

单峰子集类的分离算法由于要对概率密度函数进行估计,所以计算的工作量很大。此外,由于在进行概率估计时要选定一些参数,因此估计的结果也会受到参数的较大影响。特别是在有噪声的情况下,具有局部最大值的概率密度函数的峰点数都会发生变化,从而不能正确反映数据中的单峰子集数。在样本数较少的情况下,由于没有可能对概率密度函数进行估计,从而也就使得这种方法完全失去意义。这时分级聚类算法可能是特别有用的。

与监督模式识别相比,之所以非监督模式识别问题中存在更大的不确定性,一个主要的原因就是在非监督问题中我们没有了已知类别的样本集,甚至可能不知道类别数,可以利用的信息量大大减少了。在实际应用中,除了本章介绍的内容外,还应注意设法有效利用应用领域的专门知识,以弥补信息的不足。最终所得聚类的实际含义也往往只有依靠有关知识来解释和确定。

## 习 题

10.1 令  $x_1, \dots, x_N$  是  $d$  维样本,  $\Sigma$  是任一非奇异  $d \times d$  矩阵。证明使

$$\sum_{k=1}^N (x_k - x)^T \Sigma^{-1} (x_k - x)$$

最小的向量  $x$  是样本均值。

10.2 令  $s(x, x') = x^T x' / (\|x\| \cdot \|x'\|)$ 。若  $x$  的  $d$  个特征只取  $+1$  和  $-1$  二值,即当  $x$  具有第  $i$  个特征时,  $x_i = 1$ , 而当  $x$  没有该特征时  $x_i = -1$ , 说明  $s$  是一个相似性度量。证明对于这种情况

$$\|x - x'\|^2 = 2d(1 - s(x, x'))。$$

10.3 假使一个有  $N$  个样本的集合  $\mathcal{X}$  划分为  $c$  个不相交的子集  $\mathcal{X}_1, \dots, \mathcal{X}_c$ , 假使  $\mathcal{X}_i$  是空集, 则  $\mathcal{X}_i$  中样本的均值  $m_i$  不定义。在这种情况下, 误差平方和只和非空子集有关:

$$J_r = \sum_i \sum_{x \in \mathcal{X}_i} \|x - m_i\|^2$$

这里  $i$  是不包含空子集的子集标号。

假定  $N \geq c$ , 证明使  $J_r$  最小的划分中没有空子集。

10.4 考虑一个  $N = 2k + 1$  样本的集合, 其中有  $k$  个在  $x = -2$  重合, 有  $k$  个在  $x = 0$  上重合, 有一个在  $x = a > 0$  上。证明若  $a^2 < 2(k + 1)$ , 则使  $J_r$  最小的两类划分是  $x = 0$  的  $k$  个样本和  $x = a$  的那个样本聚为一类。若  $a^2 > 2(k + 1)$ , 则应如何聚类使  $J_r$  最小?

10.5 设  $x_1 = (4 \ 5)^T, x_2 = (1 \ 4)^T, x_3 = (0 \ 1)^T, x_4 = (5 \ 0)^T$ 。现有下列三种划分:

$$(1) \mathcal{X}_1 = \{x_1, x_2\}, \quad \mathcal{X}_2 = \{x_3, x_4\}$$

$$(2) \mathcal{X}_1 = \{x_1, x_4\}, \quad \mathcal{X}_2 = \{x_2, x_3\}$$

$$(3) \mathcal{X}_1 = \{x_1, x_2, x_3\}, \quad \mathcal{X}_2 = \{x_4\}$$

证明对于平方误差和准则, 第三种划分最好, 而若用  $|S_w|$  准则, 前两种划分为好。

10.6 若定义下列准则函数

$$J_T = \sum_{i=1}^c \sum_{x \in \mathcal{X}_i} (x - m_i)^T S_T^{-1} (x - m_i)$$

其中  $\mathbf{m}_i$  是  $\mathcal{X}_i$  中  $N_i$  个样本的均值向量,  $S_T$  是总散布矩阵,

(1) 证明  $J_T$  对数据的非奇异线性变换具有不变性。

(2) 证明把  $\mathcal{X}_i$  中的样本  $\hat{\mathbf{x}}$  转移到  $\mathcal{X}_j$  中去, 则使  $J_T$  改变为

$$J_T^* = J_T - \left[ \frac{N_j}{N_j + 1} (\hat{\mathbf{x}} - \mathbf{m}_j)^T S_T^{-1} (\hat{\mathbf{x}} - \mathbf{m}_j) - \frac{N_i}{N_i + 1} (\hat{\mathbf{x}} - \mathbf{m}_i)^T S_T^{-1} (\hat{\mathbf{x}} - \mathbf{m}_i) \right]$$

(3) 写出使  $J_T$  最小化的迭代程序。

**10.7** 令聚类  $\mathcal{X}_i$  包含  $N_i$  个样本, 令  $d_{ij}$  是两个聚类  $\mathcal{X}_i$  和  $\mathcal{X}_j$  之间的距离度量。若  $\mathcal{X}_i$  和  $\mathcal{X}_j$  合并形成一个新的聚类  $\mathcal{X}_k$ , 则  $\mathcal{X}_k$  到  $\mathcal{X}_h$  的距离一般情况下可用下式表示:

$$d_{hk} = \alpha_i d_{hi} + \alpha_j d_{hj} + \beta d_{ij} + \gamma |d_{hi} - d_{hj}|$$

证明当用不同的距离度量时,  $\alpha_i, \alpha_j, \beta, \gamma$  分别取不同的数值:

(1)  $d_{\min}: \alpha_i = \alpha_j = 0.5, \beta = 0, \gamma = -0.5$

(2)  $d_{\max}: \alpha_i = \alpha_j = 0.5, \beta = 0, \gamma = 0.5$

(3)  $d_{\text{avg}}: \alpha_i = \frac{N_i}{N_i + N_j}, \alpha_j = \frac{N_j}{N_i + N_j}, \beta = \gamma = 0$

(4)  $d_{\text{mean}}^2: \alpha_i = \frac{N_i}{N_i + N_j}, \alpha_j = \frac{N_j}{N_i + N_j}, \beta = -\alpha_i \alpha_j, \gamma = 0$ 。

**10.8** 证明任一对称集上的概率密度函数是单峰的。

**10.9** 证明对于  $C$ -均值算法, 聚类准则函数满足使算法收敛的条件。(即若  $J(\Gamma, \tilde{K}) \leq J(\Gamma, K)$ , 则有  $J(\tilde{\Gamma}, \tilde{K}) \leq J(\Gamma, \tilde{K})$ )

**10.10** 令  $\Delta(\mathbf{y}, K_i) = \frac{1}{2} (\mathbf{y} - \mathbf{m}_i)^T \Sigma_i^{-1} (\mathbf{y} - \mathbf{m}_i) + \frac{1}{2} \log |\Sigma_i|$  是点到聚类的相似性度量, 式中  $\mathbf{m}_i$  和  $\Sigma_i$  是聚类  $\Gamma_i$  的均值和协方差矩阵, 若把一点从  $\Gamma_i$  转移到  $\Gamma_j$  中去, 计算由公式 (10-44) 所示  $J_K$  的变化值。