



第四章 鉴别分析

第一节

引言

第二节

距离鉴别法

第三节

贝叶斯 (Bayes) 鉴别法

第四节

费歇 (Fisher) 鉴别法

第五节

实例分析与计算机实现

- 在我们的日常生活和工作实践中，经常会遇到鉴别分析问题，即根据历史上划分类别的有关资料和某种最优准则，拟定一种鉴别措施，鉴定一种新的样本归属哪一类。例如，某医院有部分患有肺炎、肝炎、冠心病、糖尿病等病人的资料，统计了每个患者若干项症状指标数据。目前想利用既有的这些资料找出一种措施，使得对于一种新的病人，当测得这些症状指标数据时，能够鉴定其患有哪种病。又如，在天气预报中，我们有一段较长时间有关某地域每天气象的统计资料（晴阴雨、气温、气压、湿度等），目前想建立一种用连续五天的气象资料来预报第六天是什么天气的措施。这些问题都能够应用鉴别分析措施予以处理。



多元统计

- 把此类问题用数学语言来体现，能够论述如下：设有 n 个样本，对每个样本测得 p 项指标（变量）的数据，已知每个样本属于 k 个类别（或总体） G_1, G_2, \dots, G_k 中的某一类，且它们的分布函数分别为 $F_1(x), F_2(x), \dots, F_k(x)$ 。我们希望利用这些数据，找出一种鉴别函数，使得这一函数具有某种最优性质，能把属于不同类别的样本点尽量地域别开来，并对测得一样 p 项指标（变量）数据的一种新样本，能鉴定这个样本归属于哪一类。



多元统计

- 鉴别分析内容很丰富，措施诸多。判断分析按鉴别的总体数来区别，有两个总体鉴别分析和多总体鉴别分析；按区别不同总体所用的数学模型来分，有线性鉴别和非线性鉴别；按鉴别时所处理的变量措施不同，有逐渐鉴别和序贯鉴别等。鉴别分析能够从不同角度提出问题，所以有不同的鉴别准则，如马氏距离最小准则、Fisher准则、平均损失最小准则、最小平方准则、最大似然准则、最大约率准则等等，按鉴别准则的不同又提出多种鉴别措施。本章仅简介常用的几种鉴别分析措施：距离鉴别法、Fisher鉴别法、Bayes鉴别法和逐渐鉴别法。



一 马氏距离的概念

二 距离鉴别的思想及措施

三 鉴别分析的实质



多元统计

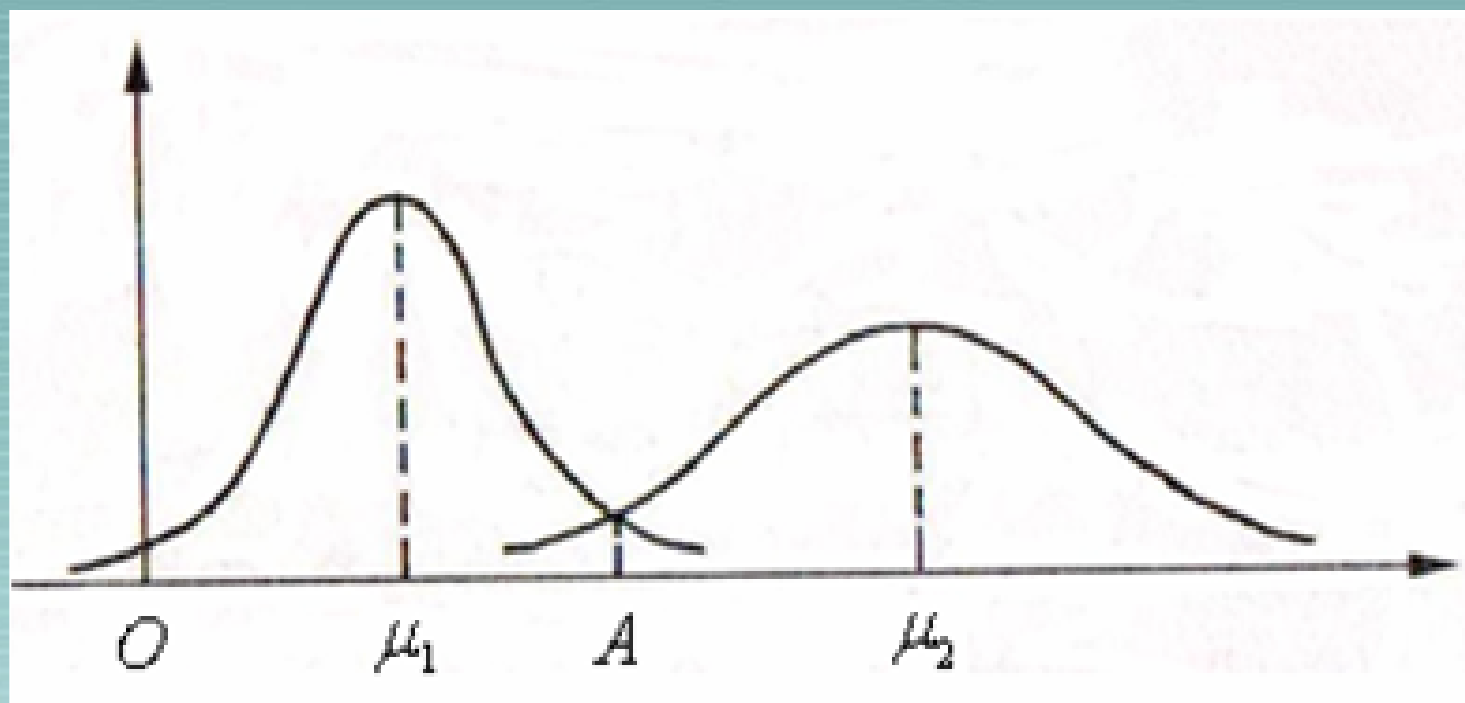


图4.1



多元统计

第二、设有量度重量和长度的两个变量 X 和 Y

$$A(0,5) \quad B(10,0) \quad C(1,0) \quad D(0,10)$$

$$AB = \sqrt{10^2 + 5^2} = \sqrt{125}$$

$$CD = \sqrt{1^2 + 10^2} = \sqrt{101}$$

$$AB = \sqrt{10^2 + 50^2} = \sqrt{2600}$$

$$CD = \sqrt{1^2 + 100^2} = \sqrt{10001}$$



多元统计

- 为此，我们引入一种由印度著名统计学家马哈拉诺比斯 (Mahalanobis, 1936) 提出的“马氏距离”的概念。

- 设 \mathbf{X} 为 p 维随机向量， \mathbf{Y} 为 G 维随机向量， $\boldsymbol{\mu}$ 为 G 维均值向量， $\boldsymbol{\Sigma}(> 0)$ 为 G 维正定协方差矩阵。

$$D^2(\mathbf{X}, \mathbf{Y}) = (\mathbf{X} - \mathbf{Y})' \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \mathbf{Y})$$

$$\mathbf{X} \quad \mathbf{Y}$$

$$D^2(\mathbf{X}, G) = (\mathbf{X} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \boldsymbol{\mu})$$

$$\boldsymbol{\Sigma} = \mathbf{I}$$



1、两个总体的距离鉴别问题

- 问题：设有协方差矩阵 Σ 相等的两个总体 G_1 和 G_2 ，其均值分别是 μ_1 和 μ_2 ，对于一种新的样品 X ，要判断它来自哪个总体。
- 一般的想法是计算新样品 X 到两个总体的马氏距离 $D^2(X, G_1)$ 和 $D^2(X, G_2)$ ，并按照如下的鉴别规则进行判断
$$\begin{cases} X \in G_1, & \text{如果 } D^2(X, G_1) \leq D^2(X, G_2) \\ X \in G_2, & \text{如果 } D^2(X, G_1) > D^2(X, G_2) \end{cases}$$
- 这个鉴别规则的等价描述为：求新样品 X 到 G_1 的距离与到 G_2 的距离之差，假如其值为正， X 属于 G_2 ；不然 X 属于 G_1 。



多元统计

■ 我们考虑

$$\begin{aligned} & D^2(\mathbf{X}, G_1) - D^2(\mathbf{X}, G_2) \\ &= (\mathbf{X} - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \boldsymbol{\mu}_1) - (\mathbf{X} - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \boldsymbol{\mu}_2) \\ &= \mathbf{X}' \boldsymbol{\Sigma}^{-1} \mathbf{X} - 2\mathbf{X}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\mu}_1' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 - (\mathbf{X}' \boldsymbol{\Sigma}^{-1} \mathbf{X} - 2\mathbf{X}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 + \boldsymbol{\mu}_2' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2) \\ &= 2\mathbf{X}' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \boldsymbol{\mu}_1' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 \\ &= 2\mathbf{X}' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ &= -2 \left(\mathbf{X} - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right)' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ &= -2(\mathbf{X} - \bar{\boldsymbol{\mu}})' \boldsymbol{\alpha} = -2\boldsymbol{\alpha}'(\mathbf{X} - \bar{\boldsymbol{\mu}}) \end{aligned}$$



多元统计

是两个总体均值的下期望。

■ 其中 $\bar{\mu} = \frac{1}{2}(\mu_1 + \mu_2)$

$$\alpha = \Sigma^{-1}(\mu_1 - \mu_2)$$

$$W(\mathbf{X}) = \alpha'(\mathbf{X} - \bar{\mu})$$

$$\begin{cases} \mathbf{X} \in G_1, & \text{如果 } W(\mathbf{X}) \geq 0 \\ \mathbf{X} \in G_2, & \text{如果 } W(\mathbf{X}) < 0 \end{cases}$$

$W(\mathbf{X})$

为两总体的判别函数，由于它是

\mathbf{X}

函数，故又称为线性判别函数。

α

称为判别向量。
在实际应用中，总体的均值和协方差矩阵一般是未知的，可由样本均值和样本协方差矩阵分别进行估计。

$\mathbf{X}_1^{(1)}, \dots, \mathbf{X}_{n_1}^{(1)}$

来自

G_1

的样本。

$\mathbf{X}_1^{(2)}, \dots, \mathbf{X}_{n_2}^{(2)}$

是来自总体

G_2

的样本。

μ_1

μ_2

和
十五章设计分类为



多元统计

$$\bar{\mathbf{X}}^{(1)} = \frac{1}{n_1} \sum_{i=1}^{n_1} \mathbf{X}_i^{(1)}$$

$$\bar{\mathbf{X}}^{(2)} = \frac{1}{n_2} \sum_{i=1}^{n_2} \mathbf{X}_i^{(2)}$$

Σ

$$\hat{\Sigma} = \frac{1}{n_1 + n_2 - 2} (\mathbf{S}_1 + \mathbf{S}_2)$$

$$\mathbf{S}_\alpha = \sum_{i=1}^{n_\alpha} (\mathbf{X}_i^{(\alpha)} - \bar{\mathbf{X}}^{(\alpha)})(\mathbf{X}_i^{(\alpha)} - \bar{\mathbf{X}}^{(\alpha)})', \quad \alpha = 1, 2$$

$$\hat{W}(\mathbf{X}) = \hat{\alpha}'(\mathbf{X} - \bar{\mathbf{X}})$$

$$\bar{\mathbf{X}} = \frac{1}{2}(\bar{\mathbf{X}}^{(1)} + \bar{\mathbf{X}}^{(2)}) \quad \hat{\alpha} = \hat{\Sigma}^{-1}(\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)})$$

$$\begin{cases} \mathbf{X} \in G_1, & \text{如果 } \hat{W}(\mathbf{X}) \geq 0 \\ \mathbf{X} \in G_2, & \text{如果 } \hat{W}(\mathbf{X}) < 0 \end{cases}$$



多元统计

■ 这里我们应该注意到:

(1) 当 $p = 1$ 时, G_1 和 G_2 部分分别为 $N(\mu_1, \sigma^2)$ 和 $N(\mu_2, \sigma^2)$. 均值为 μ_1, μ_2 , 且 $\mu_1 < \mu_2$. 系数为 $\mu_1 < \mu_2$.

$$\alpha = \frac{\mu_1 - \mu_2}{\sigma^2} < 0$$

$$W(x) = \alpha(x - \bar{\mu})$$

$$\begin{cases} x \in G_1, & \text{如果 } x \leq \bar{\mu} \\ x \in G_2, & \text{如果 } x > \bar{\mu} \end{cases}$$



多元统计

(2) 当 $\mu_1 \neq \mu_2$ 时，我们采用 (4.4) 式为判别
规则的形式，选择判别函数为 $\Sigma_1 \neq \Sigma_2$

$$\begin{aligned} W^*(\mathbf{X}) &= D^2(\mathbf{X}, G_1) - D^2(\mathbf{X}, G_2) \\ &= (\mathbf{X} - \mu_1)' \Sigma_1^{-1} (\mathbf{X} - \mu_1) - (\mathbf{X} - \mu_2)' \Sigma_2^{-1} (\mathbf{X} - \mu_2) \end{aligned}$$

\mathbf{X}

一次判别，相应的判别规则为

$$\begin{cases} \mathbf{X} \in G_1, & \text{如果 } W^*(\mathbf{X}) \leq 0 \\ \mathbf{X} \in G_2, & \text{如果 } W^*(\mathbf{X}) > 0 \end{cases}$$



多元统计

2、多种总体的距离鉴别问题

■ 问题：设有 k 个总体 G_1, G_2, \dots, G_k

$$\mu_1, \mu_2, \dots, \mu_k \quad \Sigma_1, \Sigma_2, \dots, \Sigma_k \quad \Sigma_1 = \Sigma_2 = \dots = \Sigma_k = \Sigma$$
$$\mathbf{X}$$

■

\mathbf{X}

$$\begin{aligned} D^2(\mathbf{X}, G_\alpha) &= (\mathbf{X} - \mu_\alpha)' \Sigma^{-1} (\mathbf{X} - \mu_\alpha) \\ &= \mathbf{X}' \Sigma^{-1} \mathbf{X} - 2\mu_\alpha' \Sigma^{-1} \mathbf{X} + \mu_\alpha' \Sigma^{-1} \mu_\alpha \\ &= \mathbf{X}' \Sigma^{-1} \mathbf{X} - 2(\mathbf{I}'_\alpha \mathbf{X} + C_\alpha) \end{aligned}$$

$$\mathbf{I}'_\alpha = \Sigma^{-1} \mu_\alpha \quad C_\alpha = -\frac{1}{2} \mu_\alpha' \Sigma^{-1} \mu_\alpha \quad \alpha = 1, 2, \dots, k$$



多元统计

■ 由 (4.8) 式, 可以取线性判别函数为

$$W_\alpha(\mathbf{X}) = \mathbf{I}'_\alpha \mathbf{X} + C_\alpha \quad \alpha = 1, 2, \dots, k$$

$$\mathbf{X} \in G_i \quad W_i(\mathbf{X}) = \max_{1 \leq \alpha \leq k} (\mathbf{I}'_\alpha \mathbf{X} + C_\alpha)$$

$$\mu_1, \mu_2, \dots, \mu_k \quad \Sigma$$

$$\bar{\mathbf{X}}_1^{(\alpha)}, \dots, \bar{\mathbf{X}}_{n_\alpha}^{(\alpha)} \quad G_\alpha$$

$$\alpha = 1, 2, \dots, k \quad \mu_\alpha \quad \alpha = 1, 2, \dots, k \quad \Sigma$$

$$\bar{\mathbf{X}}^{(\alpha)} = \frac{1}{n_\alpha} \sum_{i=1}^{n_\alpha} \mathbf{X}_i^{(\alpha)} \quad \alpha = 1, 2, \dots, k$$

$$\hat{\Sigma} = \frac{1}{n-k} \sum_{\alpha=1}^k \mathbf{S}_\alpha \quad n = n_1 + n_2 + \dots + n_k$$



多元统计

$$S_{\alpha} = \sum_{i=1}^{n_{\alpha}} (\mathbf{X}_i^{(\alpha)} - \bar{\mathbf{X}}^{(\alpha)}) (\mathbf{X}_i^{(\alpha)} - \bar{\mathbf{X}}^{(\alpha)})' \quad \alpha = 1, 2, \dots, k$$

■ 同样，我们注意到，如果总体 G_1, G_2, \dots, G_k

$$\bar{\Sigma}_1, \bar{\Sigma}_2, \dots, \bar{\Sigma}_k \quad \bar{\mathbf{X}}$$

$$D^2(\mathbf{X}, G_{\alpha}) = (\mathbf{X} - \bar{\boldsymbol{\mu}}_{\alpha})' \bar{\Sigma}_{\alpha}^{-1} (\mathbf{X} - \bar{\boldsymbol{\mu}}_{\alpha}) \quad \alpha = 1, 2, \dots, k$$

$$\mathbf{X} \in G_i \quad D^2(\mathbf{X}, G_i) = \min_{1 \leq \alpha \leq k} D^2(\mathbf{X}, G_{\alpha})$$

$$\bar{\boldsymbol{\mu}}_1, \bar{\boldsymbol{\mu}}_2, \dots, \bar{\boldsymbol{\mu}}_k \quad \bar{\Sigma}_1, \bar{\Sigma}_2, \dots, \bar{\Sigma}_k \quad \bar{\boldsymbol{\mu}}_{\alpha} \quad \alpha = 1, 2, \dots, k$$

$$\bar{\Sigma}_{\alpha} \quad \alpha = 1, 2, \dots, k$$

$$\hat{\Sigma}_{\alpha} = \frac{1}{n_{\alpha} - 1} S_{\alpha} \quad \alpha = 1, 2, \dots, k$$



- 我们懂得，鉴别分析就是希望利用已经测得的变量数据，找出一种鉴别函数，使得这一函数具有某种最优性质，能把属于不同类别的样本点尽量地域别开来。为了更清楚的认识鉴别分析的实质，以便能灵活的应用鉴别分析措施处理实际问题，我们有必要了解“划分”这么概念。
- 设 R_1, R_2, \dots, R_k 是 p 维空间 R^p 的 k 个子集，假如它们互不相交，且它们的和集为 R^p ，则称 R_1, R_2, \dots, R_k 为 R^p 的一种划分。



多元统计

- 在两个总体的距离判别问题中，利用 $W(\mathbf{X}) = \boldsymbol{\alpha}'(\mathbf{X} - \bar{\boldsymbol{\mu}})$

$$R^p$$

$$\begin{cases} R_1 = \{\mathbf{X} : W(\mathbf{X}) \geq 0\} \\ R_2 = \{\mathbf{X} : W(\mathbf{X}) < 0\} \end{cases}$$

$$\mathbf{X} \in R_1 \quad \mathbf{X} \in G_1 \quad \mathbf{X} \in R_2 \quad \mathbf{X} \in G_2$$

- 这么我们将会发觉，鉴别分析问题实质上就是在某种意义上，以最优的性质对 p 维空间 R^p 构造一种“划分”，这个“划分”就构成了一种鉴别规则。这一思想将在背面的各节中体现的愈加清楚。



多元统计

例 在企业的考核中，能够根据企业的生产经营情况把企业分为优异企业和一般企业。考核企业经营情况的指标有：

资金利润率=利润总额/资金占用总额

劳动生产率=总产值/职员平均人数

产品净值率=净产值/总产值

三个指标的均值向量和协方差矩阵如下。既有二个企业，观察值分别为

(7.8, 39.1, 9.6) 和 (8.1, 34.2, 6.9)，问这两个企业应该属于哪一类？



多元统计

变量	均值向量		协方差矩阵		
	■ 优	一般	■ 优	■ 一般	■ 异
资金利润率	13.5	5.4	68.39	40.24	21.41
劳动生产率	40.7	29.8	40.24	54.58	11.67
产品净值率	10.7	6.2	21.41	11.67	7.90

$$\Sigma^{-1} = \begin{bmatrix} 0.119337 & -0.02753 & -0.28276 \\ -0.02753 & 0.033129 & 0.025659 \\ -0.28276 & 0.025659 & 0.854988 \end{bmatrix}$$



多元统计

$$\mu_1 - \mu_2 = \begin{bmatrix} 8.1 \\ 10.9 \\ 4.5 \end{bmatrix} \quad (\mu_1 + \mu_2)/2 = \begin{bmatrix} 9.45 \\ 35.25 \\ 8.45 \end{bmatrix}$$

$$\text{判别函数的系数 } \Sigma^{-1}(\mu_1 - \mu_2) = \begin{bmatrix} -0.60581 \\ 0.25362 \\ 1.83679 \end{bmatrix}$$

$$\begin{aligned} & \text{判别函数的常数项 } \frac{(\mu_1 + \mu_2)'}{2} \Sigma^{-1}(\mu_1 - \mu_2) \\ &= \begin{bmatrix} 9.45 & 35.25 & 8.45 \end{bmatrix} \begin{bmatrix} -0.60581 \\ 0.25362 \\ 1.83679 \end{bmatrix} = 18.73596 \end{aligned}$$



多元统计

■ 线性鉴别函数为：

$$y = -0.60581x_1 + 0.25362x_2 + 1.83679x_3 - 18.73596$$

$$y_1 = -0.60581 \times 7.8 + 0.25362 \times 39.1 + 1.83679 \times 9.6 - 18.73596 \\ = 4.0892 > 0 \text{ (第一个新企业属于一类)}$$

$$y_2 = -0.60581 \times 8.1 + 0.25362 \times 34.2 + 1.83679 \times 6.9 - 18.73596 \\ = -2.2956 < 0 \text{ (第二个新企业属于二类)}$$



多元统计

- 错判概率：由上面的分析能够看出，马氏距离鉴别法是合理的，但是这并不意味着不会发生误判。

两总体分别服从 $N(\mu_1, \sigma^2)$ 和 $N(\mu_2, \sigma^2)$ ，
其线性判别函数为：

$$W(y) = (y - \bar{\mu}) \frac{1}{\sigma^2} (\mu_1 - \mu_2), \quad \text{其中}$$

$\bar{\mu} = \frac{\mu_1 + \mu_2}{2}$ 。不失一般性设 $\mu_1 > \mu_2$ ，这种

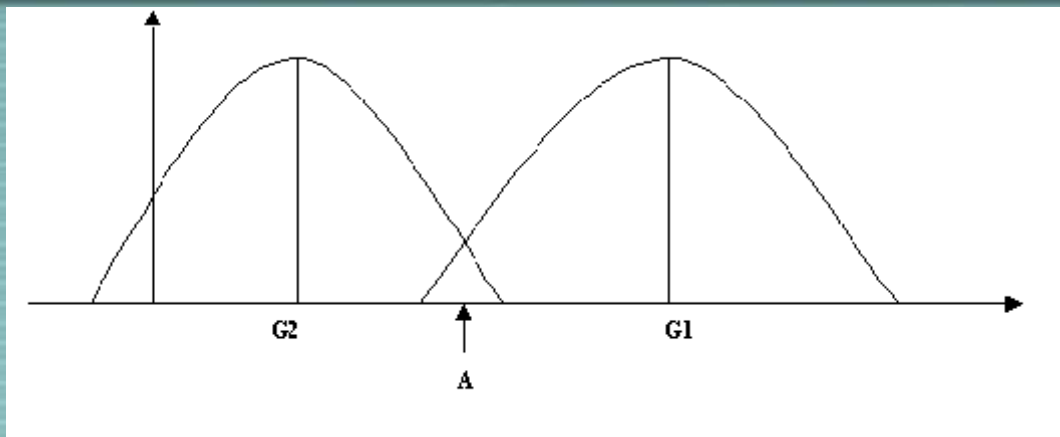
情况下线性判别函数 $W(y)$ 的符号取决于 $y > \bar{\mu}$ 还是 $y < \bar{\mu}$ 。

如果 $y > \bar{\mu}$ ，则 $y \in G_1$

如果 $y < \bar{\mu}$ ，则 $y \in G_2$ 。



多元统计



$$\text{错判概率: } P(X_2 > \bar{\mu}) = P(X_2 - \mu_2 > \frac{\mu_1 + \mu_2}{2} - \mu_2)$$

$$= P(X_2 - \mu_2 > \frac{\mu_1 - \mu_2}{2})$$

$$= 1 - \Phi\left(\frac{\mu_1 - \mu_2}{2\sigma}\right)$$



多元统计

距离鉴别只要求懂得总体的数字特征，不涉及总体的分布函数，当参数和协方差未知时，就用样本的均值和协方差矩阵来估计。距离鉴别措施简朴实用，但没有考虑到每个总体出现的机会大小，即先验概率，也没有考虑到错判的损失。贝叶斯鉴别法正是为了处理这两个问题提出的鉴别分析措施。



一 Bayes鉴别的基本思想

二 Bayes鉴别的基本措施



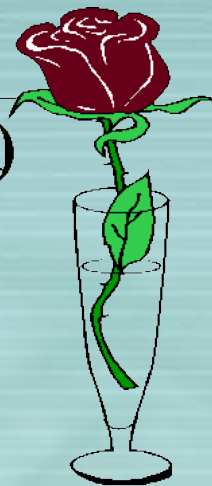
多元统计

办公室新来了一种雇员小王，小王是好人还是坏人大家都在猜测。按人们主观意识，一种人是好人或坏人的概率均为0.5。坏人总是要做坏事，好人总是做好事，偶尔也会做一件坏事，一般好人做好事的概率为0.9，坏人做好事的概率为0.2，一天，小王做了一件好事，小王是好人的概率有多大，你目前把小王判为何种人。

$$P(\text{好人}/\text{做好事})$$

$$= \frac{P(\text{好人})P(\text{做好事}/\text{好人})}{P(\text{好人})P(\text{做好事}/\text{好人}) + P(\text{坏人})P(\text{做好事}/\text{坏人})}$$

$$= \frac{0.5 \times 0.9}{0.5 \times 0.9 + 0.5 \times 0.2} = 0.82$$



多元统计

$P(\text{坏人/做好事})$

$$= \frac{P(\text{坏人})P(\text{做好事/坏人})}{P(\text{好人})P(\text{做好事/好人}) + P(\text{坏人})P(\text{做好事/坏人})}$$

$$= \frac{0.5 \times 0.2}{0.5 \times 0.9 + 0.5 \times 0.2} = 0.18$$

距离鉴别简朴直观，很实用，但是距离鉴别的措施把总体等同看待，没有考虑到总体以不同的概率（先验概率）出现，也没有考虑误判之后所造成的损失的差别。

一种好的鉴别措施，既要考虑到各个总体出现的先验概率，又要考虑到错判造成的损失，Bayes鉴别就具有这些优点，其鉴别效果愈加理想，应用也更广泛。



多元统计

- 贝叶斯公式是一种我们熟知的公式

$$P(B_i | A) = \frac{P(A | B_i)P(B_i)}{\sum P(A | B_i)P(B_i)}$$

- 贝叶斯鉴别

在各总体的概率分布及先验概率已知的前提下，分别计算待判对象属于各总体的后验概率，并以**最大后验概率**相应的总体来作为待判对象的所属总体。



■ 问题：设有 k 类 G_1, G_2, \dots, G_k

其先验概率为 q_1, q_2, \dots, q_k ，类条件概率密度函数为 $f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_k(\mathbf{x})$

其中 $q_i \geq 0$

$$\sum_{i=1}^k q_i = 1$$

G_i

误判损失为 $C(j|i)$ ， $i, j = 1, 2, \dots, k$

\mathbf{X}



多元统计

- 下面我们对这一问题进行分析。首先应该清楚 $C(i|i) = 0$

$$C(j|i) \geq 0 \quad i, j = 1, 2, \dots, k$$

$$G_1, G_2, \dots, G_k \quad p \quad R_1, R_2, \dots, R_k$$

$$\bar{R} = (R_1, R_2, \dots, R_k)$$

$$G_i$$

$$f_i(\mathbf{x})$$

$$R_j$$

$$G_j$$



多元统计

■ 故在规则 R

$$P(j | i, R) = \int_{R_j} f_i(\mathbf{x}) d\mathbf{x} \quad i, j = 1, 2, \dots, k \quad i \neq j$$

$$C(1 | i), \dots, C(i-1 | i), C(i+1 | i), \dots, C(k | i)$$

$$r(i | R) = \sum_{j=1}^k [C(j | i) P(j | i, R)] \quad i = 1, 2, \dots, k$$

$$C(i | i) = 0$$



多元统计

■ 由于 k G_1, G_2, \dots, G_k
 q_1, q_2, \dots, q_k R

$$g(R) = \sum_{i=1}^k q_i r(i, R)$$

$$= \sum_{i=1}^k q_i \sum_{j=1}^k C(j|i) P(j|i, R)$$

R_1, R_2, \dots, R_k

$g(R)$



- 设每一个总体 G_i 的概率密度函数为 $f_i(\mathbf{x})$ $i = 1, 2, \dots, k$
- 设 G_i 的先验概率为 q_i $i, j = 1, 2, \dots, k$
- 设 $C(j|i)$ 为将 G_i 误判为 G_j 的损失 $C(i|i) = 0$

$$R = (R_1, R_2, \dots, R_k)$$

$$P(j|i, R)$$

$$P(j|i, R) = \int_{R_j} f_i(\mathbf{x}) d\mathbf{x}$$

- 假如已知样品 X 来自总体 G_i 的先验概率为 q_i , $i = 1, 2, \dots, k$, 则在规则 R 下, 由 (4.12) 式知, 误判的总平均损失为



多元统计

$$\begin{aligned}g(R) &= \sum_{i=1}^k q_i \sum_{j=1}^k C(j|i)P(j|i, R) \\ &= \sum_{i=1}^k q_i \sum_{j=1}^k C(j|i) \int_{R_j} f_i(\mathbf{x}) d\mathbf{x}\end{aligned}$$

$$= \sum_{j=1}^k \int_{R_j} \left(\sum_{i=1}^k q_i C(j|i) f_i(\mathbf{x}) \right) d\mathbf{x}$$

$$\sum_{i=1}^k q_i C(j|i) f_i(\mathbf{x}) = h_j(\mathbf{x})$$

$$g(R) = \sum_{j=1}^k \int_{R_j} h_j(\mathbf{x}) d\mathbf{x}$$



多元统计

■ 如果空间 R^p

$$R^* = (R_1^*, R_2^*, \dots, R_k^*)$$

$$g(R^*) = \sum_{j=1}^k \int_{R_j^*} h_j(\mathbf{x}) d\mathbf{x}$$

$$g(R) - g(R^*) = \sum_{i=1}^k \sum_{j=1}^k \int_{R_i \cap R_j^*} [h_i(\mathbf{x}) - h_j(\mathbf{x})] d\mathbf{x}$$

$$R_i \quad R_i \quad h_i(\mathbf{x}) \leq h_j(\mathbf{x}) \quad j$$

$$R_1, R_2, \dots, R_k$$



多元统计

- 这样，我们以 Bayes 判别的思想得到的划分 $R = (R_1, R_2, \dots, R_k)$

$$R_i = \{\mathbf{x} \mid h_i(\mathbf{x}) = \min_{1 \leq j \leq k} h_j(\mathbf{x})\} \quad i = 1, 2, \dots, k$$

\mathbf{X}

k

$$h_j(\mathbf{x}) = \sum_{i=1}^k q_i C(j \mid i) f_i(\mathbf{x}) \quad j = 1, 2, \dots, k$$

k

$$h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_k(\mathbf{x})$$

\mathbf{X}



以上内容仅为本文档的试下载部分，为可阅读页数的一半内容。如要下载或阅读全文，请访问：<https://d.book118.com/638047140130006132>