e_p$$

$$Y_k = e_k^T X, \quad \text{Var}(Y_k) = e_k^T \Sigma e_k = \lambda_k.$$

$$\text{cov}(Y_k, Y_l) = e_k^T \Sigma e_l = 0 \quad \text{cov}(Y) = \text{Diag}(\lambda_1, \dots, \lambda_p)$$

总方差 $\sum_{k=1}^p \text{Var}(Y_k) = \sum_{k=1}^p \lambda_k = \text{tr}(\Sigma) = \sum_{i,j} \text{cov}(x_i, x_j)$

贡献率 $\frac{\lambda_k}{\sum_{i=1}^p \lambda_i} = \frac{\text{Var}(Y_k)}{\sum_{i=1}^p \text{Var}(Y_i)}$ 对应成分得贡献率

主成分 Y_i 与相关系数.

$$Y = P^T X \Rightarrow X = P Y$$

$$Y_i = e_{ij} Y_1 + e_{ij} Y_2 + \dots + e_{ij} Y_p \Rightarrow \text{cov}(Y_i, Y_j) = \lambda_i e_{ij}$$

主成分系数

$$\rho(Y_i, Y_j) = \frac{\text{cov}(Y_i, Y_j)}{\sqrt{\text{Var}(Y_i)} \sqrt{\text{Var}(Y_j)}} = \frac{\lambda_i e_{ij}}{\sqrt{\lambda_i} \sqrt{\lambda_j}} = \frac{\lambda_i}{\sqrt{\lambda_j}} e_{ij}$$

$$\sum_{i=1}^p \rho^2(Y_i, Y_j) = \sum_{i=1}^p \frac{\lambda_i e_{ij}^2}{\lambda_j} = 1$$

因 $P^T \Sigma P = \Lambda \Rightarrow \Sigma = P \Lambda P^T \quad P = e_1, e_2, \dots, e_p$

$\hat{e}_{ij} = (e_{ij}, e_{2j}, \dots, e_{pj})^T$ 且 $\sum_{i=1}^p e_{ij} = e_{jj}$
 现在对 因变量的贡献率
 $\eta_j^{(m)}$ 表示对 因变量 X_j 的贡献率 $\eta_j^{(m)}$ 定义为
 $\eta_j^{(m)} = \sum_{i=1}^m p_i^m (Y_i, X_j) = \sum_{i=1}^m \frac{\Delta_i e_{ij}}{e_{jj}}$
标准化
 $\chi_k^* = \frac{y_k - \bar{y}_k}{\sqrt{\text{Var}(y_k)}}$
 $\rho = \text{cov}(\chi^*)$ 相关系数
 相关性矩阵 \Rightarrow 对称矩阵解
 $\sum_{k=1}^p \text{Var}(\chi_k^*) = \sum_{k=1}^p \lambda_k^* = \sum_{k=1}^p \text{Var}(X_k^*) = p$
 贡献率 $\frac{\lambda_k}{p}$ 各是等
样本协方差
 $S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$ 代替 Σ ,
 $\rho = (r_{jk})_{p \times p} = \left(\frac{s_{jk}}{\sqrt{s_{jj}s_{kk}}} \right)_{p \times p}$ 代替 ρ .
 样本相关
 $y_{ijk} = \underbrace{\mu + \alpha_i + \beta_j + \gamma_{ij}}_{\text{总效应}} + \varepsilon_{ijk}$,
 $i=1, 2, \dots, a$
 $j=1, 2, \dots, b$
 $k=1, 2, \dots, c$
 $\rightarrow B_j$ 效应
 $\sum_{i=1}^a \alpha_i = 0 \quad \sum_{j=1}^b \beta_j = 0$.
 $\sum_{i=1}^a \gamma_{ij} = \sum_{j=1}^b \gamma_{ij} = 0$
结果:
 将 $X = (X_1, \dots, X_p)^T$ 的 n 个测值 (x_1, \dots, x_n) , $x_i =$

均方	F值	P值
MSA	MS_A/MSE	P_A
MSB	MS_B/MSE	P_B
MSAB	MS_{AB}/MSE	P_{AB}

依次代入。
 $y_k = \hat{e}_k^T X$, $k=1, 2, \dots, p$.
 得 y_{1k}, \dots, y_{nk} . 简称主成分的得分
 $y_{1k} = \hat{e}_k^T X_1, \dots, y_{nk} = \hat{e}_k^T X_n$.

得新的样本矩阵.
 $\hat{e}_k^T S \hat{e}_k = \lambda_k$ $k=1, \dots, p$.
 $\hat{e}_j^T S \hat{e}_k = 0$ $1 \leq j \neq k \leq p$.

原数据		统计量
$x_{11}, x_{12}, \dots, x_{1p}$	$y_{11}, y_{12}, \dots, y_{1p}$	
$x_{21}, x_{22}, \dots, x_{2p}$	$y_{21}, y_{22}, \dots, y_{2p}$	
\vdots	\vdots	\vdots
$x_{n1}, x_{n2}, \dots, x_{np}$	y_{n1}, \dots, y_{np}	

判别分析

典型相关分析

Y_1, \dots, Y_p 与 X_1, \dots, X_p 间相关系数

找 $U_i = a_i^T X, V_i = b_i^T Y$ 使 $a_i^T b_i$

U_i, V_i 最大可能取原相关性，各对变量提取的相关性不重叠

模型

$X = (X_1, \dots, X_p)^T, Y = (Y_1, \dots, Y_p)^T$ 假设

$(X, Y) = (X_1, \dots, X_p, Y_1, \dots, Y_p)^T$

$\sum = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$

$\Sigma_{11} = \text{cov}(X), \Sigma_{22} = \text{cov}(Y)$

$\Sigma_{12} = \Sigma_{21} = \text{cov}(X, Y)$ 请续 (18.24)

$U_i = a_i^T X, \text{Var}(U_i) = a_i^T \Sigma_{11} a_i$

$V_i = b_i^T Y, \text{Var}(V_i) = b_i^T \Sigma_{22} b_i$

$\text{cov}(U_i, V_i) = a_i^T \Sigma_{12} b_i$

$\rho_{U_i, V_i} = \frac{a_i^T \Sigma_{12} b_i}{\sqrt{a_i^T \Sigma_{11} a_i} \sqrt{b_i^T \Sigma_{22} b_i}}$

判别问题

$\max_{U_i, V_i} \rho_{U_i, V_i} = a_i^T \Sigma_{12} b_i, \quad \left\{ \begin{array}{l} \max_{U_i, V_i} P_{U_i, V_i} = a_i^T \Sigma_{12} b_i \\ a_i^T \Sigma_{11} a_i = 1 \end{array} \right.$

$\left. \begin{array}{l} \max_{U_i, V_i} P_{U_i, V_i} = a_i^T \Sigma_{12} b_i \\ a_i^T \Sigma_{11} a_i = 1 \end{array} \right. \quad \left. \begin{array}{l} \max_{U_i, V_i} P_{U_i, V_i} = a_i^T \Sigma_{12} b_i \\ a_i^T \Sigma_{11} a_i = 1 \end{array} \right.$

判别条件

$X = (X_1, \dots, X_p)^T$

C_i : 分布 $F_i(\bar{x}) = F_i(x_1, \dots, x_p)$

$f_i(\bar{x})$

找规则，使得某在某种规则下最佳。

距离

欧式 距离 没有距离

马氏:

$\left\{ \begin{array}{l} x, y, \text{来自 } \mu, \Sigma \text{ 的条件} \text{ 为 } \\ x \in \mathbb{R}^n \\ d^2(x, y) = (x - y)^T \Sigma^{-1} (x - y) \\ x \in \mathbb{C}^n \\ d^2(x, \mu) = (x - \mu)^T \Sigma^{-1} (x - \mu) \end{array} \right.$

μ_1, μ_2, Σ 的关系

$d^2(\mu_1, \mu_2) = (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2)$

判别规则

特征值分解
 $A = \sum_{ii}^{-\frac{1}{2}} \sum_{12} \sum_{22}^{-1} \sum_{21} \sum_{11}^{-\frac{1}{2}}$, $B = \sum_{22}^{-\frac{1}{2}} \sum_{11} \sum_{11}^{-1} \sum_{12} \sum_{22}^{-\frac{1}{2}}$

特征值 $\rho_1^2 \geq \rho_2^2 \geq \dots \geq \rho_p^2$ 特征向量 e_1, e_2, \dots, e_p .
 f_1, f_2, \dots, f_p

$U_k = e_k^T \sum_{ii}^{-\frac{1}{2}} X$, $V_k = f_k^T \sum_{22}^{-\frac{1}{2}} Y$.

$P_{uk|vk} = \rho_k$, $k=1, 2, \dots, p$

4.5. $b_i^T \sum_{ii} b_i = 1$
 $\text{cov}(u_i, u_j) = \text{cov}(V_i, V_j) = 0, \forall i \neq j$
 $\text{cov}(u_i, v_j) = \text{cov}(V_i, U_j) = 0$

容易判断:
 $\begin{cases} x \in G_1, & d(x, G_1) \leq d(x, G_2) \\ x \notin G_1, & d(x, G_1) > d(x, G_2) \end{cases}$

归到最小的子空间分离法.

$\sum_{i=1}^n = \sum_{k=1}^p k$:
 $\sum_{i=1}^n = \frac{1}{n-k} \left[\sum_{i=1}^k (n_i - 1) S_i \right]$

$W_i = G_i^T x + b$, $a_i = \sum_{j=1}^m y_{ij}$

二次就算, $d^2(x, G_i)$ 取最小的.

判别准则:
 p_i : 样本属于 G_i 的概率 (类)

$p^* = P(C \text{ 将 } x \text{ 判错}) = p_i P(C(1, R))$

$P(C(j|i, R))$: R 将属于 G_i 的样本 j

1. 回代法.
 $n = n_1 + n_2$ 训练 \rightarrow 判别.

	189			
实验	C ₁	C ₂	Slt	
C ₁	n ₁₁	n ₁₂	n ₁	
C ₂	n ₂₁	n ₂₂	n ₂	
Slt	\hat{n}_1	\hat{n}_2	n	

2. 对别脂的进行抑制.
3. 又对高分样品重量, 混合比例作出限制.

n_{12}^* : 对 G_1 依次混合.

Bayes, 裁决分析.

先验: G_1 和 G_2 出现的几率

$$P_1 = P(G_1), \quad P_2 = P(G_2)$$

$$P_1 + P_2 = 1 \quad \text{比例} \quad P_1 = \frac{n_1}{N}, \quad P_2 = \frac{n_2}{N}$$

利用 R: 空间的线性分. $\mathbb{R}^P = \{R_1, R_2\}$

$$P(C(1|R)) = \int_{R_2} f_1(x) dx. \quad \text{损失为 } C(2|1)$$

$$P(C(1|R)) = \int_{R_1} f_2(x) dx. \quad \text{损失为 } C(1|2)$$

G_1 和 G_2 的后验: $P(G_i|x) = \frac{C(C(i|x))}{P_f(x) + P_g(x)}$

$$\begin{cases} C(1|x) = C(2|1) \\ R_1 = \{x: p_f(x) > p_g(x)\} \\ R_2 = \{x: p_g(x) > p_f(x)\} \end{cases}$$

$$P^* = P_1 P(C(1|R)) + P_2 P(C(1|R)) \quad \text{证明时选择一项}$$

$\hat{x}: C(1|R) \neq C(2|R)$.

$$R_1 = \{x: C(1|R) p_f(x) > C(2|R) p_g(x)\}$$

$$R_2 = \{x: C(1|R) p_g(x) > C(2|R) p_f(x)\}$$

对正态:

损失相同

(1) 当 $\Sigma_1 = \Sigma_2 = \Sigma$ 时.

1. 例 2. 练习 3. Mink 4. Che 5. 附录 6. 方差 7. 附录

$R_1 = \{x : W_1(x) \geq W_2(x)\}$,
 $R_2 = \{x : W_1(x) < W_2(x)\}$.

$W_1(x) = a^T x + b$, $a_1 = \Sigma^{-1} \mu_1$, $b = -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \ln p_1$
 $W_2(x) = a_2^T x + b$, $a_2 = \Sigma^{-1} \mu_2$, $b = -\frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \ln p_2$

μ_1, μ_2, Σ 估计如何.

(2) 若 $\Sigma_1 \neq \Sigma_2$ 时, $C(1|2) \neq C(2|1)$.

$R_1 = \{x : -\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) - \frac{1}{2} \ln |\Sigma_1| + \ln(C(1|2)p_1)\}$
 $\geq -\frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) - \frac{1}{2} \ln |\Sigma_2| + \ln(C(1|2)p_2)\}$

$R_2 = \{\dots\}$

④ 重心距离:
 $D_{pq}^2 =$

⑤ 均差:
 $D_{pq}^2 = W_p$

⑥ 集中度:
 $D_{pq}^2 = d$

⑦ 离散度:
 $D_{pq}^2 = W_q$

⑧ 期望距离:
 $D_{pq}^2 = W_r$

⑨ 方差:
 $D_{pq}^2 = W_s$

⑩ 相似度:
 $D_{pq}^2 = W_t$

$$\widehat{\chi_p} = \frac{1}{n_p} \sum_{i=1}^{n_p} \chi_i^{(p)}, \quad \chi_1^{(p)}, \dots, \chi_{n_p}^{(p)} \text{ が } \infty$$

並列.

$$G_p, G_q, G_k \xrightarrow{\text{合算}} G_r, G_k.$$

$$\begin{cases} D_{pk} \\ D_{qk} \end{cases} \longrightarrow D_{rk}$$

$$D_{rk}^2 = \alpha_p D_{pk}^2 + \alpha_q D_{qk}^2 + \beta D_{pq}^2 + \gamma |D_{pk}^2 - D_{qk}^2|$$

$$\max_{i \in G_p, j \in G_q} \{d_{ij}\} \quad D_{rk} = \min \{D_{pk}, D_{qk}\}$$

$$\frac{1}{n_p n_q} \sum_{i \in G_p} \sum_{j \in G_q} d_{ij} \quad D_{pk} = \frac{n_q}{n_r} D_{pk} + \frac{n_q}{n_r} D_{qk}.$$

$$\frac{1}{n_p n_q} \sum_{i \in G_p} \sum_{j \in G_q} d_{ij}^2 \quad D_{rk}^2 = \frac{n_p}{n_r} D_{pk}^2 + \frac{n_q}{n_r} D_{qk}^2$$

$$(\bar{\chi}_p, \bar{\chi}_q) = \sqrt{(n_p - n_q)(\bar{\chi}_p - \bar{\chi}_q)^T (\bar{\chi}_p - \bar{\chi}_q)}$$

$$D_{rk}^2 = \frac{n_p}{n_r} D_{pk}^2 + \frac{n_q}{n_r} D_{qk}^2 - \frac{n_p n_q}{n_r n_r} D_{pq}^2$$

$$W_p - W_q = \frac{n_p n_q}{n_p + n_q} (\bar{\chi}_p - \bar{\chi}_q)^T (\bar{\chi}_p - \bar{\chi}_q)$$

$$D_{rk}^2 = \frac{n_p + n_q}{n_r + n_k} D_{pk}^2 + \frac{n_q + n_k}{n_r + n_k} D_{qk}^2 - \frac{n_k}{n_r + n_k} D_{pq}^2$$

$$\text{设 } D_{pq} = \begin{bmatrix} 0 & d_{12} & \cdots & d_{1n} \\ d_{21} & 0 & \cdots & d_{2n} \\ \vdots & \ddots & \ddots & \vdots \\ d_{n1} & d_{n2} & \cdots & 0 \end{bmatrix} \quad D_{pq} = D_{pp}$$