

第十章 非监督学习与聚类

苏智勇

可视计算研究组

南京理工大学

suzhiyong@njust.edu.cn

<https://zhiyongsu.github.io>

主要内容

10.1 引言

10.2 基于模型的聚类方法

10.3 混合模型的估计

10.4 动态聚类算法

10.5 模糊聚类方法

10.6 分级聚类方法

10.7 一致聚类方法

10.1 引言

- 监督模式识别
 - (已知) 样本集 → 训练 (学习, 分类器设计) → 识别 (分类)
- 非监督模式识别
 - (未知) 样本集 → 非监督学习 (聚类分析) → 后处理
 - 分类
 - ✓ 基于模型的方法
 - ✓ 基于相似性度量的方法

10.2 基于模型的聚类方法：单峰子集分离法

- 基本假设

- 每个聚类的样本分布是单峰的，根据总体分布中的单峰来划分子集

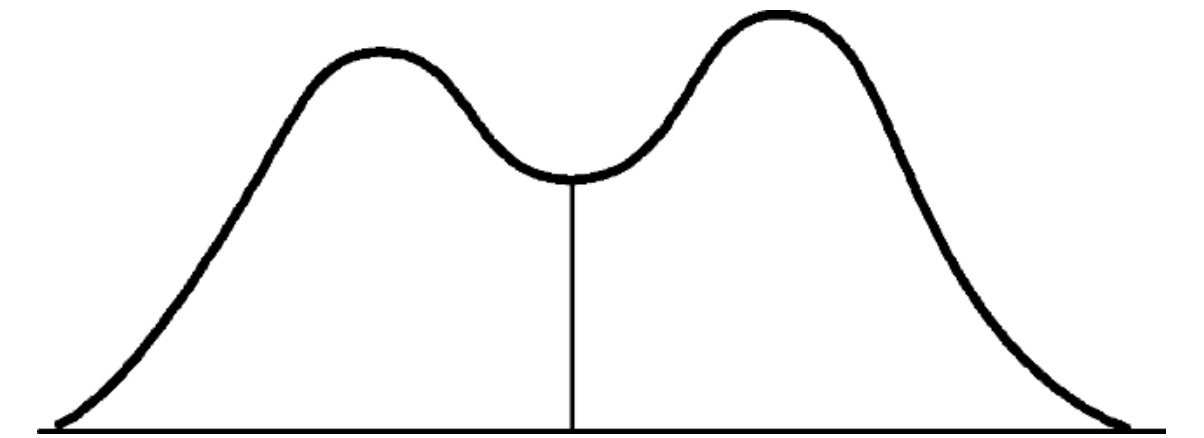
- 投影方法

- 基本思路

把样本按照某种准则投影到某个一维坐标上，在这一维度上估计样本的概率密度，在其中寻找单峰并进行聚类划分(如果这一维上只有一个峰，则寻找下一个投影方向)

- 投影方向

使方差最大的方向，即协方差矩阵本征值最大的本征向量方向



10.2 基于模型的聚类方法：单峰子集分离法

- 算法步骤

(1)主成分分析：计算所有样本 $\{\mathbf{x}\}$ 的协方差矩阵的最大本征值对应的本征向量 \mathbf{u}_j ，把样本投影到 \mathbf{u}_j 上 $v_j = \mathbf{u}_j^T \mathbf{x}$

(2)用非参数法估计投影后样本 $v_j = \mathbf{u}_j^T \mathbf{x}$ 的概率密度函数 $P(v_j)$ （用直方图方法或其它方法）

(3)求 $P(v_j)$ 中的极小点（波谷），在这些极小点上作垂直于 \mathbf{u}_j 的超平面作为分类超平面，得到子集划分

(4)如果 $P(v_j)$ 上没有这种极小点，则用下一个本征值对应的本征向量作为投影方向，重复(2) ~ (3)

(5)对划分出的每一个子集重复上述过程，直到不能再分（所有方向上都是单峰）

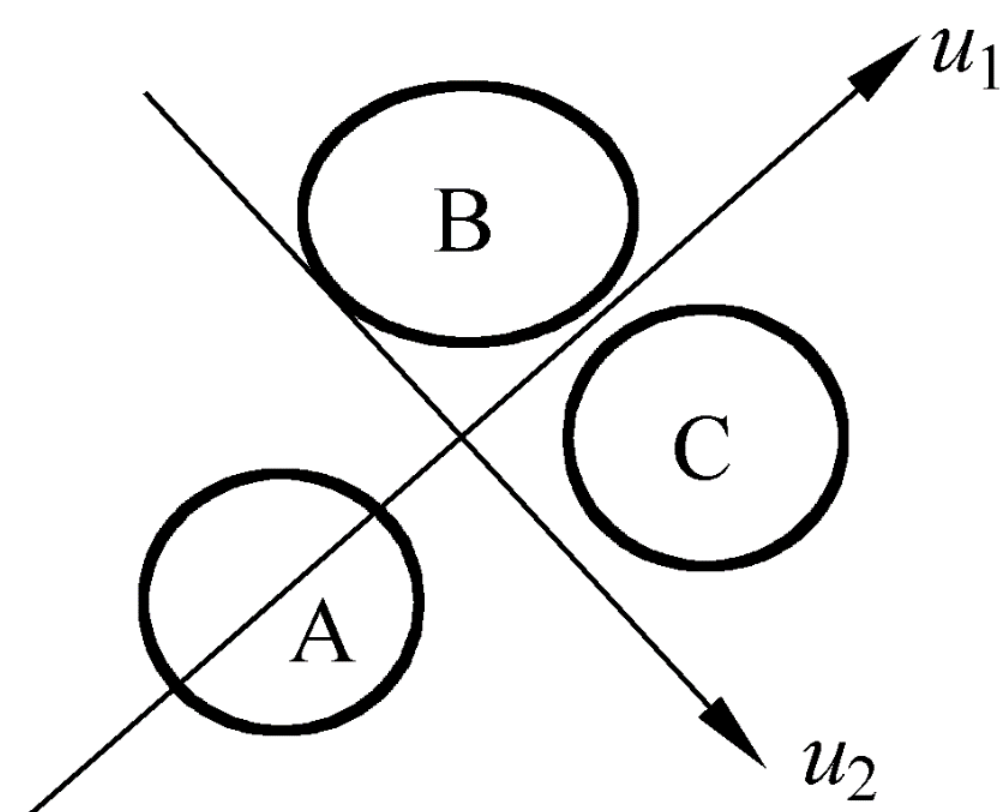
10.2 基于模型的聚类方法：单峰子集分离法

- 问题

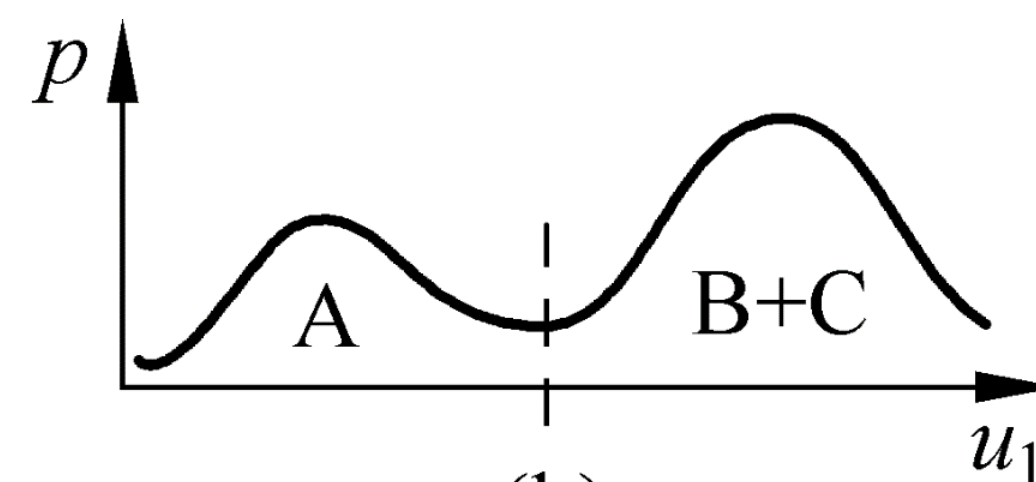
- 如何选择投影方向？

- 方差最大的准则（经验准则）有时并不一定最有利于聚类

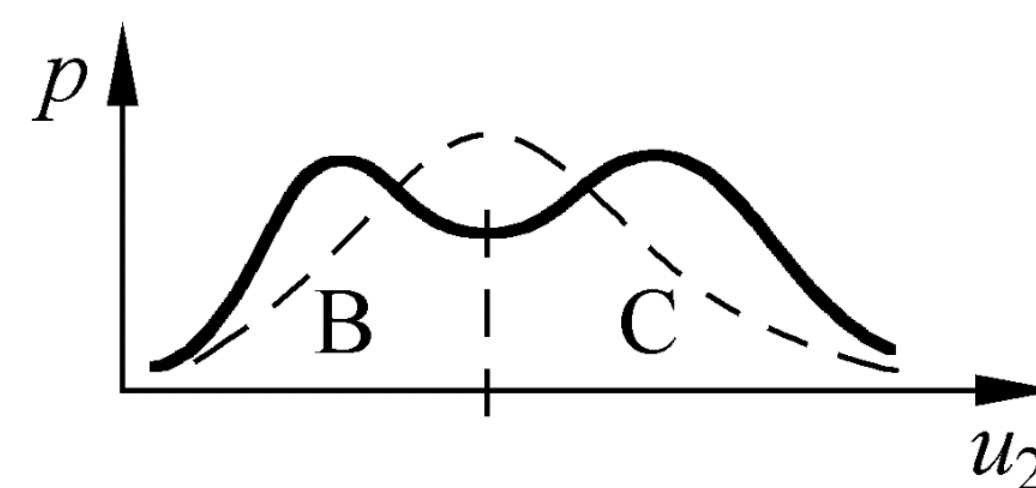
- 有效的聚类



(a)



(b)

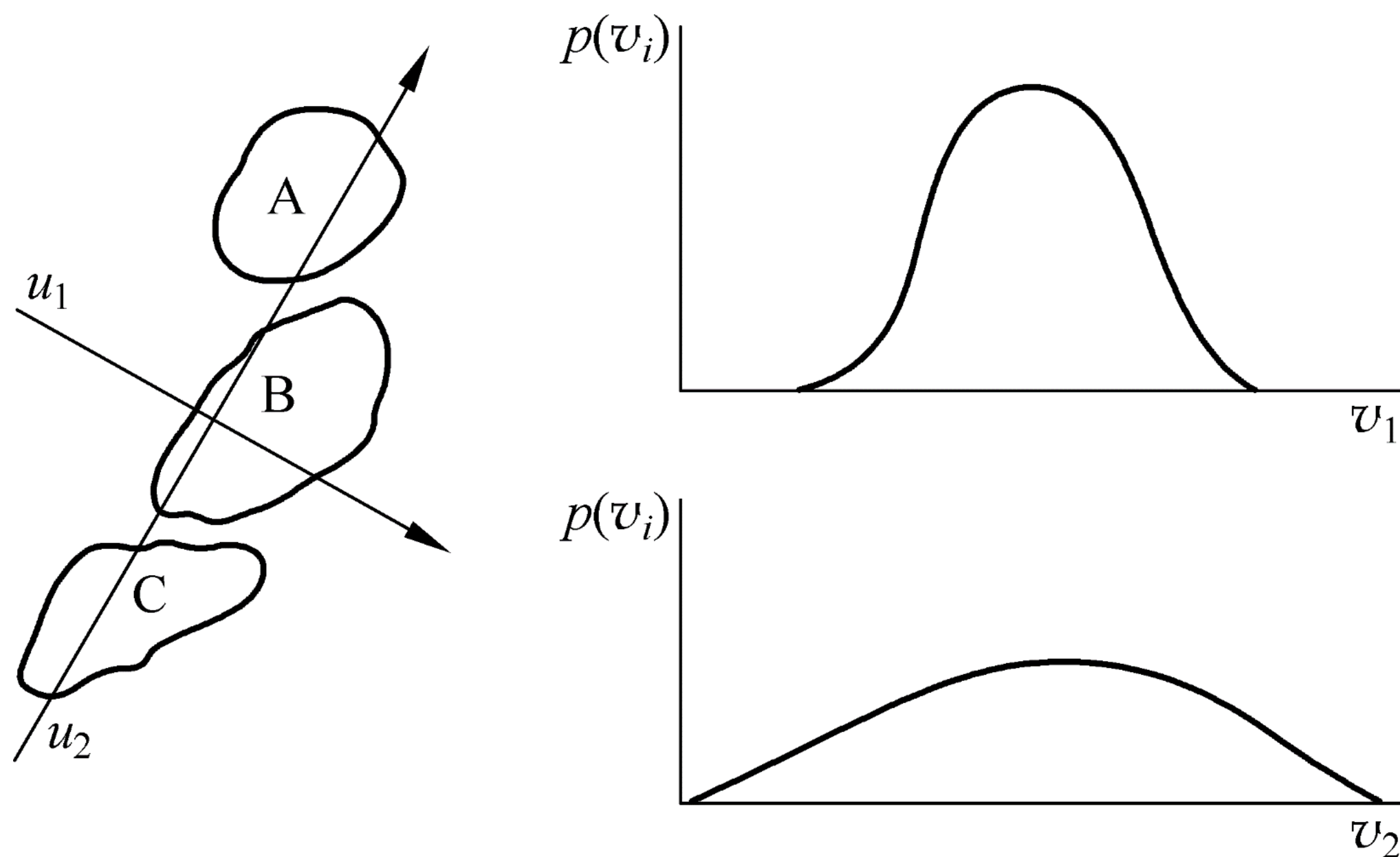


(c)

10.2 基于模型的聚类方法：单峰子集分离法

- 失败的聚类

- 例：在两个主成分方向上都得不到单峰子集



10.3 混合模型的估计：非监督参数估计

- 定义

- 非监督参数估计指样本类别未知，但各类条件概率密度函数的形式已知，根据所有样本估计各类密度函数中的参数。

- 本节只介绍非监督最大似然估计的思路

10.3.1 混合密度的最大似然估计

- 假设条件

- 样本集 $X = \{x_1, \dots, x_N\}$ 中的样本属于 c 个类别，但不知各样本属哪类。
- 类先验概率 $P(\omega_i), i = 1, \dots, c$ 已知
- 类条件概率密度形式已知 $P(x | \omega_i, \theta_i), i = 1, \dots, c$
- 未知的仅是 c 个参数向量 $\theta_1, \theta_2, \dots, \theta_c$ 的值，所有未知参数组成的向量记为 $\theta = [\theta_1, \theta_2, \dots, \theta_c]^T$

10.3.1 混合密度的最大似然估计

- 似然函数

- 混合密度: $p(x|\theta) = \sum_{i=1}^c p(x|\omega_i, \theta_i) P(\omega_i)$

- 分量密度: 类条件密度 $p(x|\omega_i, \theta_i)$

- 混合参数: 先验概率 $P(\omega_i)$ (有时也可未知, 一起参与估计)

- 设样本集 X 中的样本是从混合密度为 $p(x|\theta)$ 的总体中独立抽取的, 即满足独立同分布条件, θ 确定但未知, 则

- ✓ 似然函数: $l(\theta) = p(X|\theta) = \prod_{i=1}^N p(x_i|\theta)$

- ✓ 对数似然函数: $H(\theta) = \ln [l(\theta)] = \sum_{i=1}^N \ln p(x_i|\theta)$

- ✓ 最大似然估计 $\hat{\theta}$ 就是使 $l(\theta)$ 或 $H(\theta)$ 取最大的 θ 值

10.3.1 混合密度的最大似然估计

- 可识别性问题

- 求出 $\hat{\theta}$ ，就得到了 $\hat{\theta}_1, \dots, \hat{\theta}_c$ ，即从混合密度函数中恢复出了分量密度函数。可能吗？什么条件下可能？
- 可识别性定义：若对 $\theta \neq \theta'$ ，混合分布中总存在 x 使 $p(x | \theta) \neq p(x | \theta')$ ，则密度 $p(x | \theta)$ 是可识别的
 - 大部分常见连续随机变量的分布密度函数都是可识别的
 - 离散随机变量的混合概率函数则往往是不可识别的

10.3.1 混合密度的最大似然估计

- 计算问题

- 对于可识别的似然函数，如何求最大似然估计？
- 思路同监督情况，即如果 $p(x|\theta)$ 对 θ 可微，则令 $\nabla_{\theta}H(\theta) = 0$
- 得一系列方程组，它们是最大似然估计的必要条件，若存在唯一极值则就是解

$$\begin{aligned}\nabla_{\theta_i}H(\theta) &= \sum_{k=1}^N \frac{1}{p(x_k|\theta)} \nabla_{\theta_i} \left[\sum_{j=1}^c p(x_k|\omega_j, \theta_j) P(\omega_j) \right] \\ &= \sum_{k=1}^N \frac{1}{p(x_k|\theta)} \nabla_{\theta_i} \left[p(x_k|\omega_i, \theta_i) P(\omega_i) \right] = \sum_{k=1}^N \frac{1}{p(x_k|\theta)} \nabla_{\theta_i} \ln p(x_k|\omega_i, \theta_i) \text{ (设 } \theta_i, \theta_j \text{ 独立)}\end{aligned}$$

10.3.1 混合密度的最大似然估计

- 计算问题

- 其中后验概率
$$P(\omega_i | x_k, \theta_i) = \frac{p(x_k | \omega_i, \theta_i) P(\omega_i)}{p(x_k | \theta)}$$

- 有微分方程组:
$$\nabla_{\theta_i} H(\hat{\theta}) = 0, i = 1, 2, \dots, c$$

- 另, 若 $p(\omega_i)$ 也未知, 则可引入限制条件: $P(\omega_i) > 0, i = 1, 2, \dots, c, \sum_{i=1}^c P(\omega_i) = 1$

- 可用拉格朗日法求条件极值问题, 定义拉格朗日函数:
$$H' = H + \lambda \left[\sum_{i=1}^c P(\omega_i) - 1 \right]$$

10.3.1 混合密度的最大似然估计

- 计算问题

- 可得

$$\hat{P}(\omega_i) = \frac{1}{N} \sum_{k=1}^N \hat{P}(\omega_i | x_k, \hat{\theta}_i), i = 1, 2, \dots, c$$

$$\sum_{k=1}^N P(\omega_i | x_k, \omega_i) \nabla_{\theta_i} \ln p(x_k | \omega_i, \theta_i) = 0, i = 1, 2, \dots, c$$

$$\text{其中, } \hat{P}(\omega_i | x_k, \hat{\theta}_i) = \frac{p(x_k | \omega_i, \hat{\theta}_i) \hat{P}(\omega_i)}{\sum_{j=1}^c p(x_k | \omega_j, \hat{\theta}_j) \hat{P}(\omega_j)}, i = 1, 2, \dots, c$$

原则上可以从上述微分方程组中求解出最大似然估计 $\hat{\theta}$ 和 $\hat{P}(\omega_i)$ ，但实际上多数问题中只能采用某种迭代方法求解

10.3.2 混合正态分布的参数估计

- 混合高斯模型 (mixture of Gaussian models)

- 混合模型中的各个分布都是多维正态分布，即

$$p(x | \omega_i, \theta_i) \sim N(\mu_i, \Sigma_i)$$

- 主要有三种情况： (“?”表示未知, “✓”表示已知)

情况	μ_i	Σ_i	$P(\omega_i)$	c
1	?	✓	✓	✓
2	?	?	?	✓
3	?	?	?	?

10.3.2 混合正态分布的参数估计

- 情况1: 均值向量 μ_i 未知, 其他参数已知

- 由上节知最大似然估计满足方程组

$$\sum_{k=1}^N \hat{P}(\omega_i | x_k, \hat{\theta}_i) \nabla_{\theta_i} \ln p(x_k | \omega_i, \hat{\theta}_i) = 0, i = 1, \dots, c$$

- 代入正态分布公式, 可得

$$\sum_{k=1}^N P(\omega_i | x_k, \hat{\mu}_i) \Sigma_i^{-1} (x_k - \hat{\mu}_i) = 0$$

10.3.2 混合正态分布的参数估计

- 情况1: 均值向量 μ_i 未知, 其他参数已知

- 即

$$\hat{\mu}_i = \frac{\sum_{k=1}^N P(\omega_i | x_k, \hat{\mu}_i) x_k}{\sum_{k=1}^N P(\omega_i | x_k, \hat{\mu}_i)}$$

- 样本的加权平均, 物理意义明确; 但是权值中包含未知参数, 其中

$$\hat{P}(\omega_i | x_k, \hat{\mu}_i) = \frac{p(x_k | \omega_i, \hat{\mu}_i) P(\omega_i)}{\sum_{j=1}^c p(x_k | \omega_j, \hat{\mu}_j) P(\omega_j)}$$

10.3.2 混合正态分布的参数估计

- 情况1: 均值向量 μ_i 未知, 其他参数已知
 - 迭代法求解: 用某种方法得到一个较好的初值 $\hat{\mu}_i(0)$, 然后用下式迭代:

$$\hat{\mu}_i(j+1) = \frac{\sum_{k=1}^N P(\omega_i | x_k, \hat{\mu}_i(j)) x_k}{\sum_{k=1}^N P(\omega_i | x_k, \hat{\mu}_i(j))}$$

- 梯度法, 可能不是全局最优解, 受初值影响大

10.3.2 混合正态分布的参数估计

- 情况2：类别数目 c 已知，其他参数均未知
 - 思路与情况 1类似，将有关分布公式代入上小节方程即可，只是公式复杂一些，也可得到物理意义明确的方程式，但一般也只能用迭代法求解
- 情况3：参数均未知
 - 无法用最大似然法求解

10.4 动态聚类算法

- **基于模型的方法**
 - 估计概率密度函数：困难
 - 寻找密度函数中的单峰：需要较多样本或先验知识，在非监督学习中不易满足
- **基于相似性度量的聚类方法**
 - 考查样本之间的相似性，根据相似性把样本集划分为若干子集，使某种表示聚类质量的准则函数最优
 - 相似性度量：以某种距离定义
 - 直观理解：同一类的样本的特征向量应是相互靠近的（前提：特征选取合理，能反映所感兴趣的关系）

10.4 动态聚类算法

- 动态聚类方法

- 多次迭代，逐步调整类别划分，最终使某准则达到最优
- 三个要点：
 - ✓ 选某种**距离**作为样本相似性度量
 - ✓ 定义某个**准则函数**，用于评价聚类质量
 - ✓ 初始分类方法及迭代算法

10.4.1 C均值算法 (K均值算法)

- 误差平方和准则

$$J_e = \sum_{i=1}^c \sum_{y \in \Gamma_i} \|y - m_i\|^2 = \sum_{i=1}^c J_i$$

其中, Γ_i 是第*i*个聚类, $i = 1, 2, \dots, c$, 其中样本数为 N_i , Γ_i 中样本均值为 $m_i = \frac{1}{N_i} \sum_{y \in \Gamma_i} y$

- J_e 反映了用*c*个聚类中心代表*c*个样本子集所带来的总的误差平方和
- J_e 是样本集 Y 与类别集 Ω 的函数
- C均值算法的目标: 最小化 —— 最小方差划分
- 用*c*个码本来代表整个样本集, 使这种表示带来的总体误差最小 —— 向量量化
- 无法用解析法求解, 只能用迭代法, 通过不断调整样本的类别归属来求解

10.4.1 C均值算法 (K均值算法)

- 算法研究

设已有一个划分方案，考查 Γ_k 中的样本 y ，若把 y 移入 Γ_j ，此时 Γ_k 变成 $\tilde{\Gamma}_k$ ，此时 Γ_j 变成 $\tilde{\Gamma}_j$ ，有

$$\tilde{m}_k = m_k + \frac{1}{N_k - 1} [m_k - y], \quad \tilde{m}_j = m_j + \frac{1}{N_j + 1} [y - m_j]$$

$$\tilde{J}_k = J_k - \frac{N_k}{N_k - 1} \|y - m_k\|^2, \quad \tilde{J}_j = J_j + \frac{N_j}{N_j + 1} \|y - m_j\|^2$$

若 $\frac{N_j}{N_j + 1} \|y - m_j\|^2 < \frac{N_k}{N_k - 1} \|y - m_k\|^2$ ，则把 y 从 Γ_k 移入 Γ_j 会使 J_e 减小

10.4.1 C均值算法 (K均值算法)

- C均值算法

- (1) 初始划分 c 个聚类 Γ_i , 计算 m_i 和 J_e , $i = 1, 2, \dots, c$
- (2) 任取一个样本 y , 设 $y \in \Gamma_i$
- (3) 若 $N_i = 1$, 则转 (2); 否则继续
- (4) 计算 $\rho_i = \frac{N_i}{N_i - 1} \|y - m_i\|^2$, $\rho_j = \frac{N_j}{N_j + 1} \|y - m_j\|^2, i \neq j$
- (5) 选 ρ_j 中的最小者, 即若 $\rho_k < \rho_j, \forall j$, 则把 y 从 Γ_i 移到 Γ_k 中: 有利于 J_e 的减少
- (6) 重新计算 m_i 和 J_e , $i = 1, 2, \dots, c$
- (7) 若连续 N 次迭代 J_e 不改变, 则停止; 否则转 (2)

10.4.1 C 均值算法 (K 均值算法)

- 初始划分：一般先选代表点，再进行初始分类
 - 代表点选择方法
 - ✓ 经验选择
 - ✓ 随机分成 c 类，选各类重心作为代表点
 - ✓ “密度”法：计算每个样本的一定球形邻域内的样本数作为“密度”，选“密度”最大的样本点作为第一个代表点，在离它一定距离之外选最大“密度”点作为第二个代表点，...，依次类推。
 - ✓ 用前 c 个样本点作为代表点
 - ✓ 用 $c-1$ 聚类求 c 个代表点：各类中心外加离它们最远的样本点，从1类开始...

10.4.1 C均值算法 (K均值算法)

- 初始分类方法

- 最近距离法：离哪个代表点近就归入哪一类
- 最近距离法归类，但每次都重新计算该类代表点
- 直接划分初始分类：第一个样本自成一类，第二个样本若离它小于某距离阈值则归入此类，否则建新类，……
- 将特征归一化，用样本各特征之和作为初始分类依据

$$SUM(i) = \sum_{j=1}^D y_{ij}, i = 1, \dots, N. \quad y_i \in R^D$$

$$MA = \max_i SUM(i), \quad MI = \min_i SUM(i), \quad k_{\text{取整}} = \frac{(c-1)}{MA - MI} [SUM(i) - MI] + 1$$

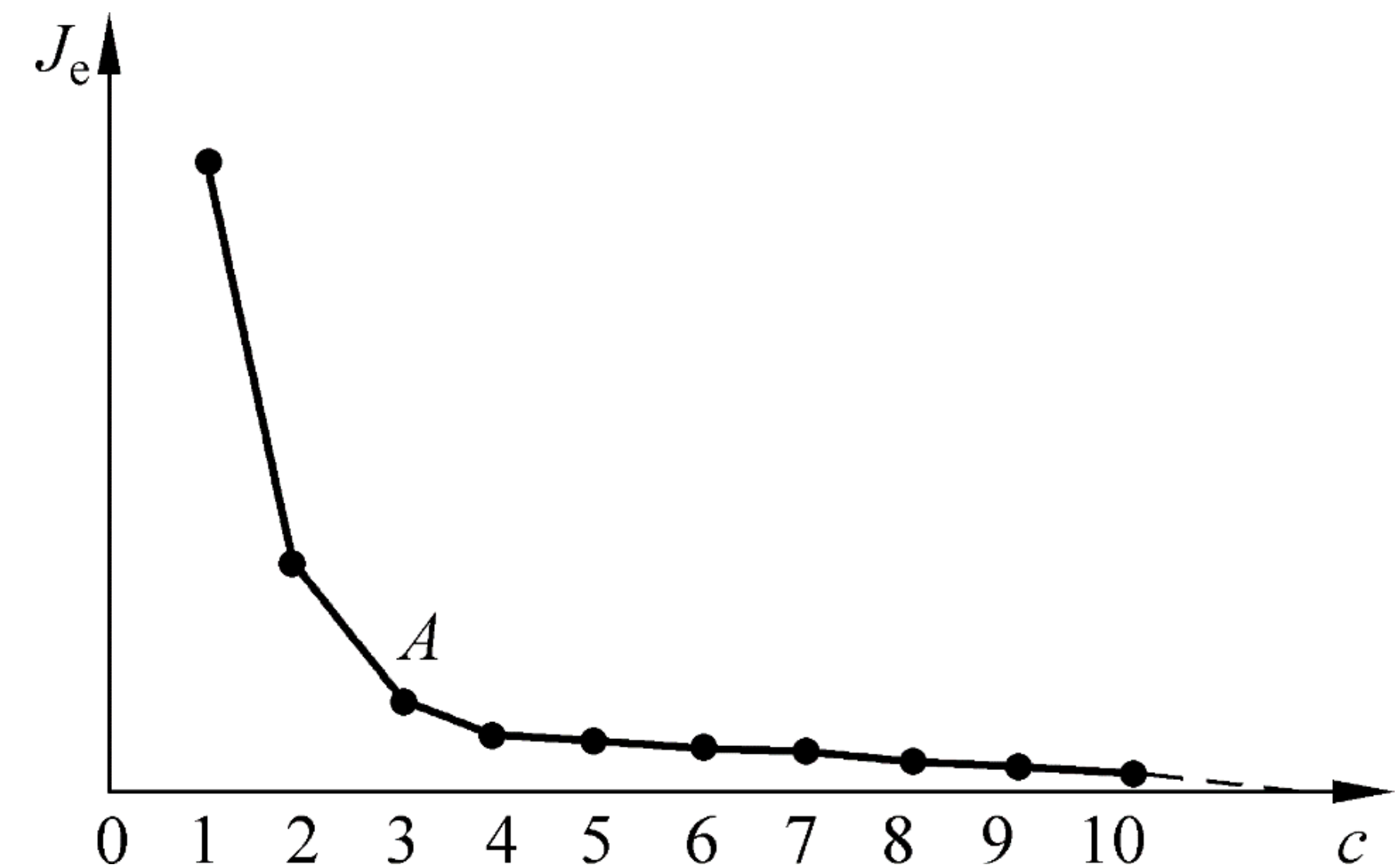
10.4.1 C 均值算法 (K 均值算法)

- 说明

- 初始划分无一定之规，多为启发式方法
- C 均值方法结果受初值影响，是局部最优解

关于 C 均值方法的类别数 c ，多假定已知，由先验知识确定

- 一种实验确定方法
 - 对 $c = 1, 2, 3, \dots$ 聚类，求各自的 $J_e(c)$ 。
 - 找到其中的拐点 (图中 $\hat{c} = 3$) (注：实际问题中并不一定有明确的拐点)



10.4.1 C 均值算法 (K 均值算法)

- C 均值聚类方法用于非监督模式识别的问题：
 - 要求类别数已知
 - 是最小方差划分，并不一定能反映内在分布
 - 与初始划分有关，不保证全局最优

10.4.2 ISODATA方法

- **ISODATA**

- iterative self-organizing data analysis techniques: 迭代自组织数据分析技术
- 一种改进的 C 均值算法
- 特点: 和 C 均值算法的区别
 - 成批样本修正, 把所有样本调整完后才重新计算均值
 - 引入对类别的评判准则, 可进行类别合并与分裂, 突破事先给定类别数的限制

10.4.2 ISODATA方法

- 算法步骤

设有 N 个样本组成的样本集 $\{\mathbf{y}_1, \mathbf{y}_1, \dots, \mathbf{y}_N\} \subset R^d$, 事先设定如下参数:

- K : 期望得到的聚类数
- θ_N : 一个聚类中的最少样本数
- θ_s : 标准偏差参数
- θ_c : 合并参数
- L : 每次迭代允许合并的最大聚类对数
- I : 允许迭代的次数

10.4.2 ISODATA方法

• 算法步骤

- ① 初始化，设初始聚类数 c ，中心 \mathbf{m}_i ， $i = 1, 2, \dots, c$ （不一定等于期望聚类数 K ）
- ② 把所有样本分到距离最近的类中， Γ_i ， $i = 1, 2, \dots, c$
- ③ 若某个类 Γ_j 中样本数过少（ $N_j < \theta_N$ ），则去掉这一类（根据各样本到其他类中心的距离分别合入其他类），置 $c = c - 1$

- ④ 重新计算均值：
$$\mathbf{m}_j = \frac{1}{N_j} \sum_{y \in \Gamma_j} \mathbf{y}, \quad j = 1, \dots, c$$

- ⑤ 计算第 j 类样本与其中心的平均距离和总平均距离

$$\bar{\delta}_j = \frac{1}{N_j} \sum_{y \in T_j} \|y - \mathbf{m}_j\|, \quad \bar{\delta} = \frac{1}{N} \sum_{j=1}^c N_j \bar{\delta}_j, \quad j = 1, \dots, c$$

10.4.2 ISODATA方法

• 算法步骤

⑥ 若是最后一次迭代（由总迭代次数 I 确定），则程序停止；

- 若 $c \leq K/2$ ，则转⑦（分裂）；

- 若 $c \geq 2K$ ，或是偶数次迭代，则转⑧（合并）

⑦ 分裂

对每个类，求各维标准偏差 $\sigma_j = [\sigma_{j1}, \sigma_{j2}, \dots, \sigma_{jd}]^T$,

$$\sigma_{ji} = \sqrt{\frac{1}{N_j} \sum_{y_k \in T_j} (y_{ki} - m_{ji})^2}, \quad j = 1, \dots, c, \quad i = 1, \dots, d$$

10.4.2 ISODATA方法

- 算法步骤

- ⑦ 分裂

- 对每个类，求出标准偏差最大的分量 σ_{jmax} ， $j = 1, \dots, c$

- 在 σ_{jmax} 对应的特征分量上分裂：对各类的 σ_{jmax} ，存在 $\sigma_{jmax} > \theta_s$ （标准偏差参数），且 $\bar{\delta}_j > \bar{\delta}$ ， $N_j > 2(\theta_N + 1)$ ，或 $c \leq k/2$ ，则 Γ_j 分裂为两类，中心分别为 m_j^+ 和 m_j^- ，置 $c = c + 1$ ， $m_j^+ = m_j + \gamma_j$ ， $m_j^- = m_j - \gamma_j$ ，其中 $\gamma_j = k\sigma_{jmax}$ ， $0 < k \leq 1$

10.4.2 ISODATA方法

- 算法步骤

- ⑧ 合并

- 计算各类中心之间的距离算 $\delta_{ij} = \| m_i - m_j \|$, $i, j = 1, \dots, c, i \neq j$

- 比较 δ_{ij} 与 θ_c (合并参数), 对小于 θ_c 者排序 $\delta_{i_1j_1} < \delta_{i_2j_2} < \dots < \delta_{i_lj_l}$

- 从最小的 $\delta_{i_lj_l}$ 开始, 把每个 $\delta_{i_lj_l}$ 对应的 m_{i_l} 和 m_{j_l} 合并:

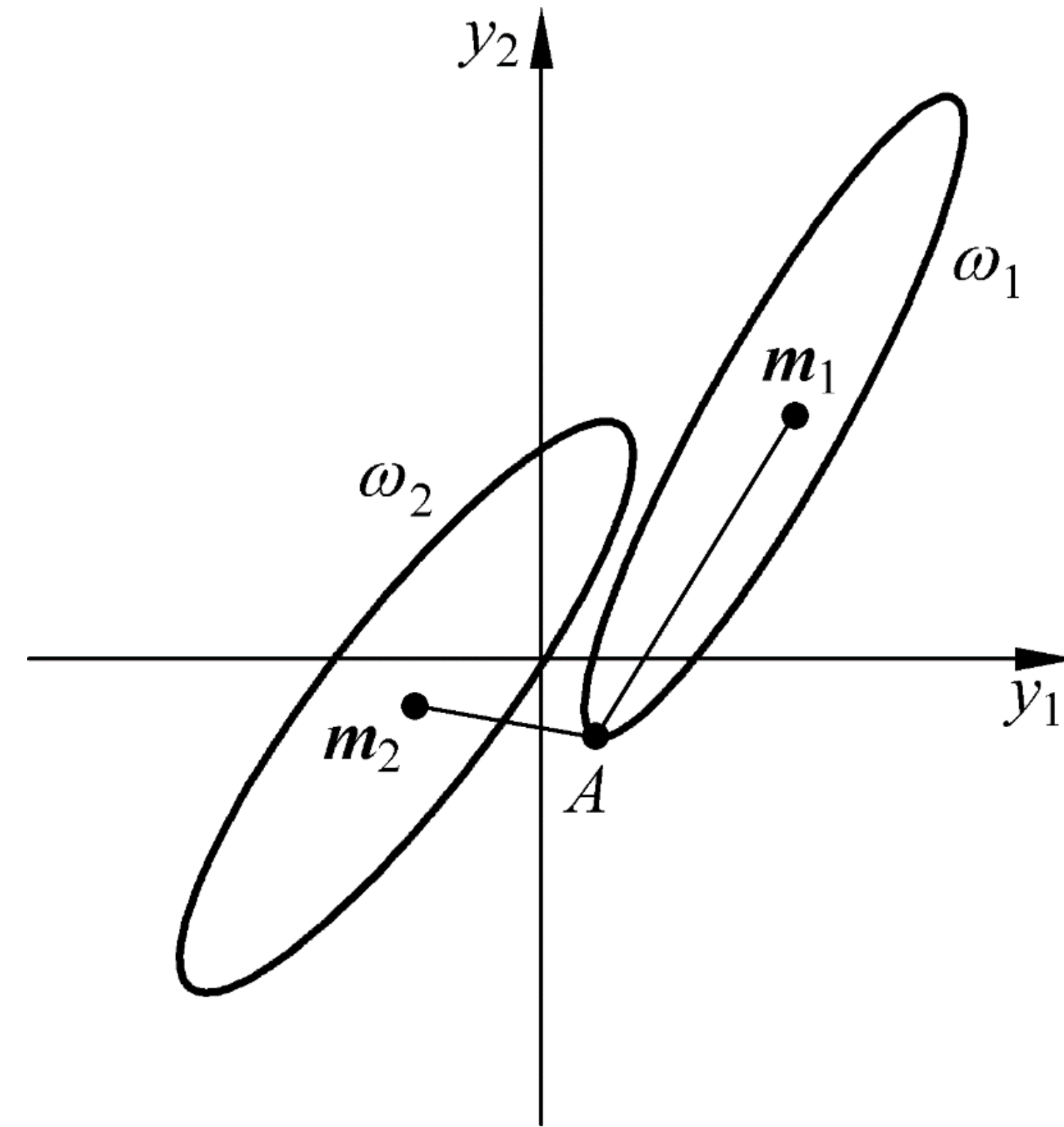
- $$m_l = \frac{1}{N_{i_l} + N_{j_l}} [N_{i_l} m_{i_l} + N_{j_l} m_{j_l}]$$
, 并置 $c = c - 1$ 。每次迭代中避免同一类被合并

- 两次。

- ⑨ 若是最后一次迭代, 则终止。否则迭代次数加1, 转② (必要时可调整算法参数)

10.4.3 基于核的动态聚类算法

- C 均值方法的缺点
 - 用均值代表类，只适用于近似球状分布的类
- 改进
 - 用核 $K_j = K(\mathbf{y}, V_j)$ 来代表一个类 Γ_i ， V_j 是参数集
 - 核 K 可以是一个函数、一个点集或某种分类模型
 - 定义样本 \mathbf{y} 到类 Γ_i (核 K_j) 之间的距离: $\Delta(\mathbf{y}, K_j)$



10.4.3 基于核的动态聚类算法

- 准则函数

$$J_k = \sum_{i=1}^c \sum_{y \in \Gamma_j} \Delta(y, K_j)$$

- 算法流程

- ① 选择初始划分，将样本集划分为 c 类，得到初始核 K_j , $j = 1, 2, \dots, c$
- ② 按以下规则把各样本分类：若 $\Delta(y, K_j) = \min_{k=1, \dots, c} \Delta(y, K_k)$ ，则 $y \in \Gamma_j$
- ③ 更新 K_j , $j = 1, 2, \dots, c$ ，若 K_j 不变，则终止；否则转②
- ④ C 均值可看作基于核的动态聚类算法的特例： K_j 为 m_j ， Δ 为欧氏距离

10.4.3 基于核的动态聚类算法

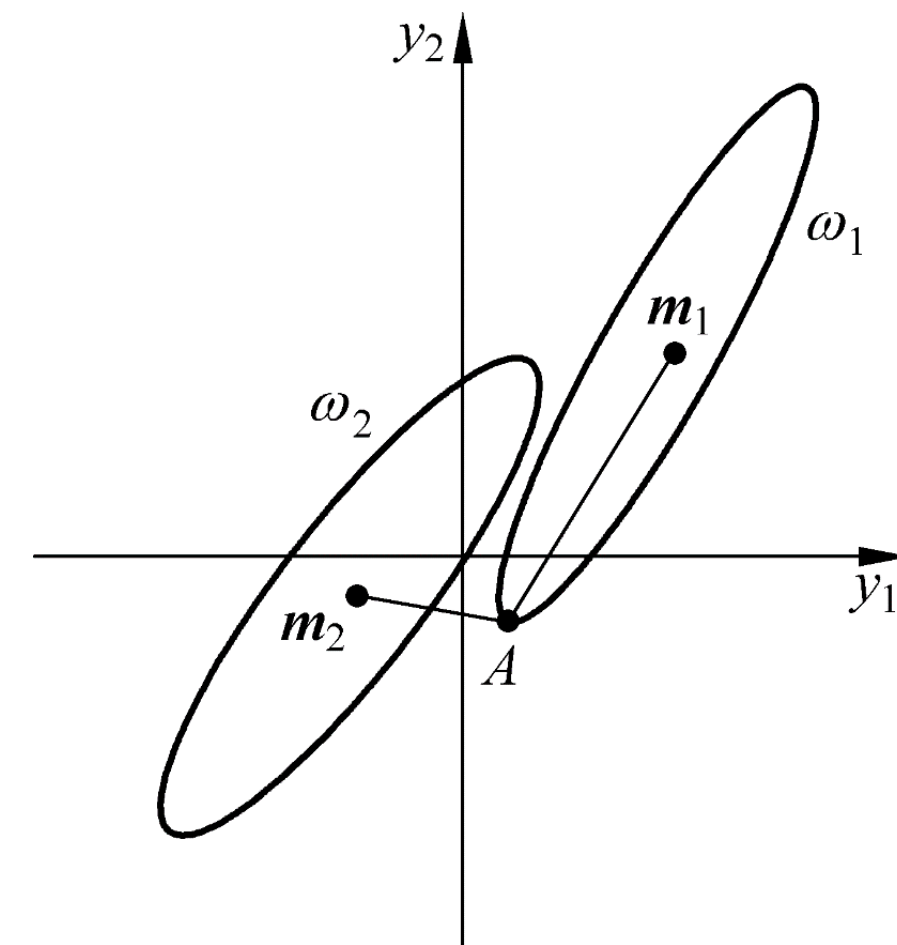
- 两种核函数举例

- 正态核函数

$$K_k(\mathbf{y}, V_j) = \frac{1}{(2\pi)^{d/2} \left| \Sigma_j \right|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{m}_j)^T \hat{\Sigma}_j^{-1} (\mathbf{y} - \mathbf{m}_j) \right\}$$

$$V_j = \left\{ \mathbf{m}_j, \hat{\Sigma}_j \right\}$$

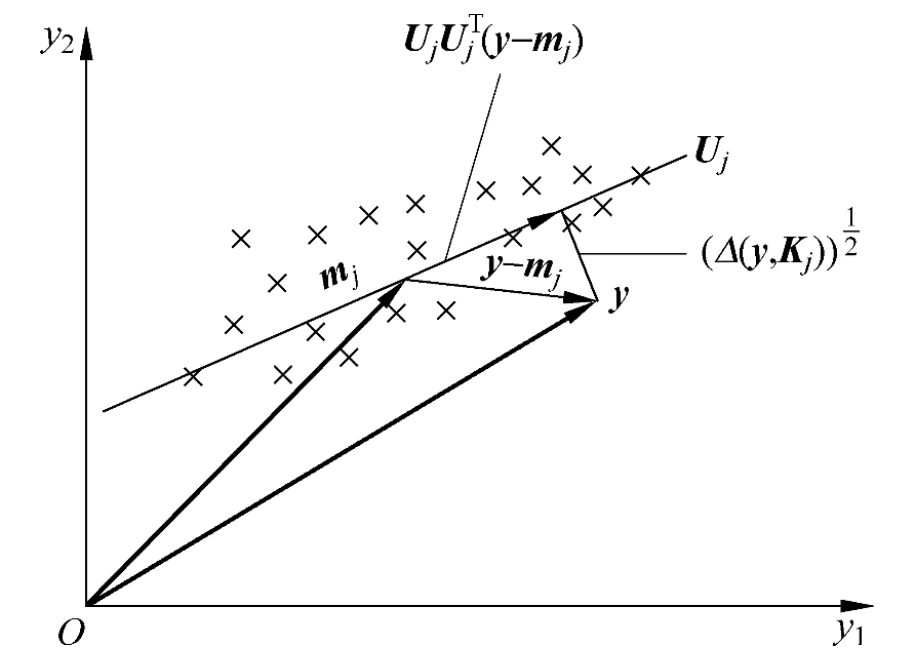
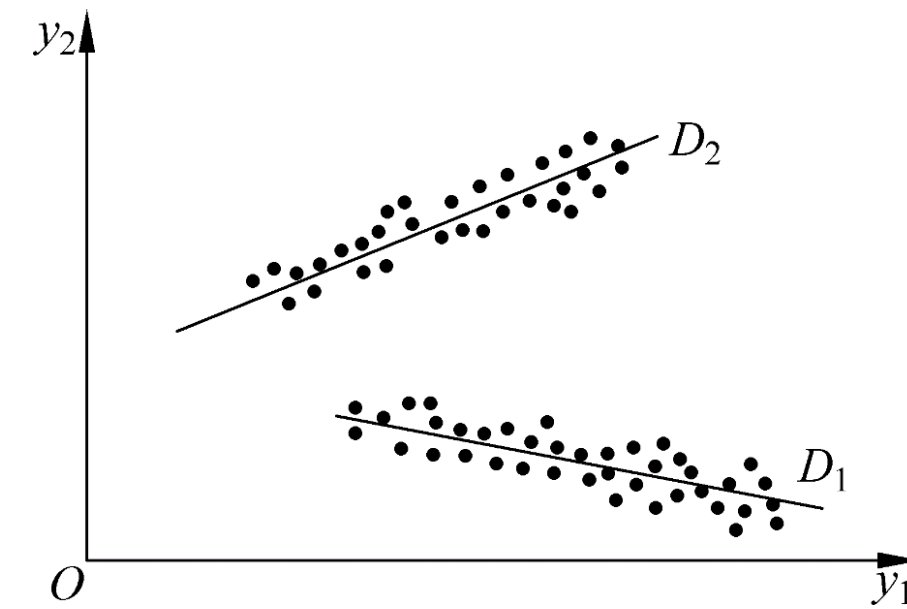
$$\Delta(\mathbf{y}, K_j) = \frac{1}{2} (\mathbf{y} - \mathbf{m}_j)^T \hat{\Sigma}_j^{-1} (\mathbf{y} - \mathbf{m}_j) + \frac{1}{2} \log \left| \hat{\Sigma}_j \right|$$



10.4.3 基于核的动态聚类算法

- 两种核函数举例

- 主轴核函数



- 用K-L变换得到样本子集的主轴方向作为核: $K(y, V_j) = U_j^T y$,

- $U_j^T = [u_1, u_2, \dots, u_{d_j}]$ 是 $\hat{\Sigma}_j$ 的 d_j 个最大本征值的本征向量系统

- 计算样本 y 到 Γ_j 类主轴之间的欧氏距离的平方来度量

- $$\Delta(y, K_j) = \left[(y - m_j) - U_j U_j^T (y - m_j) \right]^T \left[(y - m_j) - U_j U_j^T (y - m_j) \right]$$

10.5 模糊聚类算法

10.5.1 模糊集的基本知识

- **确定集合**（脆集合）：元素或者属于或者不属于集合
- **模糊集合**：元素以一定的程序属于某集合 —— 适于表达自然语言变量和常识性知识
- 几种叫法：模糊集、模糊逻辑、模糊数学、模糊系统、模糊技术、模糊方法
- 与其它技术相结合：模糊神经网络、模糊控制、模糊模式识别

10.5.1 模糊集的基本知识

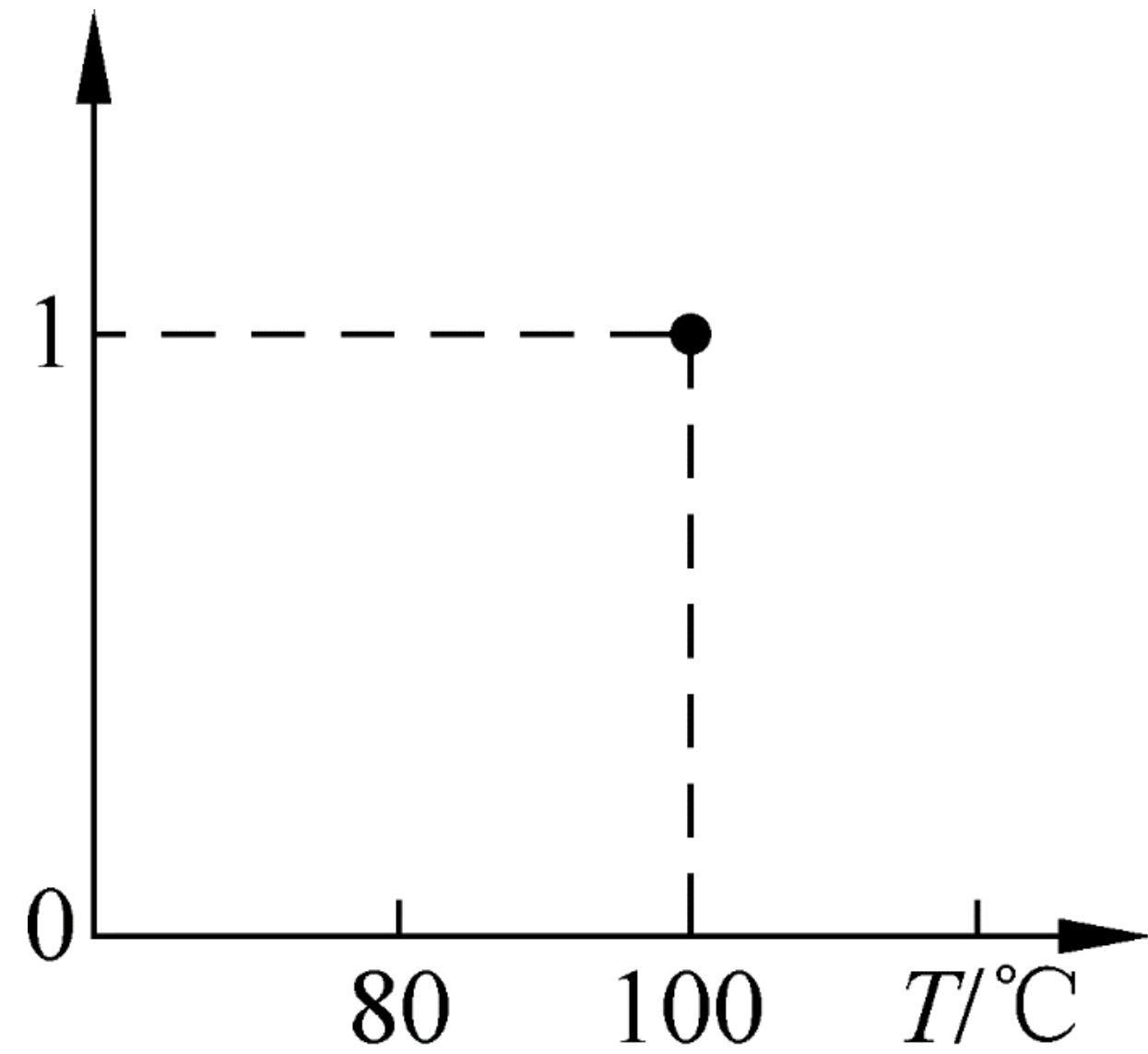
- 隶属度函数 $\mu_A(x)$: x 属于集合 A 的程度
 - 自变量: 所有可能属于 A 的对象, 空间 $X = \{x\}$
 - 值域: $[0, 1]$, $0 \leq \mu_A(x) \leq 1$,
 $\mu_A(x) = 1 \Leftrightarrow x \in A$
 $\mu_A(x) = 0 \Leftrightarrow x \notin A$
- 模糊集合: 一个定义在 $X = \{x\}$ 上的隶属度函数就定义了一个模糊集合 A , 或称定义在空间 $X = \{x\}$ 上的模糊子集, 表示为:

$$A = \left\{ (\mu_A(x_i), x_i) \right\} \text{ 或 } A = \bigcup_i \mu_i/x_i$$

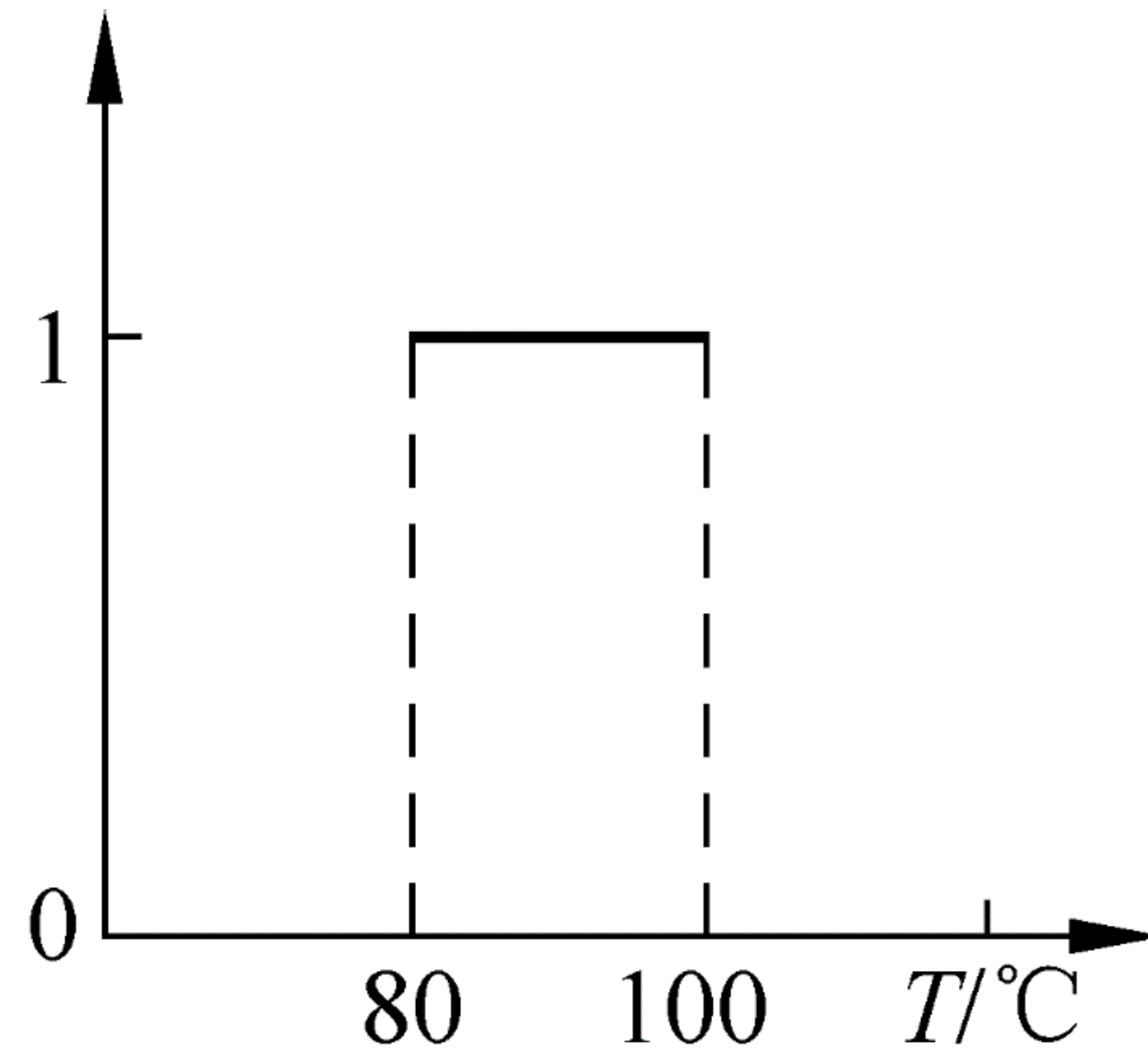
- 模糊集 A 的支持集: $S(A) = \{x, x \in X, \mu_A(x) > 0\}$

10.5.1 模糊集的基本知识

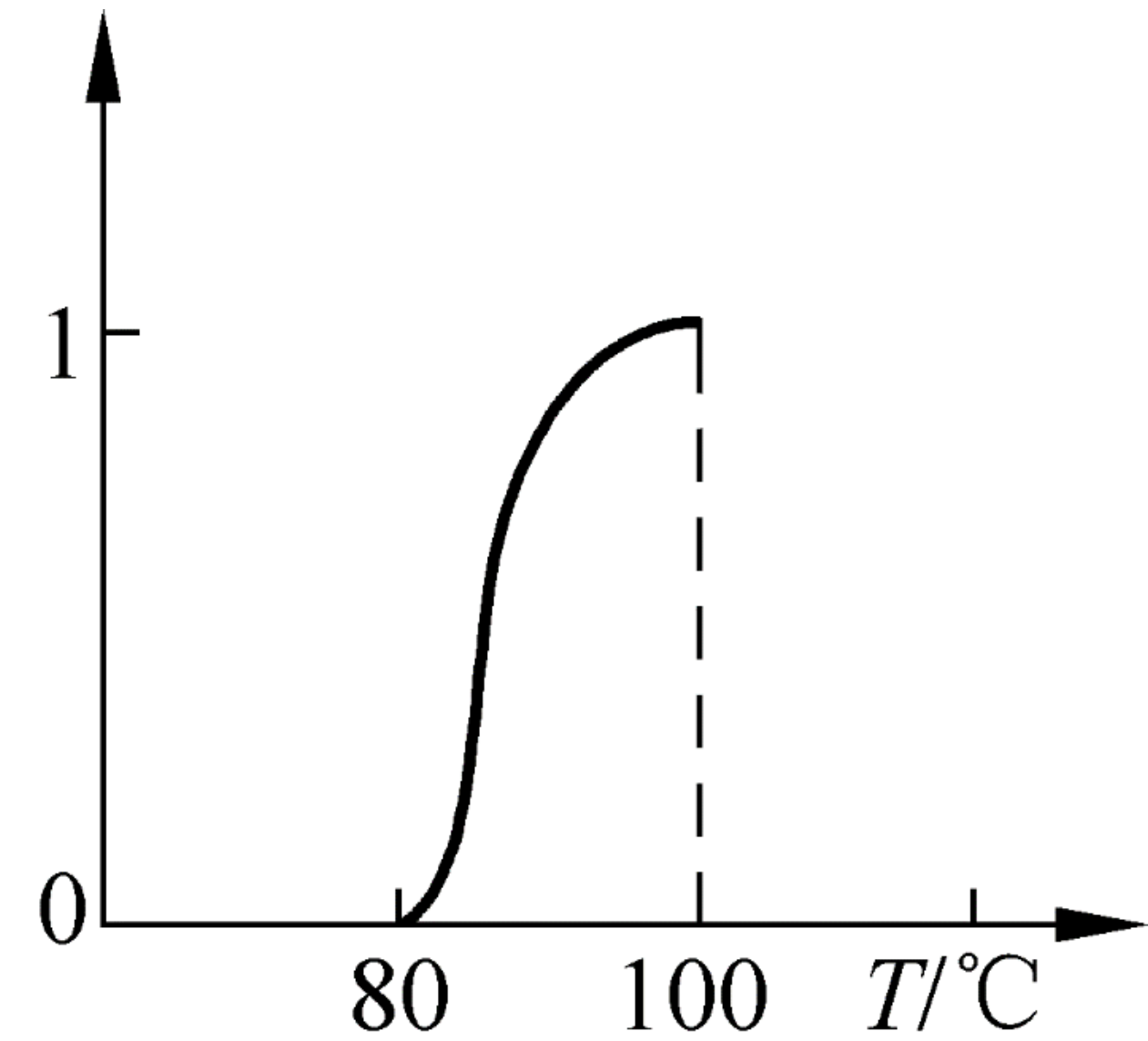
- 例：“开水”概念的表达



(a)



(b)



(c)

10.5.2 模糊C均值算法

- C均值方法

- 把 n 个样本划分到 c 个类中，使各样本与其所在类的均值的误差平方和最小

$$J_e = \sum_{i=1}^c \sum_{y \in I_i} \|y - m_i\|^2$$

- 把硬分类变成模糊分类，即得模糊 c 均值方法

10.5.2 模糊C均值算法

- 模糊C均值方法

- 符号约定

- ✓ 样本集 $\{x_i, i = 1, 2, \dots, n\}$

- ✓ 聚类中心 $m_j, j = 1, 2, \dots, c$

- ✓ $\mu_j(x_i)$: 第 i 个样本对于第 j 类的隶属度函数

- 聚类损失函数

$$J_f = \sum_{j=1}^c \sum_{i=1}^n \left[\mu_j(x_i) \right]^b \left\| x_i - m_j \right\|^2$$

其中, $b > 1$ 可控制聚类结果的模糊程度

10.5.2 模糊C均值算法

- 模糊聚类: $\min J_f$, 对不同的隶属度定义, 就得到不同的模糊聚类方法

- 模糊C均值: $\sum_{j=1}^c \mu_j(x_i) = 1, i = 1, \dots, n$

- 令 $\partial J_f / \partial m_j = 0$ 和 $\partial J_f / \partial \mu_j(x_i) = 0$, 可得

- $$m_j = \frac{\sum_{i=1}^n [\mu_j(x_i)]^b x_i}{\sum_{i=1}^n [\mu_j(x_i)]^b}, j = 1, \dots, c$$

- $$\mu_j(x_i) = \frac{\left(1 / \|x_i - m_j\|^2\right)^{1/(b-1)}}{\sum_{k=1}^c \left(1 / \|x_i - m_k\|^2\right)^{1/(b-1)}}, j = 1, \dots, c, i = 1, \dots, n$$

10.5.2 模糊C均值算法

- 模糊C均值算法 (FCM) :
 - 设定聚类数目 c 和常数 b
 - 初始化各聚类中心 m_j , $j = 1, \dots, c$ (可参考上两节中的方法)
 - 重复下面的运算, 直到各 $\mu_j(x_i)$ 稳定
 - ✓ 用当前的 m_j 计算 $\mu_j(x_i)$
 - ✓ 用当前的 $\mu_j(x_i)$ 计算 m_j
 - 如需要, 可对所有模糊聚类去模糊化

10.5.3 改进的模糊C均值算法 (AFC)

- 模糊C均值算法的缺点

- 由于 $\sum_{j=1}^c \mu_j(x_i) = 1$ ，故对某些野值（本应属于各类程度都很小），隶属度可能较大

- 改进的模糊C均值算法

- 放松归一化条件为 $\sum_{j=1}^c \sum_{i=1}^n \mu_j(x_i) = n$

- $$m_j = \frac{\sum_{i=1}^n [\mu_j(x_i)]^b x_i}{\sum_{i=1}^n [\mu_j(x_i)]^b}, \quad j = 1, \dots, c \text{ (不变)}$$

- $$\mu_j(x_i) = \frac{n \left(1 / \|x_i - m_j\|^2 \right)^{1/(b-1)}}{\sum_{k=1}^c \sum_{l=1}^n \left(1 / \|x_l - m_k\|^2 \right)^{1/(b-1)}}, \quad j = 1, \dots, c, \quad i = 1, \dots, n$$

10.5.3 改进的模糊C均值算法 (AFC)

- 改进的模糊C均值算法步骤与模糊C均值相同
 - 改进的模糊C均值算法所得到的 $\mu_j(x_i)$ 可能大于1，不是严格意义下的隶属度函数。必要时可做归一化
 - 改进的模糊C均值算法有更好的鲁棒，且对给定的聚类数目不十分敏感
 - 但有时可能会出现一个类中只包含一个样本的情况，可通过在距离计算中引入非线性使之不会小于某值来改进
 - 改进的模糊C均值和C均值一样，依赖于初值

10.6 分级聚类算法

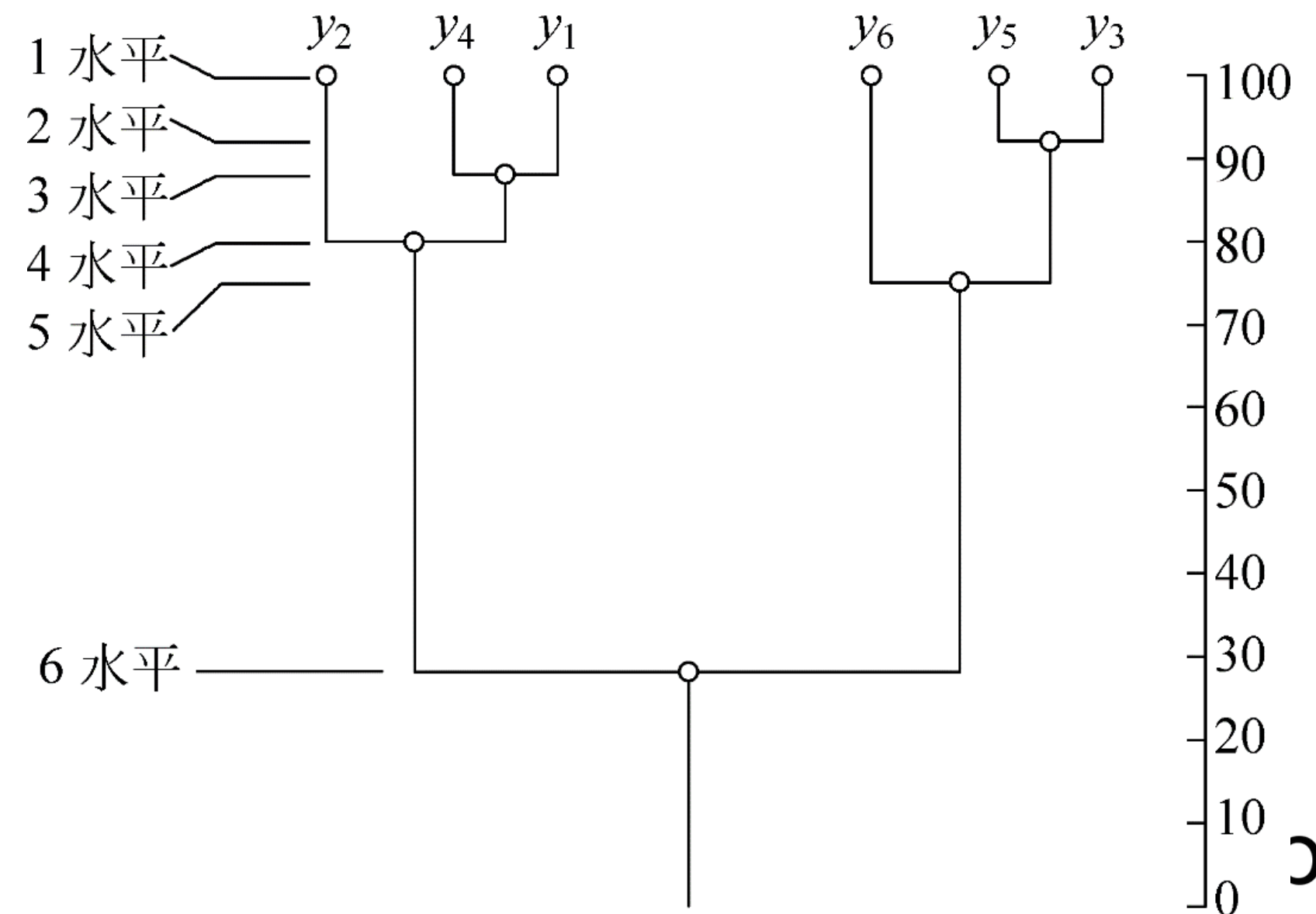
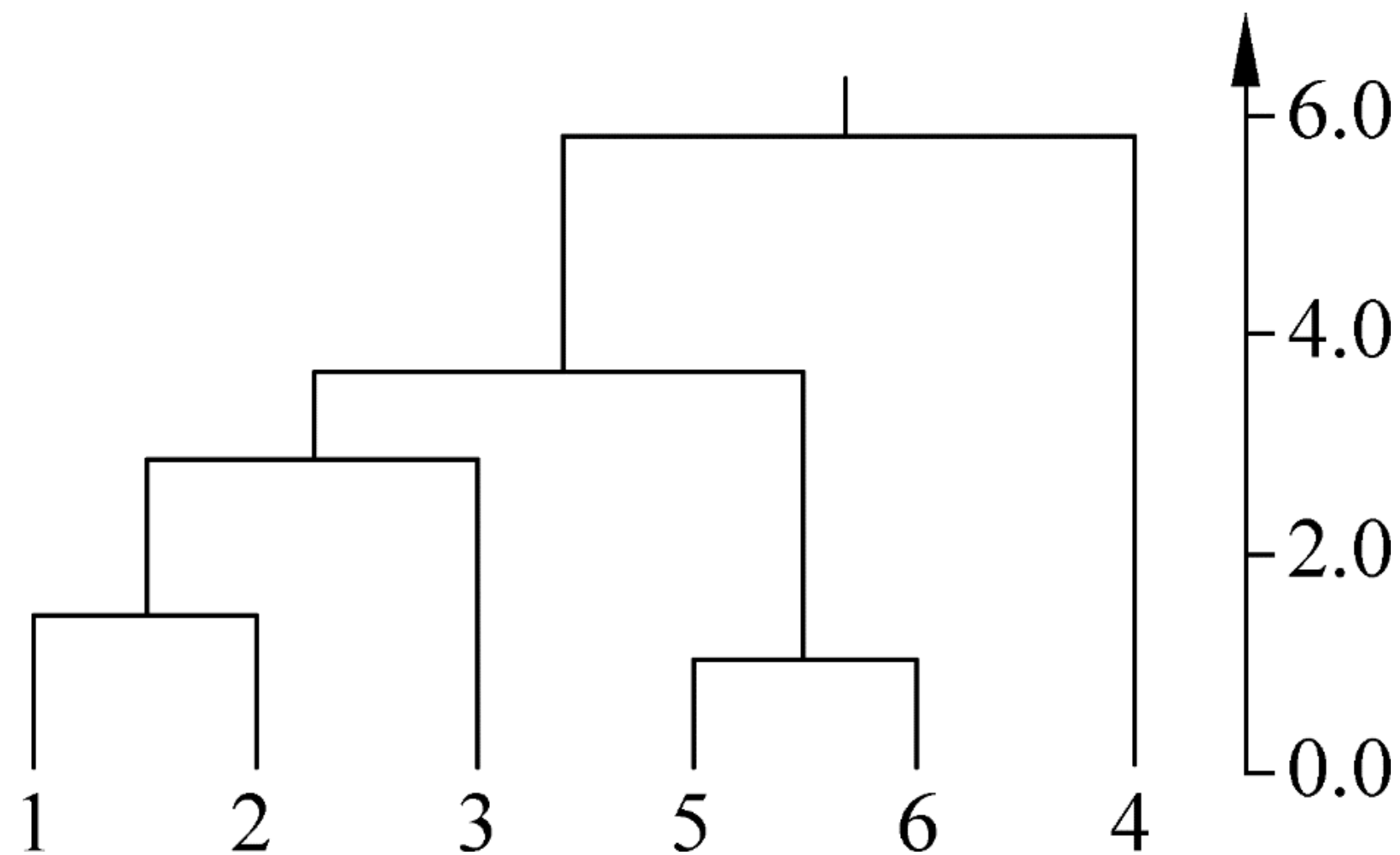
- **思想**

- 从各类只有一个样本点开始，逐级合并，每级只合并两类，直到最后所有样本都归到一类
- 聚类过程中逐级考查类间相似度，依此决定类别数

10.6 分级聚类算法

- 思想

- 树枝长度：反映结点/树枝之间的相似度或距离
- 树枝位置：在不改变树的结构的情况下可以任意调整，调整方法需研究
- 距离/相似性度量：多种选择，如欧式距离、相关、City Block、...



10.6 分级聚类算法

- 距离（相似性度量）
 - 样本之间的度量
 - 聚类之间的度量
- 两种聚类策略
 - bottom-up: agglomerative algorithm
 - top-down: divisive algorithm

10.6 分级聚类算法

- 算法（从底向上）

- (1) 初始化，每个样本形成一类
- (2) 把相似性最大（距离最小）的两类合并
- (3) 重复（2），直到所有样本合并为两类

- 常用的几种类间似性度量

- 最近距离（single-link）， $\Delta(\Gamma_i, \Gamma_j) = \min_{y \in \Gamma_i, \tilde{y} \in \Gamma_j} \delta(y, \tilde{y})$

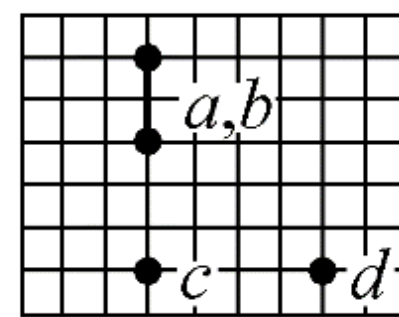
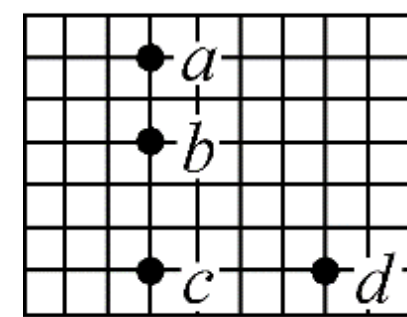
- 最远距离（complete-link）， $\Delta(\Gamma_i, \Gamma_j) = \max_{y \in \Gamma_i, \tilde{y} \in \Gamma_j} \delta(y, \tilde{y})$

- 均值距离（average-link）， $\Delta(\Gamma_i, \Gamma_j) = \delta(m_i, m_j)$

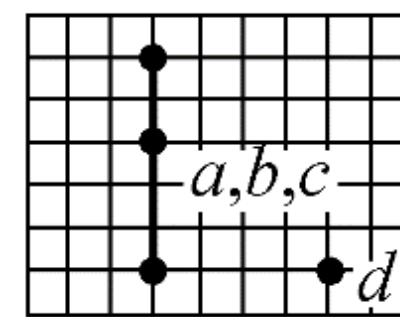
10.6 分级聚类算法

- 不同相似性度量对结果的影响

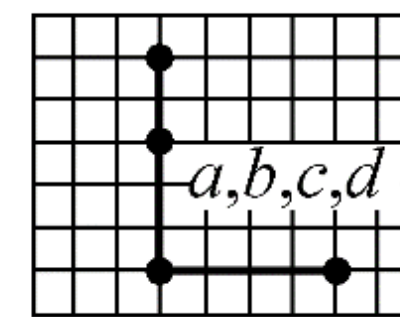
最近距离连接



(1)



(2)



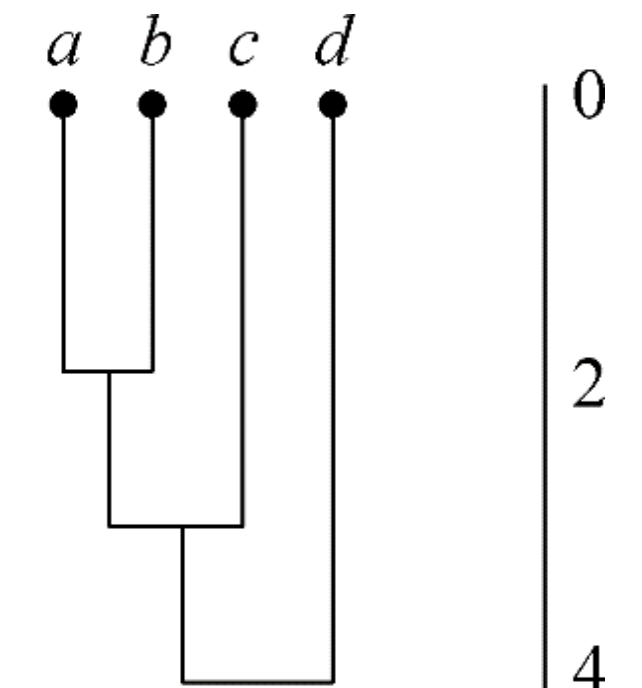
(3)

	b	c	d
a	2	5	6
b		3	5
c			4

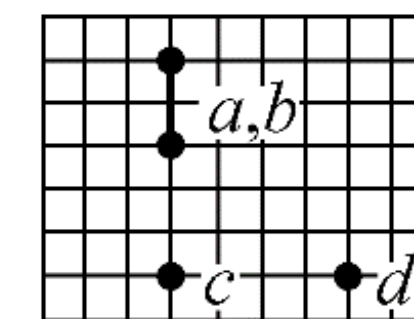
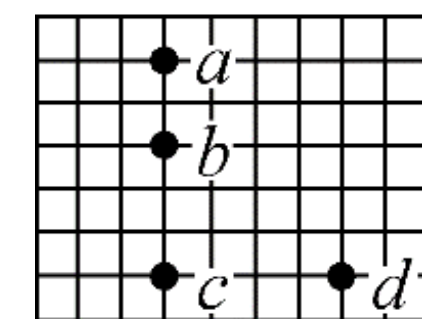
	b	c	d
a	2	5	6
b		3	5
c			4

	c	d
a,b	3	5
c		4

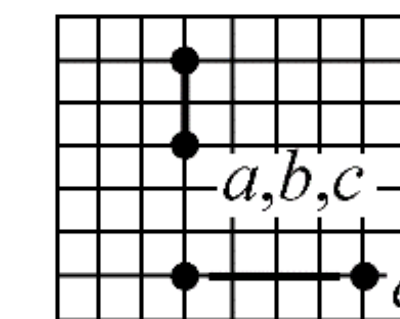
	d
a,b,c	4



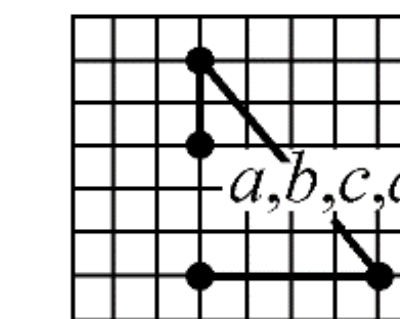
最远距离连接



(1)



(2)



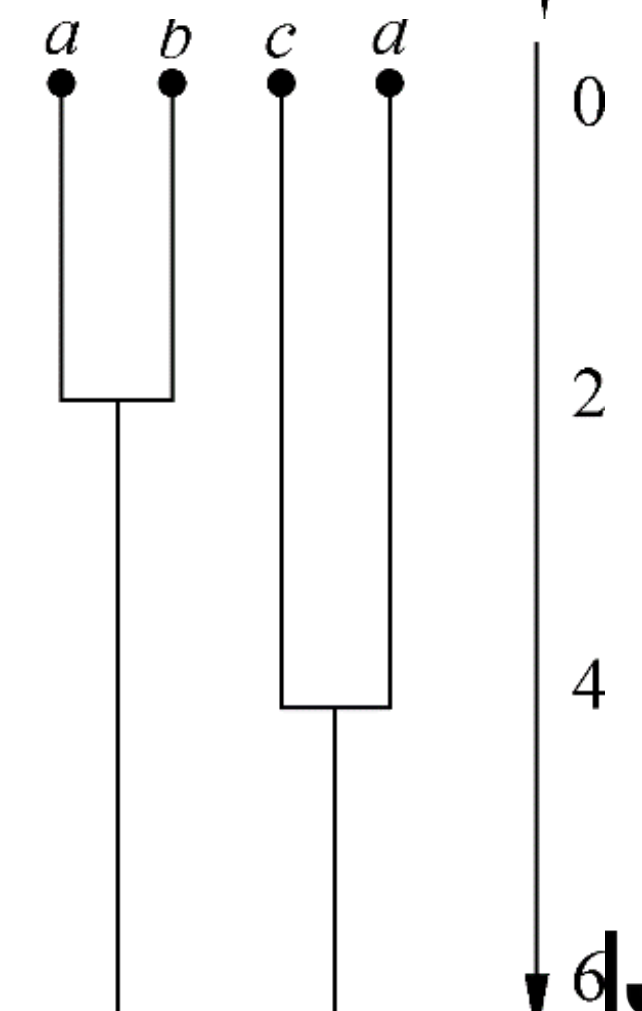
(3)

	b	c	d
a	2	5	6
b		3	5
c			4

	b	c	d
a	2	5	6
b		3	5
c			4

	c	d
a,b	5	6
c		4

	c,d
a,b	6



↓ d JUST

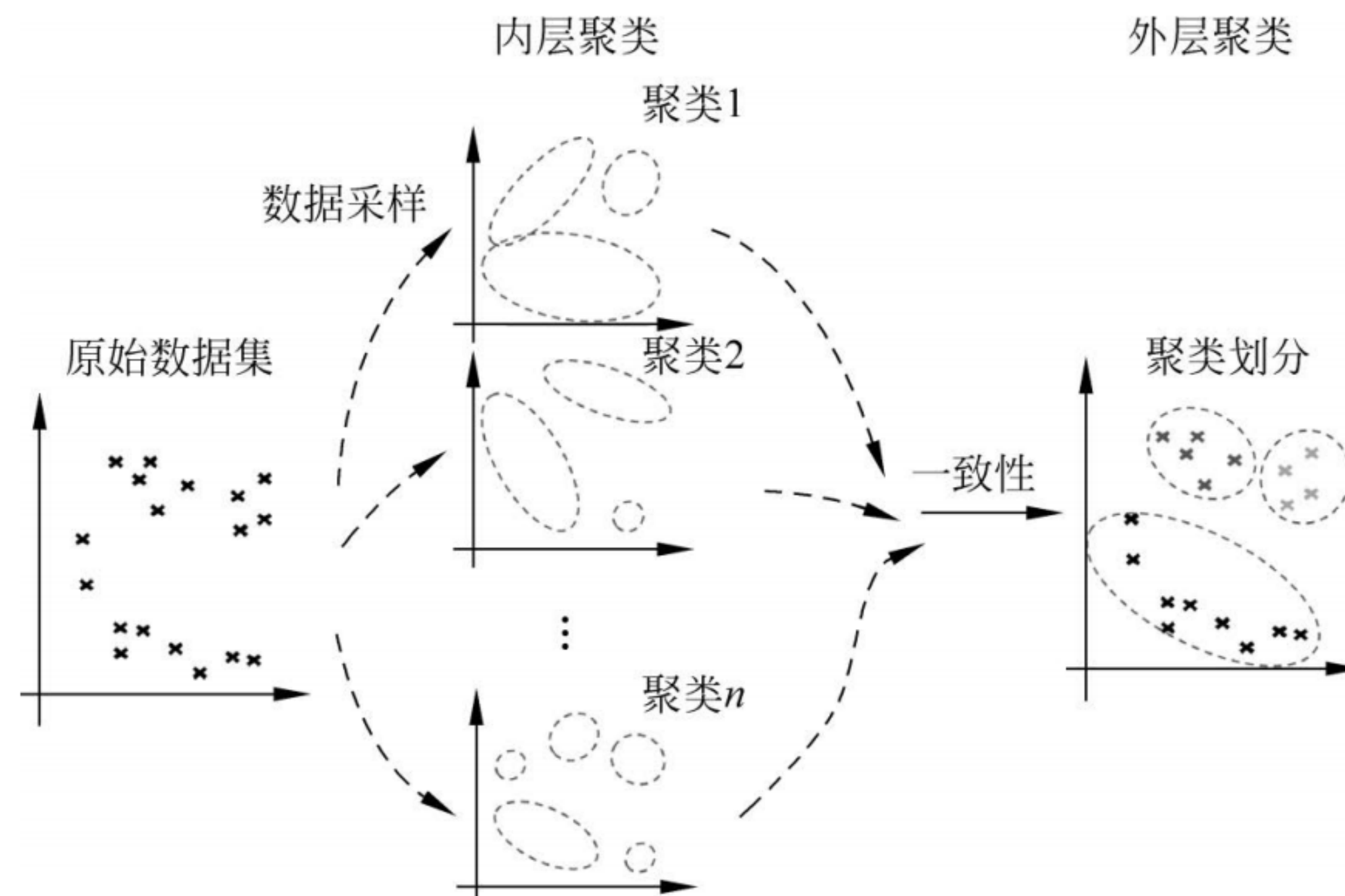
10.6 分级聚类算法

- 几点说明：
 - 分级聚类是一种局部搜索，对样本中的噪声敏感
 - 聚类树的画法不是唯一的。同一类中的两个分支可以左右互换而不改变聚类结果，但会改变树的外观和分析者的判断

10.7 一致聚类方法

- 基本思路

通过不同的数据抽样和不同方法进行多次聚类，再对结果进行合并，将在大多数结果中一致的结果作为最终的聚类划分依据。

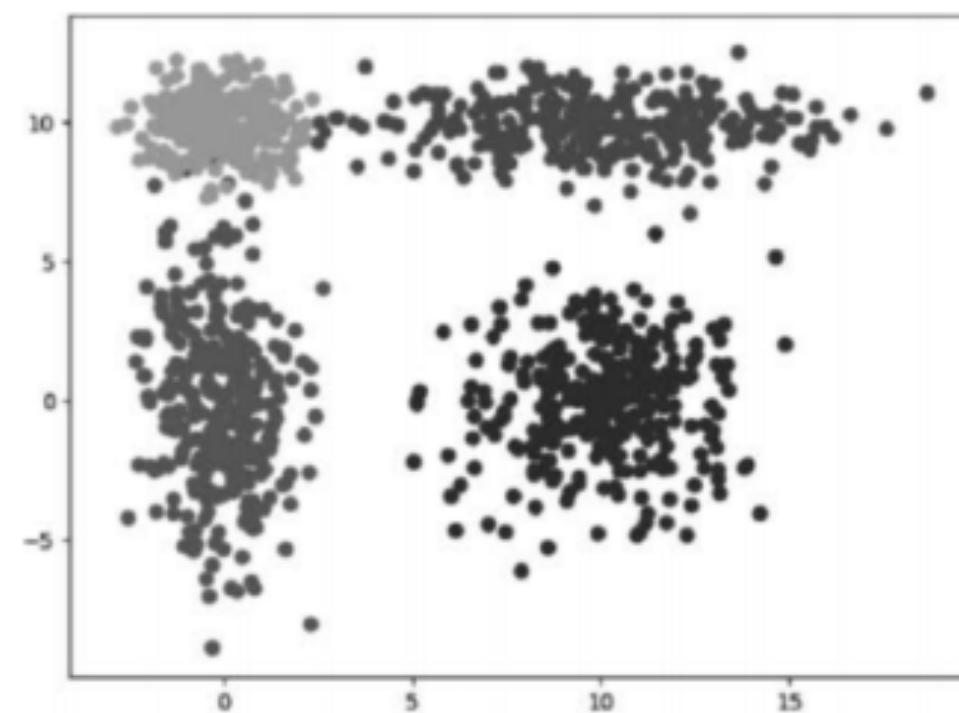


10.7 一致聚类方法

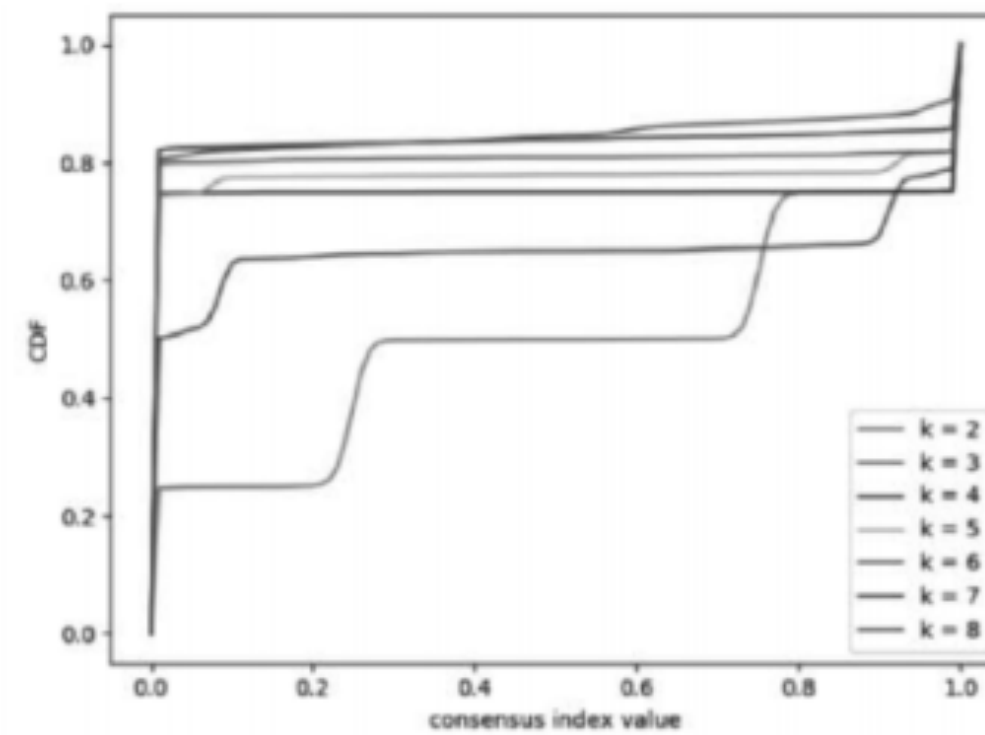
- 原始数据集: $D = \{x_1, x_2, \dots, x_N\}$
- 数据子集: $D^{(1)}, \dots, D^{(S)}$ (S次采样, 无放回重采样)
- 内层聚类分析: 用不同的 K 对每一个数据子集进行聚类
 - 连接矩阵 $M_{(s)}^{(K)}$: 表示在 $D^{(s)}$ 上将数据聚类为 K 类时, 样本 i 和样本 j 处于同一个类中, 则 $M_{(s)}^{(K)}(i, j) = 1$; 否则 $M_{(s)}^{(K)}(i, j) = 0$
 - 聚类数为 K 时的整体一致性矩阵 $M^{(K)}$: $M^{(K)}(i, j) = \frac{\sum_s M_{(s)}^{(K)}(i, j)}{\sum_s I_s(i, j)}$, 其中 $I_{(s)}$ 为指示矩阵, 但样本 i 和 j 都出现在数据集 $D^{(s)}$ 中时, $I_s(i, j) = 1$

10.7 一致聚类方法

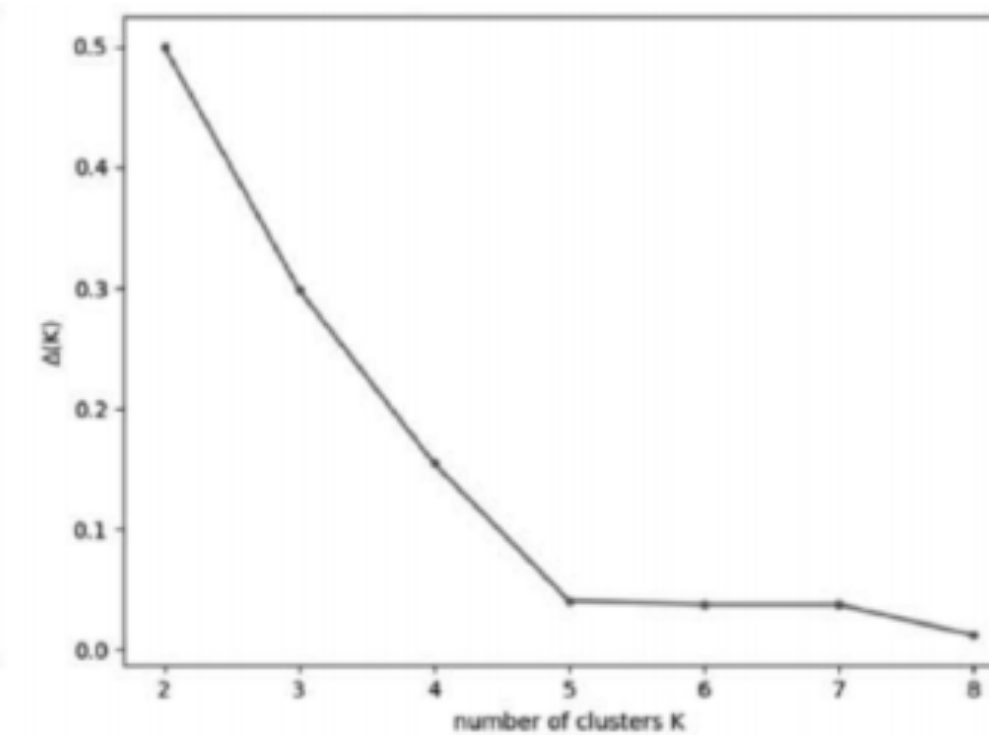
- 一致性矩阵 $M^{(K)}$ 是对称矩阵，元素取值在 $[0,1]$ 之间；可以可视化观察聚类结果的稳定性



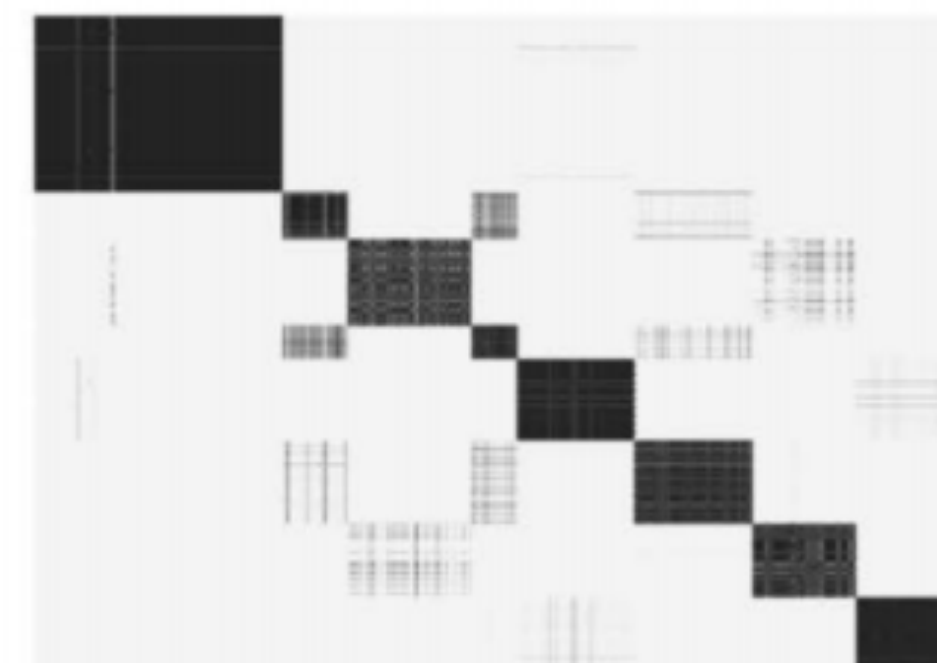
(a) 仿真数据
 $k=2$



(b) CDF曲线
 $k=4$



(c) $\Delta(K)$ 随 K 的变化
 $k=8$



(d) $M^{(K)}$ 矩阵的可视化分析

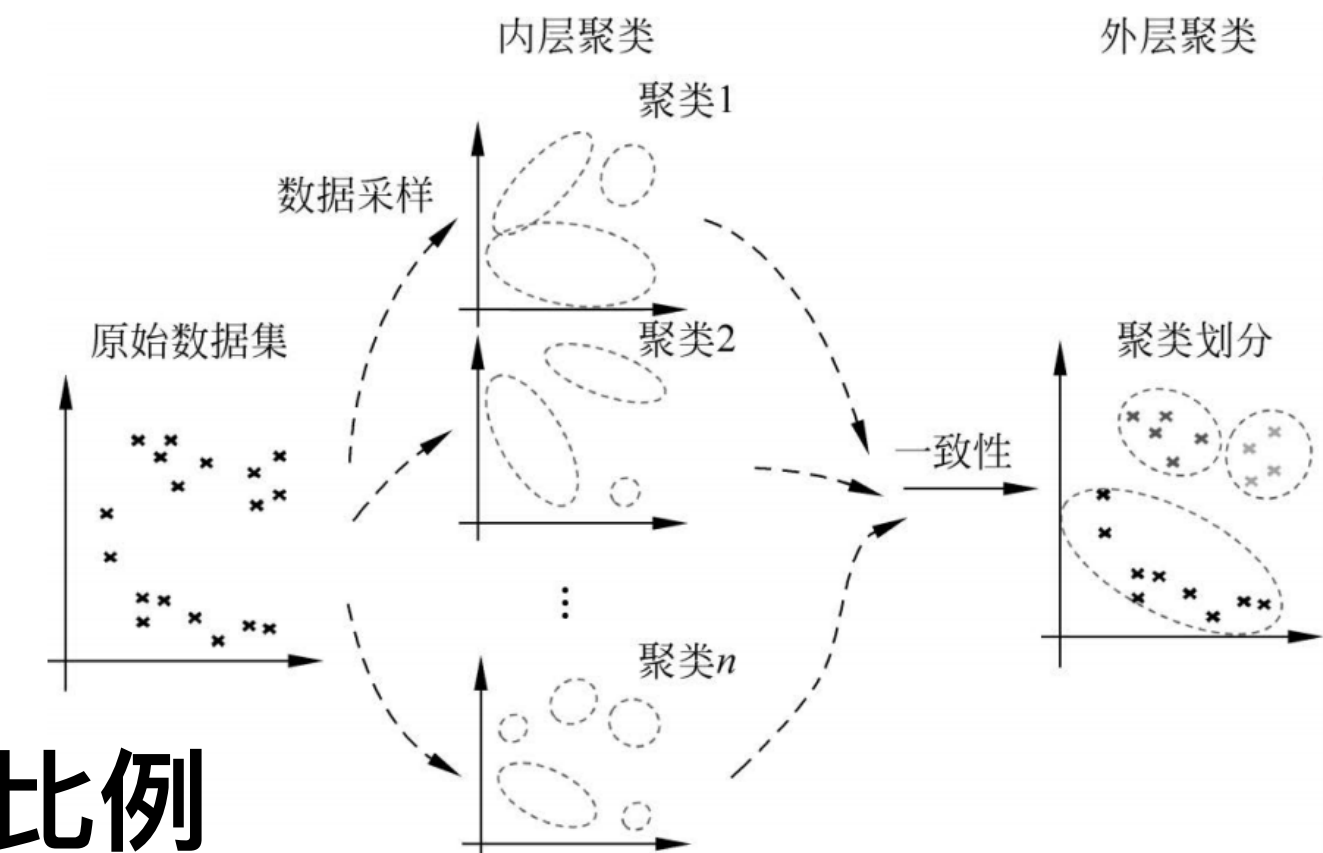
10.7 一致聚类方法

- 外层聚类分析

- 距离度量矩阵: $Dist^{(K)} = (1 - M^{(K)})$, 在此距离矩阵进行聚类 (如层次聚类), 将在大多数聚类结果中一致的结果最为最终的聚类划分依据, 得到最终的聚类结果

- 类别数 K

- CDF函数: $CDF^{(K)}(t) = \frac{\sum_{i < j} I \{ M^{(K)} \leq t \}}{N(N-1)/2}, t \in [0, 1]$



- 表示一致性矩阵 M 中取值小于阈值 t 的样本对占总样本对数量的比例
- 若一致性高, 则CDF的取值在 t 靠近0和1附近变化较大, 中间平稳。
- 若一致性差, 则CDF的取值随着 t 增大缓慢上升。

10.7 一致聚类方法

- 类别数 K

- 比较CDF函数曲线的线下面积AUC来确定聚类数 K

$$A(K) = \sum_{i=2}^{N(N-1)/2} [x_i - x_{i-1}] CDF^{(K)}(x_i)$$

其中 $\left\{ x_1, x_2, \dots, x_{\frac{1}{2}N(N-1)} \right\}$ 是对 $M^{(K)}$ 元素取值从小到大的排序

定义衡量 $A(K)$ 变化的指标: $\Delta(K) = \begin{cases} A(K), K = 2 \\ \frac{A(K+1) - \hat{A}^{(K)}}{A(K)}, K > 2 \end{cases}$, 其中

$\hat{A}^{(K)} = \max_{k \in \{2, \dots, K\}} A(k)$, 可以取 $\Delta(K)$ 趋近稳定前一时刻的 K