



河北师范大学软件学院  
Software College of Hebei Normal University

# 机器学习

非监督式机器学习

聚类(Clustering)

河北师范大学软件学院

2018.05.10

# 主要内容

## 1. 聚类的引入

## 2. 动态聚类

## 3. 基于模型的聚类

## 4. 密度聚类

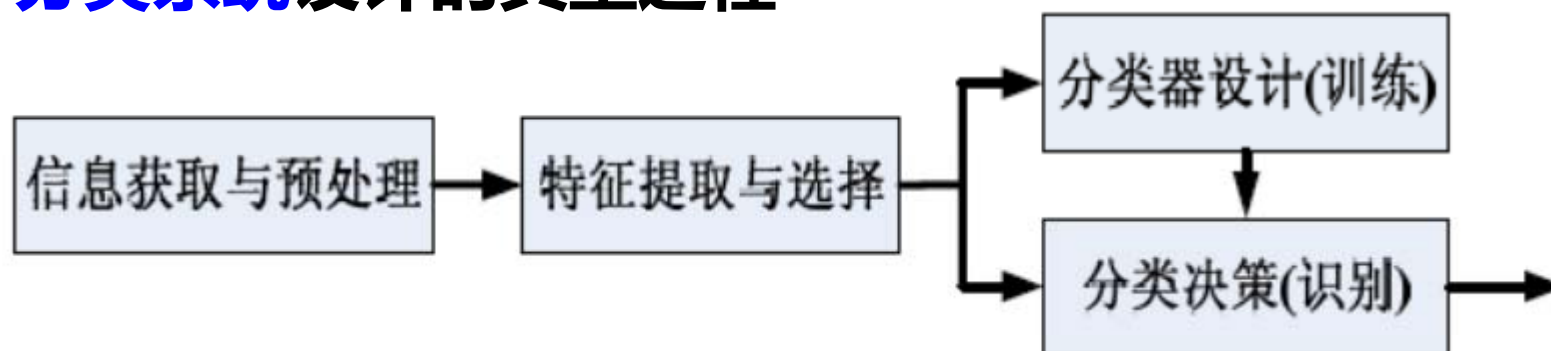
## 5. 层次聚类

## 6. 其它聚类算法



# 问题的引入

## 分类系统设计的典型过程



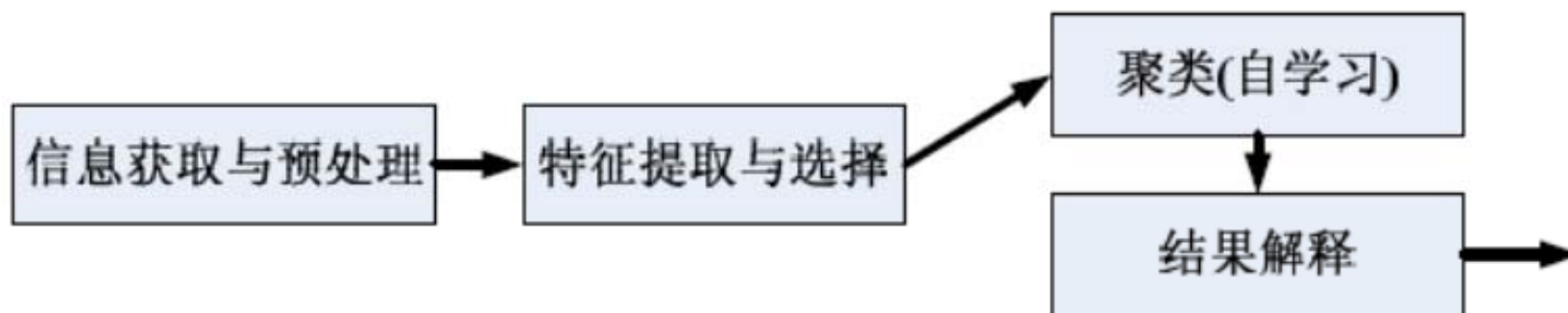
## 监督式学习

**不足：**需要大量已知类别标签的样本集；  
收集与标记费时、费力

**实际面临的问题：**样本数据多是无标签数据

# 非监督式学习(密度函数估计、聚类、降维...)

例：聚类系统设计的典型过程

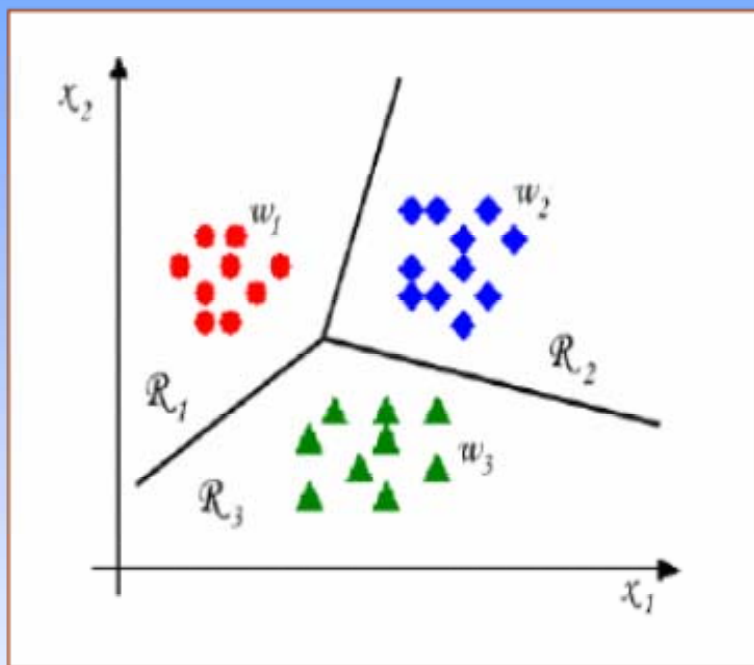


适用：

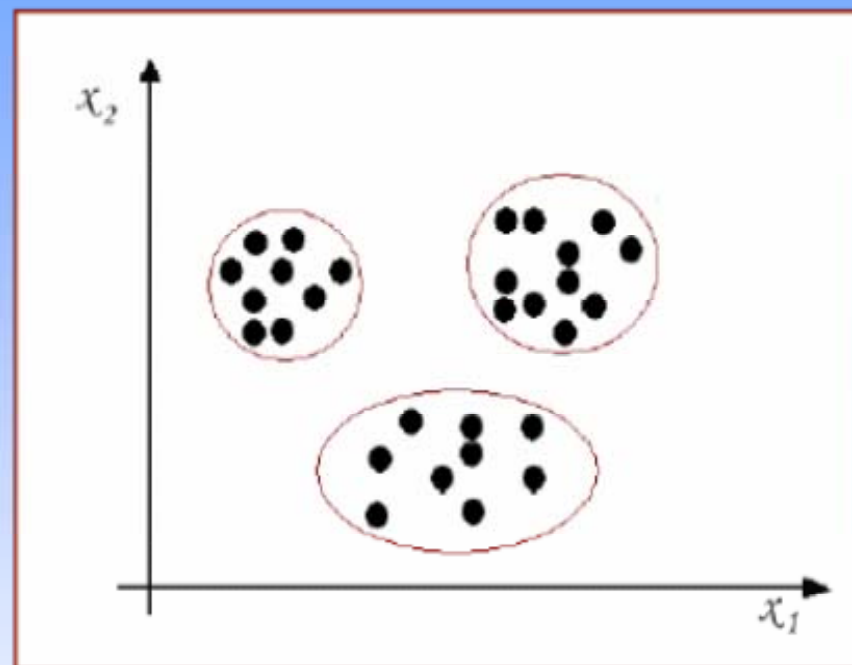
- (1) 大型数据挖掘：大量未标记数据训练分类器；  
人工标记分组结果
- (2) 揭示观测数据的内在结构特性
- (3) 是分类或其它学习任务的前驱阶段  
提取数据的基本特征，进一步用于分类...



## 分类 (Classification) 与 聚类 (Clustering)



Given labeled training patterns,  
construct decision boundaries  
or partition the feature space



Given some patterns, discover the  
underlying structure (categories)  
in the data

# 什么是聚类(clustering)



[http://home.dei.polimi.it/matteucc/Clustering/tutorial\\_html/](http://home.dei.polimi.it/matteucc/Clustering/tutorial_html/)

*To find a structure in a **collection** of **unlabeled data**.*

- *A loose definition of **clustering** could be “**the process of organizing objects into groups whose members are similar in some way**”.*
- *A **cluster** is therefore **a collection of objects which are “similar” between them and are “dissimilar” to the objects belonging to other clusters.***
- *High intra-class(cluster) similarity  
Low inter-class(cluster) similarity*

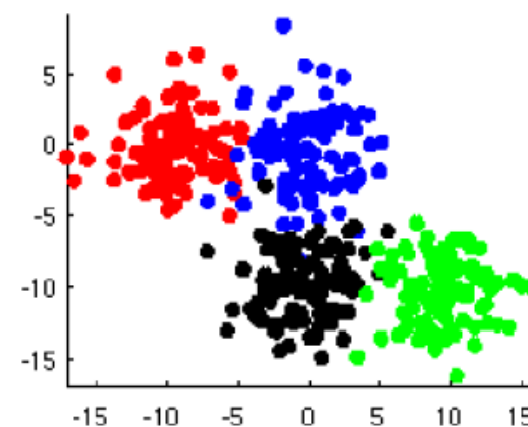
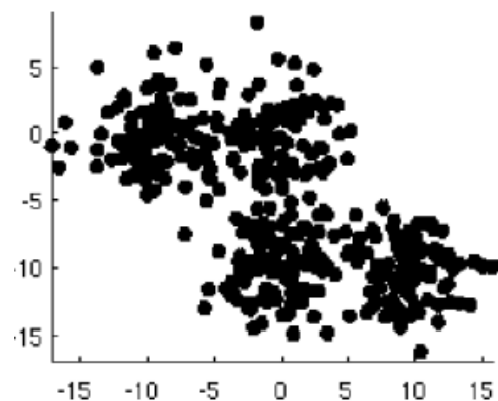
# 聚类问题的描述

**输入：**无标签数据集  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ ,  $\mathbf{x}_i = [x_{i1} \quad \dots \quad x_{id}]^T \in R^d$

要生成的簇的数目  $k$

**输出：** $k$ 个互不相交的簇  $\{C_l \mid l = 1, \dots, k\}$

$$\begin{cases} C_i \cap_{i \neq j} C_j = \Phi \\ \bigcup_{l=1}^k C_l = D \end{cases}$$



$\lambda_j$  -- 样本  $\mathbf{x}_j \in D$  的簇标记,  $j \in \{1, 2, \dots, k\}$

数据集  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  的标签集合  $\lambda = \{\lambda_1, \dots, \lambda_m\}$



聚类有很多典型应用，如：

- 相似功能的基因分组
- 相似政见的个体划分
- 相似主题的文档划分

...



# 聚类任务的遵循步骤及有关问题

## ➤ 特征选择及样本描述

选择什么样的特征？是否需要规范化预处理？

## ➤ 近邻测度

如何度量样本之间的“相似”或“相异”

## ➤ 聚类准则

依赖于专家对“可判别”的解释，聚类准则应以蕴涵于数据集内类的类型为基础。

## ➤ 聚类算法

选择特定的算法，用于揭示数据集的聚类结构

## ➤ 聚类性能评价、结果的解释





## 近邻测度与聚类准则

### A. 向量(点)之间的测度

#### (1) Dissimilarity Metric

-- dist

$$d: X \times X \rightarrow \mathbb{R}$$

对于  $\forall x, y, z \in X$ ,  $dist(\cdot, \cdot)$  须满足:

$$\left\{ \begin{array}{l} \text{非负性: } \forall x, y, dist(x, y) \geq 0 \\ \text{同一性: } dist(x, y) = 0 \text{ 当且仅当 } x = y \\ \text{对称性: } dist(x, y) = dist(y, x) \\ \text{直递性: } dist(x, y) \leq dist(x, z) + dist(y, z) \end{array} \right.$$

#### (2) Similarity Metric -- s

$$s: X \times X \rightarrow \mathbb{R}$$

对于  $\forall x, y, z \in X$ ,  $s(\cdot)$  须满足:

$$\left\{ \begin{array}{l} \exists s_0 \quad -\infty < s(x, y) \leq s_0 < +\infty \quad \forall x, y \\ s(x, x) = s_0 \\ s(x, y) = s(y, x) \\ s(x, z) s(y, z) \leq [s(x, z) + s(y, z)] s(x, y) \end{array} \right.$$

## 例：向量之间的几种典型距离度量

## 有序属性之间的距离

对于  $\forall \mathbf{x}, \mathbf{y} \in X \subset \mathbf{R}^n$       $\mathbf{x} = [x_1 \ \cdots \ x_n]^T$     $\mathbf{y} = [y_1 \ \cdots \ y_n]^T$

### A. 闵可夫斯基距离 (*Minkowski distance*)

$$dist_{mk}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_p = \left( \sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}}$$

### B. 曼哈顿距离 (*Manhattan distance*)     $dist_{man}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_1 = \sum_{i=1}^n |x_i - y_i|$

### C. 欧式距离 (*Euclidean distance*)     $dist_{ed}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2 = \left( \sum_{i=1}^n |x_i - y_i|^2 \right)^{\frac{1}{2}}$

### D. 切氏距离 (*Chebyshev distance*)     $dist_{che}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_\infty = \max_{1 \leq i \leq n} |x_i - y_i|$

### E. 马氏距离 (*Mahalanobis distance*)     $d_{mah}(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T \Sigma^{-1} (\mathbf{x} - \mathbf{y})}$

### F. Canberra距离 (*Lance 距离, Williams 距离*)

$$dist_{cam}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n \frac{|x_i - y_i|}{|x_i + y_i|} \quad x_i, y_i \geq 0 \text{ 且 } x_i + y_i \neq 0$$



## 例：向量之间的几种典型距离度量



河北师范大学软件学院  
Software College of Hebei Normal University

$m_{u,a}$  —— 属性 $u$ 上取值为 $a$ 的样本数

$m_{u,a,i}$  —— 第 $i$ 簇中，属性 $u$ 上取值为 $a$ 的样本数

$k$  —— 样本集划分的聚类簇数目

无序属性之间的距离

[1] 无序属性 $u$ 上两个离散值 $a, b$ 之间的**VDM距离** (*Value Difference Metric*)

$$VDM_p(a, b) = \sum_{i=1}^k \left| \frac{m_{u,a,i}}{m_{u,a}} - \frac{m_{u,b,i}}{m_{u,b}} \right|^p$$

[2] 基于混合属性 ( $n_c$  个有序属性以及  $n - n_c$  个无序属性) 的距离

$$MinkovDM_p(\mathbf{x}_i, \mathbf{x}_j) = \left( \sum_{u=1}^{n_c} |x_{iu} - x_{ju}|^p + \sum_{u=n_c+1}^n VDM_p(x_{iu}, x_{ju}) \right)^{\frac{1}{p}}$$

[3] **加权闵可夫斯基距离** (*Weighted Minkowski distance*)

$$dist_{wmk}(\mathbf{x}_i, \mathbf{x}_j) = \left( \sum_{u=1}^n w_u |x_{iu} - x_{ju}|^p \right)^{\frac{1}{p}} \quad w_u \geq 0, \quad \sum_{u=1}^n w_u = 1$$

## 例：向量之间的相似性度量

对于  $\forall \mathbf{x}, \mathbf{y} \in X \subset \mathbf{R}^n$       $\mathbf{x} = [x_1 \quad \cdots \quad x_n]^T$       $\mathbf{y} = [y_1 \quad \cdots \quad y_n]^T$

**A. 内积**      $s_{inner}(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = \sum_{i=1}^n x_i y_i$

**B. 余弦相似度**      $s_{cosine}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$

**C. Pearson** 相关系数      $r_{Pearson}(\mathbf{x}, \mathbf{y}) = \frac{(\mathbf{x} - \bar{x}\mathbf{1})^T (\mathbf{y} - \bar{y}\mathbf{1})}{\|\mathbf{x} - \bar{x}\mathbf{1}\| \|\mathbf{y} - \bar{y}\mathbf{1}\|}$

其中      $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$       $\bar{y} = \frac{\sum_{i=1}^n y_i}{n}$

**D. Tanimoto** 测度      $s_T(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - \mathbf{x}^T \mathbf{y}}$



### B. 样本与集合之间的测度

**方式1.**集合中所有样本对近邻测度 $\mathcal{D}(x, C)$ 均有贡献。

设观测样本 $x, y$ 之间近邻测度为 $\mathcal{D}(x, y)$

则样本 $x$ 与聚类(或簇) $C$ 之间的近邻函数 $\mathcal{D}(x, C)$ ,可以是:

**最大近邻函数**  $\mathcal{D}_{\max}^{ps}(x, C) = \max_{y \in C} \mathcal{D}(x, y)$

**最小近邻函数**  $\mathcal{D}_{\min}^{ps}(x, C) = \min_{y \in C} \mathcal{D}(x, y)$

**平均近邻函数**  $\mathcal{D}_{\text{avg}}^{ps}(x, C) = \frac{1}{n_C} \sum_{y \in C} \mathcal{D}(x, y)$   $n_C$ 为集合 $C$ 的势

**方式2.**近邻性以样本 $x$ 与集合 $C$ 的表示之间的近邻性度量。

--集合的表示通常有点、超平面、超球面等。

## C. 两集合之间的近邻函数



河北师范大学软件学院  
Software College of Hebei Normal University

设两观测样本 $x, y$ 之间近邻测度为 $\mathcal{D}(x, y)$

对于给定的两个向量集合 $D_i, D_j$ , 近邻函数常见:

**最大近邻函数**  $\mathcal{D}_{\max}^{ss}(D_i, D_j) = \max_{x \in D_i, y \in D_j} \mathcal{D}(x, y)$

**最小近邻函数**  $\mathcal{D}_{\min}^{ss}(D_i, D_j) = \min_{x \in D_i, y \in D_j} \mathcal{D}(x, y)$

**平均近邻函数**  $\mathcal{D}_{\text{avg}}^{ss}(D_i, D_j) = \frac{1}{n_{D_i} n_{D_j}} \sum_{x \in D_i} \sum_{y \in D_j} \mathcal{D}(x, y)$

其中 $n_{D_i} n_{D_j}$ 为集合 $D_i, D_j$ 的势

**均值近邻函数**  $\mathcal{D}_{\text{mean}}^{ss}(D_i, D_j) = \mathcal{D}(m_{D_i}, m_{D_j})$

$m_{D_i}, m_{D_j}$ 是关于集合 $D_i, D_j$ 的点描述, 如均值点、中值等。

**其它**  $\mathcal{D}_e^{ss}(D_i, D_j) = \sqrt{\frac{n_{D_i} n_{D_j}}{n_{D_i} + n_{D_j}}} \mathcal{D}(m_{D_i}, m_{D_j})$



## D. 聚类准则

- 类内距离准则
- 类间距离准则
- 基于类内、类间距离的准则函数
- 基于模式与类核的距离的准则函数



# 聚类的典型实现方式

- (1) 基于**模型**的方法
  - 基于**概率密度函数**估计的方法
- (2) 基于**数据**的方法
  - 基于样本间**距离**或**相似性度量**的聚类

## 共性：

- (1) 将样本集划分为若干子集；
- (2) 直接解决分类问题，或者将其作为“标注的”  
**训练样本集**进行后续监督式分类器设计

# 聚类(clustering)算法的分类

根据算法设计不同，可以是：

## *(1) Exclusive Clustering*

*K-means Clustering*

## *(2) Overlapping Clustering*

*Fuzzy C-means*

## *(3) Hierarchical Clustering*

*Agglomerative Clustering*

*Divisive Clustering*

## *(4) Probabilistic Clustering*

*Mixture of Gaussian*





## 1. 聚类的引入

## 2. 动态聚类

### 2.1 K-Means Clustering

### 2.2 学习向量量化

## 4. 基于模型的聚类

## 5. 密度聚类

## 6. 层次聚类

## 7. 其它聚类算法



### 动态聚类的三个要点:

- ① 选定某种**距离度量**作为样本间的相异性度量
- ② 确定某个**准则函数**, 用于评价聚类结果的质量
- ③ 给定初始分类方法, 以迭代算法找出使准则函数取极值的最好聚类结果

### 动态聚类过程:

多次迭代, 逐步调整类别划分, 最终使准则最优。

- C-Means Clustering
- Fuzzy C-Means Clustering
- ...

## K-Means Clustering 聚类问题描述

**输入：** 样本集  $D = \{x_1, \dots, x_m\}$ ,  $x_i = [x_{i1} \quad \dots \quad x_{id}]^T \in R^d$   
要生成的簇的数目  $k$

**输出：**  $k$  个互不相交的簇  $C = \{C_l \mid l = 1, \dots, k\}$

# 1. 聚类准则--“最小误差平方和”准则

将样本集 $D$ 划分成 $k$ 个簇： $D = C_1 \cup \dots \cup C_k$

误差平方和目标函数  $E(\mu_1, \dots, \mu_k) = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$

其中  $C_i$  --第 $i$ 个簇(聚类),  $i = 1, \dots, k$   $k \ll m$

$N_i$  --第 $i$ 个簇的样本数目

$\mu_i$  --第 $i$ 个聚类中心,  $\mu_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$

注： $\{\mu_i, i = 1, \dots, k\}$  --*codebook*

$\sum_{x \in C_i} \|x - \mu_i\|^2$  -- 簇内样本紧致性的度量



## 2. 算法实现

### 初始化方式

### 交叉迭代的动态更新

$$\left\{ \mu_i^{(j)}, \quad i = 1, \dots, k \right\}$$

$$\rightarrow \left\{ C_i^{(j+1)}, \quad i = 1, \dots, k \right\}$$

$$\rightarrow \left\{ \mu_i^{(j+1)}, \quad i = 1, \dots, k \right\}$$

$\rightarrow \dots$

### 算法终止条件



输入: 样本集  $D = \{x_1, x_2, \dots, x_m\}$ ;  
聚类簇数  $k$ .

过程:

- 初始化方式
- 交叉迭代的动态更新方式
- 迭代终止条件

```
1: 从  $D$  中随机选择  $k$  个样本作为初始均值向量  $\{\mu_1, \mu_2, \dots, \mu_k\}$ 
2: repeat
3:   令  $C_i = \emptyset$  ( $1 \leq i \leq k$ )
4:   for  $j = 1, 2, \dots, m$  do
5:     计算样本  $x_j$  与各均值向量  $\mu_i$  ( $1 \leq i \leq k$ ) 的距离:  $d_{ji} = \|x_j - \mu_i\|_2$ ;
6:     根据距离最近的均值向量确定  $x_j$  的簇标记:  $\lambda_j = \arg \min_{i \in \{1, 2, \dots, k\}} d_{ji}$ ;
7:     将样本  $x_j$  划入相应的簇:  $C_{\lambda_j} = C_{\lambda_j} \cup \{x_j\}$ ;
8:   end for
9:   for  $i = 1, 2, \dots, k$  do
10:    计算新均值向量:  $\mu'_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$ ;
11:    if  $\mu'_i \neq \mu_i$  then
12:      将当前均值向量  $\mu_i$  更新为  $\mu'_i$ 
13:    else
14:      保持当前均值向量不变
15:    end if
16:  end for
17: until 当前均值向量均未更新
```

输出: 簇划分  $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$



河北师范大学软件学院  
Software College of Hebei Normal University



## A. 聚类中心的初始化

### “代表点”选择

样本集**D**划分之前，先选择代表点作为初始聚类核心，再将其余样本初始分类。

迭代结果与初始代表点选择有关

几种“代表点”的选择方法：

- 经验选择
- “密度法”选择代表点
- 用 $(k - 1)$ 聚类划分求 $k$ 个代表点

## K-Means ++ 中更为有效的聚类中心初始化方式

1. Initialize an empty set  $\mathbf{M}$  to store the  $k$  centroids being selected.
2. Randomly choose the first centroid  $\mu^{(j)}$  from the input samples and assign it to  $\mathbf{M}$ .
3. For each sample  $x^{(i)}$  that is not in  $\mathbf{M}$ , find the minimum squared distance  $d(x^{(i)}, \mathbf{M})^2$  to any of the centroids in  $\mathbf{M}$ .
4. To randomly select the next centroid  $\mu^{(p)}$ , use a weighted probability distribution equal to 
$$\frac{d(\mu^{(p)}, \mathbf{M})^2}{\sum_i d(x^{(i)}, \mathbf{M})^2}.$$
5. Repeat steps 2 and 3 until  $k$  centroids are chosen.
6. Proceed with the classic k-means algorithm.

## K-Means ++ 中更为有效的聚类中心初始化算法

1. 聚类中心集合 $M$ 的初始化:  $M \leftarrow \Phi$

2. 在样本集 $D$ 内随机选择1个样本初始化聚类中心 $\mu_1$ :

$$M \leftarrow M \cup \{\mu_1\}$$

3. *for*  $i = 2, \dots, K$  *do*

3.1 对于样本集 $D$ 内任意不同于 $M$ 内  $i-1$ 个中心的样本  $x$ ,  
计算该样本关于 $M$ 内 $i-1$ 个中心的距离平方最小值,

$$\text{记为 } d(x, M) = \min_{j=1, \dots, i-1} [d(x, \mu_j)]^2$$

3.2 按照概率  $\frac{d(x^*, M)}{\sum_{x_i \in D \setminus M} d(x_i, M)}$  随机选择一个样本 $x^*$ , 将其

作为初选的第 $i$ 个聚类中心 $\mu_i$

3.3  $M \leftarrow M \cup \{\mu_i\}$

4. 返回 $M$

## B. 聚类数 $k$ 的确定



河北师范大学软件学院  
Software College of Hebei Normal University

通常要求事先给定**聚类数 $k$** 。

若类别数目未知，可按如下方法确定

(a) 一般根据领域先验知识确定；

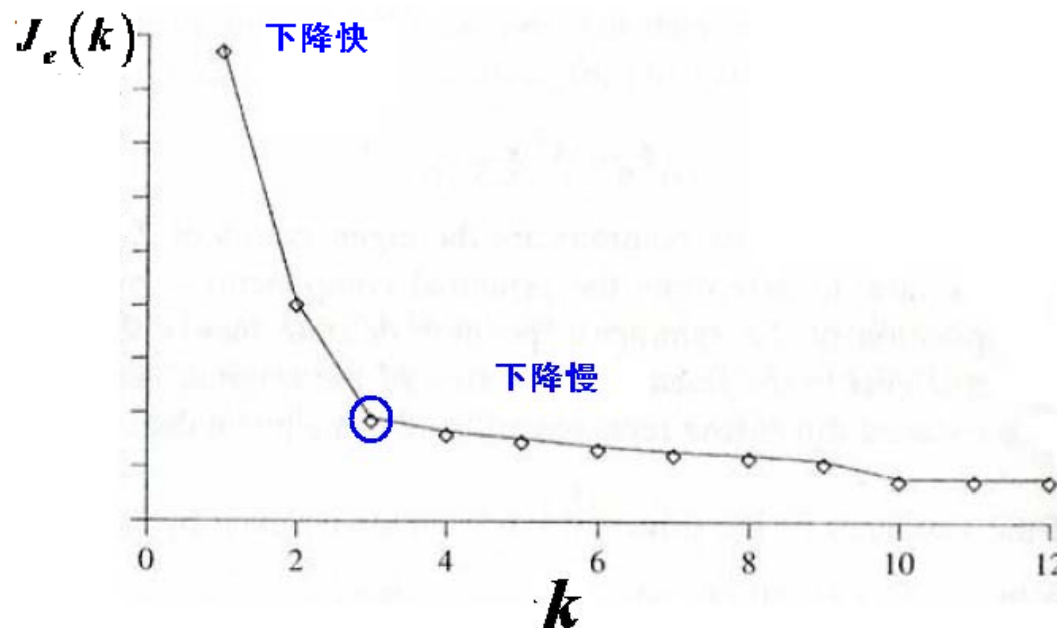
(b) 实验确定：

令 $k = 1, 2, 3, \dots$ ，分别进行聚类，得 $J_e(k)$ ，

绘制 $J_e(k)$ - $k$ 曲线图；

找出拐点，对应聚类数目为最终类别数。

该方法并不总是有效。



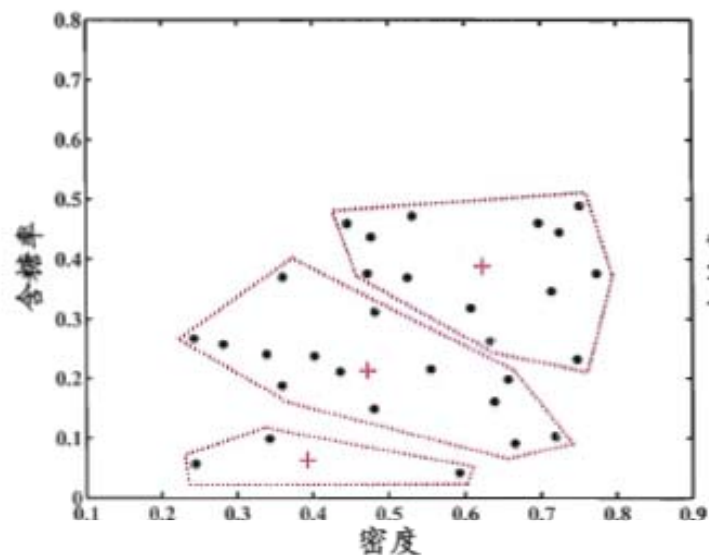
## C. 其它

是否需要~~进行~~样本集的规范化预处理？--**是**

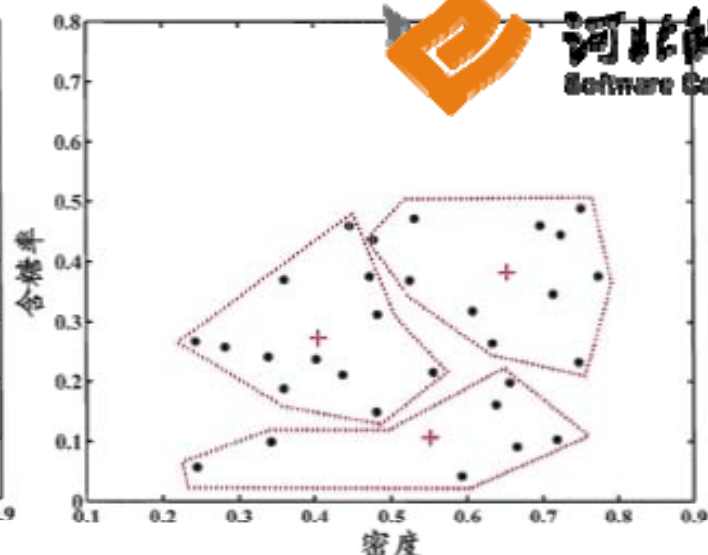
是否需要重复多次？

--从多个聚类结果中选择最好的那个

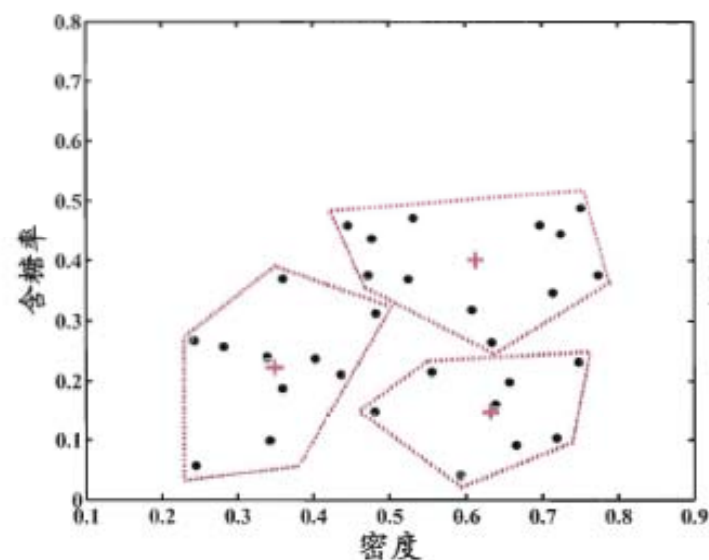
任意形状的聚类都可处理吗？--**不**



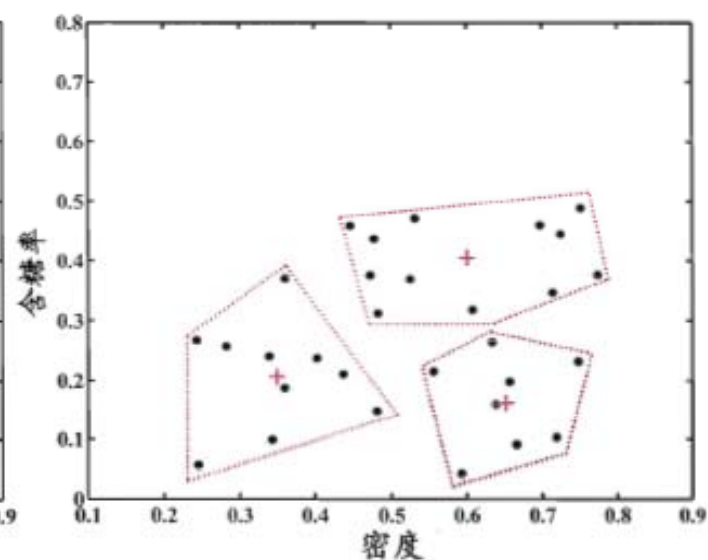
(a) 第一轮迭代后



(b) 第二轮迭代后



(c) 第三轮迭代后



(d) 第四轮迭代后

图 9.3 西瓜数据集 4.0 上  $k$  均值算法( $k = 3$ )在各轮迭代后的结果. 样本点与均值向量分别用“●”与“+”表示, 红色虚线显示出簇划分.



## 1. 聚类的引入

## 2. 性能度量

## 3. 动态聚类

### 3.1 K-Means Clustering

### 3.2 学习向量量化

## 4. 基于模型的聚类

## 5. 密度聚类

## 6. 层次聚类

## 7. 其它聚类算法

# 学习向量量化 (*Learning Vector Quantization, LVQ*)

## 问题描述

**输入：** 样本集  $D = \{(x_1, y_1), \dots, (x_m, y_m)\}$

要生成的原型向量数目  $q$

其中  $x_i = [x_{i1} \quad \dots \quad x_{id}]^T \in X \subset R^d$

$y_i \in Y$

**输出：**

(1)  $q$  个原型向量  $\{p_l \mid l = 1, \dots, q\}$ ,

每个原型向量代表一个聚类簇，簇的标记值  $t_l \in Y$

(2) 基于学得的原型向量，对原始空间进行划分



## 第一阶段, 原型向量的学习



输入: 样本集  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$ ;  
原型向量个数  $q$ , 各原型向量预设的类别标记  $\{t_1, t_2, \dots, t_q\}$ ;  
学习率  $\eta \in (0, 1)$ .

过程:

- 1: 初始化一组原型向量  $\{p_1, p_2, \dots, p_q\}$
- 2: **repeat**
- 3:   从样本集  $D$  随机选取样本  $(x_j, y_j)$ ;
- 4:   计算样本  $x_j$  与  $p_i$  ( $1 \leq i \leq q$ ) 的距离:  $d_{ji} = \|x_j - p_i\|_2$ ;
- 5:   找出与  $x_j$  距离最近的原型向量  $p_{i^*}$ ,  $i^* = \arg \min_{i \in \{1, 2, \dots, q\}} d_{ji}$ ;
- 6:   **if**  $y_j = t_{i^*}$  **then**
- 7:      $p' = p_{i^*} + \eta \cdot (x_j - p_{i^*})$
- 8:   **else**
- 9:      $p' = p_{i^*} - \eta \cdot (x_j - p_{i^*})$
- 10:   **end if**
- 11:   将原型向量  $p_{i^*}$  更新为  $p'$
- 12: **until** 满足停止条件

输出: 原型向量  $\{p_1, p_2, \dots, p_q\}$

若 $p_{i^*}$ 与 $x_j$ **类别一致**, 则 $p_{i^*}$ 更新后为  $p' = p_{i^*} + \eta(x_j - p_{i^*})$

$$\begin{aligned}\|p' - x_j\|_2 &= \|p_{i^*} + \eta(x_j - p_{i^*}) - x_j\|_2 \\ &= \|(1 - \eta)p_{i^*} - (1 - \eta)x_j\|_2 \\ &= (1 - \eta)\|p_{i^*} - x_j\|_2 < \|p_{i^*} - x_j\|_2\end{aligned}$$

若 $p_{i^*}$ 与 $x_j$ **类别不一致**, 则 $p_{i^*}$ 更新后为  $p' = p_{i^*} - \eta(x_j - p_{i^*})$

$$\begin{aligned}\|p' - x_j\|_2 &= \|p_{i^*} - \eta(x_j - p_{i^*}) - x_j\|_2 \\ &= \|(1 + \eta)p_{i^*} - (1 + \eta)x_j\|_2 \\ &= (1 + \eta)\|p_{i^*} - x_j\|_2 > \|p_{i^*} - x_j\|_2\end{aligned}$$

## 第二阶段, 基于原型向量的样本空间簇划分。

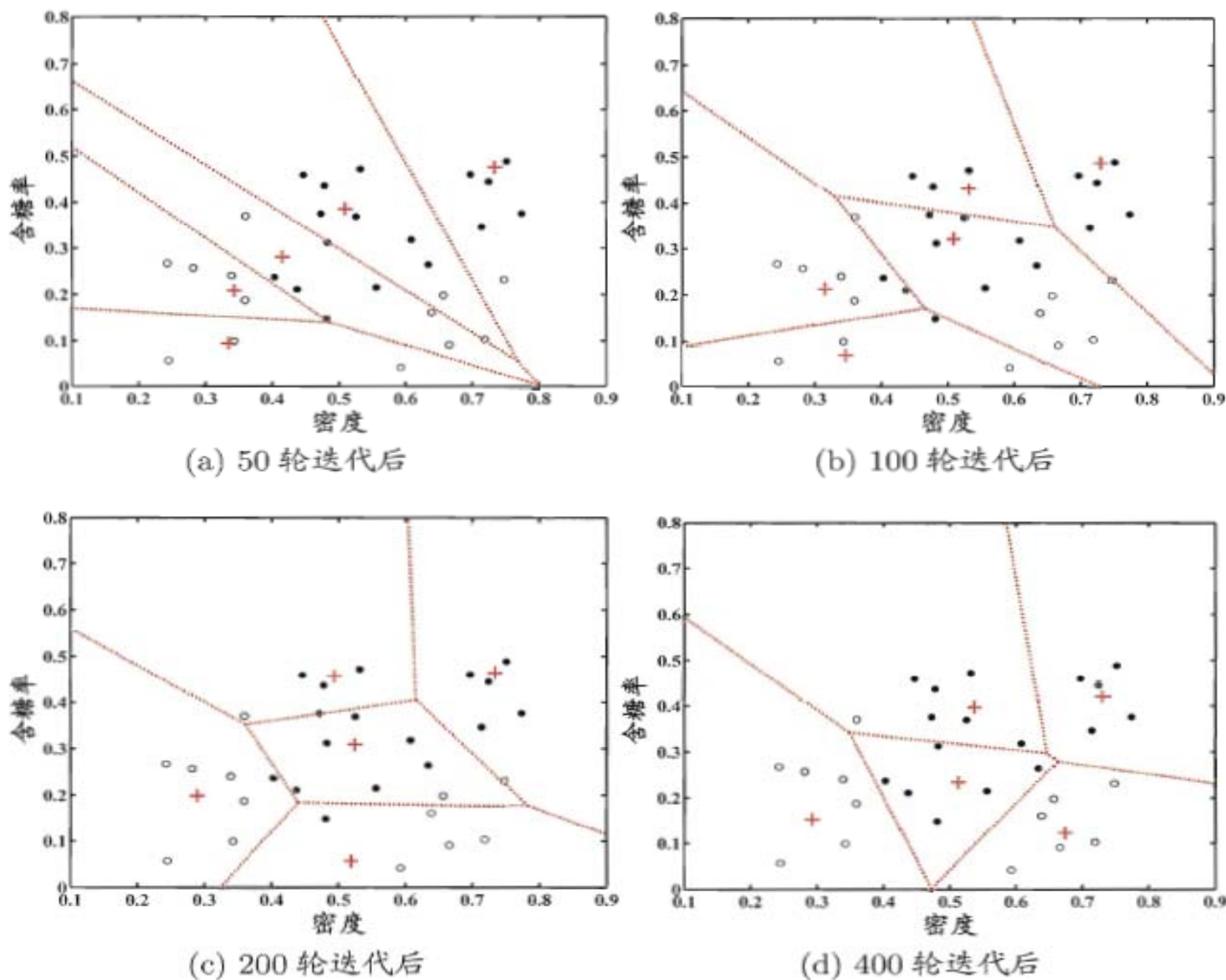
$q$  个原型向量  $\{p_l \mid l = 1, \dots, q\}$

$$\forall x \in X \subset R^d$$

$$R_i = \{x \in X \mid \|x - p_i\|_2 \leq \|x - p_{i'}\|_2, i' \neq i\}$$

$$X = R_1 \cup R_2 \cup \dots \cup R_q$$





**图 9.5** 西瓜数据集 4.0 上 LVQ 算法( $q = 5$ )在不同轮数迭代后的聚类结果.  $c_1$ ,  $c_2$  类样本点与原型向量分别用“●”, “○”与“+”表示, 红色虚线显示出聚类形成的 Voronoi 划分.



## 1. 聚类的引入

## 2. 动态聚类

## 3. 基于模型的聚类

高斯混合聚类(高斯混合模型+贝叶斯决策)

Gaussian Mixture Model, GMM

## 5. 密度聚类

## 6. 层次聚类

## 7. 其它聚类算法

## 问题描述

**输入：** 样本集  $D = \{x_1, \dots, x_m\}$ ,  $x_i = [x_{i1} \ \dots \ x_{id}]^T \in R^d$

要生成的簇的数目  $k$

**输出：** 以  $k$  个**高斯模型的混合**描述的簇  $C = \{C_l \mid l = 1, \dots, k\}$

**混合高斯模型**  $p_M(x) = \sum_{j=1}^k \alpha_i p(x \mid \mu_i, \Sigma_i)$

第  $i$  个**高斯成分**的概率密度函数 (第  $i$  类的条件概率密度函数)

$$p(x \mid \mu_i, \Sigma_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left[ -\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) \right]$$

**高斯成分的混合系数：** 第  $i$  个高斯成分的先验概率  $\alpha_i$   $\sum_{j=1}^k \alpha_j = 1$

**参数集合：**  $\{(\alpha_i, \mu_i, \Sigma_i) \mid i = 1, \dots, k\}$





D中的每个观测样本按照相同的混合概率密度函数独立抽取，样本集D满足**独立、同分布**。

混合概率密度函数的形式已知，但参数未知。

利用样本集D估计已知形式的概率密度函数的未知参数，就叫做**概率密度函数的参数估计问题**。

由于D内每个样本都是无标签的，若采用无类标签的样本数据，进行概率密度函数的参数估计，就是概率密度函数的非监督式参数估计，常采用最大似然估计法。

然后利用估计得到的混合高斯模型，对数据集D划分。



## 高斯混合聚类的问题分解

### STEP1. 高斯混合模型的非监督式参数估计.

基于样本集  $D = \{x_1, \dots, x_m\}$  估计混合高斯模型的参数:

$$\alpha = \{\alpha_1, \dots, \alpha_k\}$$

$$\mu = \{\mu_1, \dots, \mu_k\} \quad \Sigma = \{\Sigma_1, \dots, \Sigma_k\}$$

$$p_M(x; \alpha, \mu, \Sigma) = \sum_{j=1}^k \alpha_j p(x | \mu_j, \Sigma_j)$$

### STEP2. 基于最大后验概率的样本集划分.

$$P_M(z_j = i | x_j) = \frac{P(z_j = i) p(x_j | z_j = i)}{p_M(x_j)} = \frac{\alpha_i p(x_j | \mu_i, \Sigma_i)}{\sum_{l=1}^k \alpha_l p(x | \mu_l, \Sigma_l)}$$

$$\lambda_j = \arg \max_{i \in \{1, 2, \dots, k\}} P_M(z_j = i | x_j) \quad j = 1, \dots, m$$



## 1. 样本集 $D = \{x_1, \dots, x_m\}$ 的随机抽取



假设样本集  $D = \{x_1, \dots, x_m\}$  中的各样本按照独立**同分布**的方式随机抽取得到.

$$p_M(x; \alpha, \mu, \Sigma) = \sum_{j=1}^k \alpha_j p(x | \mu_j, \Sigma_j)$$

抽取过程可以分解为:

**for**  $i = 1, \dots, m$  抽取样本  $x_i \in D$ :

按照  $\alpha_1, \dots, \alpha_k$  定义的先验分布, 随机选择相应的高斯成分, 其中  $\alpha_j$  为第  $j$  个高斯成分出现的先验概率;

按照被选中成分的高斯分布, 进行随机抽取得到观测样本  $x_i$ .

## 2. 似然函数及**对数**似然函数



河北师范大学软件学院  
Software College of Hebei Normal University

(1) 观测样本  $\mathbf{x}_i \in D$  的似然值

$$p_M(\mathbf{x}_i; \alpha, \mu, \Sigma) = \sum_{j=1}^k \alpha_j p(\mathbf{x}_i | \mu_j, \Sigma_j)$$

(2) 样本集  $D$  的似然值 (*likelihood*, 或  $n$  个样本的联合似然值)

$D = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  内各样本独立抽取, 似然值:

$$L(D; \alpha, \mu, \Sigma) = \prod_{i=1}^m p_M(\mathbf{x}_i; \alpha, \mu, \Sigma) = \prod_{i=1}^m \left\{ \sum_{j=1}^k \alpha_j p(\mathbf{x}_i | \mu_j, \Sigma_j) \right\}$$

(3) 样本集  $D$  的**对数**似然值 (*loglikelihood*)

$$\begin{aligned} H(\alpha, \mu, \Sigma) &= LL(D; \alpha, \mu, \Sigma) = \ln L(D; \alpha, \mu, \Sigma) = \ln \left\{ \prod_{i=1}^m p_M(\mathbf{x}_i; \alpha, \mu, \Sigma) \right\} \\ &= \sum_{i=1}^m \ln \left\{ p_M(\mathbf{x}_i; \alpha, \mu, \Sigma) \right\} = \sum_{i=1}^m \ln \left\{ \sum_{j=1}^k \alpha_j p(\mathbf{x}_i | \mu_j, \Sigma_j) \right\} \end{aligned}$$

### 3.高斯混合模型的学习--最大似然估计



河北师范大学软件学院  
Software College of Hebei Normal University

样本集 $D$ 的对数似然值( $loglikelihood$ )

$$LL(D; \alpha, \mu, \Sigma) = \sum_{i=1}^m \ln \{ p_M(\mathbf{x}_i; \alpha, \mu, \Sigma) \} = \sum_{i=1}^m \ln \left\{ \sum_{j=1}^k \alpha_j p(\mathbf{x}_i | \mu_j, \Sigma_j) \right\}$$

因样本集 $D$ 为给定信息，对数似然值的大小取决于参数集合 $\alpha, \mu, \Sigma$ ，因此称 $H(\alpha, \mu, \Sigma)$ 为基于样本集 $D$ 的**对数似然函数**。

**最大似然估计** (*Maximum Likelihood Estimation, MLE*):

$$\begin{aligned} \{\hat{\alpha}, \hat{\mu}, \hat{\Sigma}\} &= \arg \max_{\alpha, \mu, \Sigma} L(D; \alpha, \mu, \Sigma) = \arg \max_{\alpha, \mu, \Sigma} LL(D; \alpha, \mu, \Sigma) \\ &= \arg \max_{\alpha, \mu, \Sigma} \sum_{i=1}^m \ln \left\{ \sum_{j=1}^k \alpha_j p(\mathbf{x}_i | \mu_j, \Sigma_j) \right\} \end{aligned}$$

## 最优化问题：

$$\max_{\alpha, \mu, \Sigma} \sum_{i=1}^m \ln \left\{ \sum_{j=1}^k \alpha_j p(\mathbf{x}_i | \mu_j, \Sigma_j) \right\}$$
$$\text{s.t. } \alpha_i \geq 0 \quad \sum_{j=1}^k \alpha_j = 1$$



河北师范大学软件学院  
Software College of Hebei Normal University



选择何种最优化方法估计上述目标函数的最优值解？  
EM算法(交叉迭代，“E步+M步”)

## 4.最大似然估计的具体实现过程



$$\text{最优化问题} \quad \begin{cases} \max_{\alpha, \mu, \Sigma} \sum_{j=1}^m \ln \left\{ \sum_{i=1}^k \alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i) \right\} \\ \text{s.t. } \alpha_i \geq 0 \quad \sum_{i=1}^k \alpha_i = 1 \end{cases}$$

$$\begin{aligned} LL(D; \alpha, \mu, \Sigma) &= \sum_{j=1}^m \ln \left\{ \sum_{i=1}^k \alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i) \right\} \\ &= \sum_{j=1}^m \ln \left\{ \sum_{i=1}^k \frac{\alpha_i}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \mu_i) \right] \right\} \end{aligned}$$

$$\nabla LL(D; \alpha, \mu, \Sigma) = ?$$

$$\frac{\partial}{\partial \mu_i} LL(D; \alpha, \mu, \Sigma) = ? \quad \frac{\partial}{\partial \Sigma_i} LL(D; \alpha, \mu, \Sigma) = ? \quad \frac{\partial}{\partial \alpha_i} LL(D; \alpha, \mu, \Sigma) = ?$$



$$\frac{\partial}{\partial \mu_i} LL(D; \alpha, \mu, \Sigma) = ?$$

由于  $LL(D; \alpha, \mu, \Sigma) = \sum_{j=1}^m \ln \left\{ \sum_{i=1}^k \alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i) \right\}$

并且  $p(\mathbf{x}_j | \mu_i, \Sigma_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \mu_i) \right]$

$$\frac{\partial}{\partial \mu_i} p(\mathbf{x}_j | \mu_i, \Sigma_i) = p(\mathbf{x}_j | \mu_i, \Sigma_i) \Sigma_i^{-1} (\mathbf{x}_j - \mu_i)$$

所以 
$$\begin{aligned} \frac{\partial}{\partial \mu_i} LL(D; \alpha, \mu, \Sigma) &= \sum_{j=1}^m \frac{\partial}{\partial \mu_i} \ln \left\{ \sum_{i=1}^k \alpha_j p(\mathbf{x}_j | \mu_i, \Sigma_i) \right\} \\ &= \sum_{j=1}^m \frac{\sum_{i=1}^k \alpha_i \frac{\partial}{\partial \mu_i} p(\mathbf{x}_j | \mu_i, \Sigma_i)}{\sum_{i=1}^k \alpha_j p(\mathbf{x}_j | \mu_i, \Sigma_i)} = \sum_{i=1}^m \frac{\alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i) [\Sigma_i^{-1} (\mathbf{x}_j - \mu_i)]}{\sum_{i=1}^k \alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i)} \\ &= \sum_{i=1}^m P_M(z_j = i | \mathbf{x}_j) [\Sigma_i^{-1} (\mathbf{x}_j - \mu_i)] \end{aligned}$$

$$\text{令 } \frac{\partial}{\partial \mu_i} LL(D; \alpha, \mu, \Sigma) = \sum_{j=1}^m P_M(z_j = i | x_j) [\Sigma_i^{-1} (x_j - \mu_i)] = 0$$

$$\text{两边左乘 } \Sigma_i, \text{ 有 } \sum_{j=1}^m P_M(z_j = i | x_j) [(x_j - \mu_i)] = 0$$

$$\text{即 } \sum_{j=1}^m P_M(z_j = i | x_j) x_j = \sum_{j=1}^m P_M(z_j = i | x_j) \mu_i$$

$$\sum_{j=1}^m P_M(z_j = i | x_j) x_j = \left[ \sum_{j=1}^m P_M(z_j = i | x_j) \right] \mu_i$$

所以

$$\hat{\mu}_i = \frac{\sum_{j=1}^m P_M(z_j = i | x_j) x_j}{\sum_{j=1}^m P_M(z_j = i | x_j)} = \frac{\sum_{j=1}^m \gamma_{ji} x_j}{\sum_{j=1}^m \gamma_{ji}} \quad i = 1, \dots, k$$

西瓜书209页式(9.34)



河北师范大学软件学院  
Software College of Hebei Normal University





$$\frac{\partial}{\partial \Sigma_i} LL(D; \alpha, \mu, \Sigma) = ?$$

由于  $LL(D; \alpha, \mu, \Sigma) = \sum_{j=1}^m \ln \left\{ \sum_{i=1}^k \alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i) \right\}$

并且  $p(\mathbf{x}_j | \mu_i, \Sigma_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \mu_i) \right]$

$$\frac{\partial}{\partial \Sigma_i} (|\Sigma_i|) = |\Sigma_i| (\Sigma_i^{-1})^T = |\Sigma_i| \Sigma_i^{-T} = |\Sigma_i| \Sigma_i^{-1}$$

$$\frac{\partial}{\partial \Sigma_i} (a^T \Sigma_i^{-1} b) = \frac{\partial}{\partial \Sigma_i} \text{tr}(a^T \Sigma_i^{-1} b) = -\Sigma_i^{-1} b a^T \Sigma_i^{-1}$$

$$\frac{\partial}{\partial \Sigma_i} ((\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \mu_i)) = -\Sigma_i^{-1} (\mathbf{x}_j - \mu_i) (\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1}$$





$$\frac{\partial}{\partial \Sigma_i} p(\mathbf{x}_j | \mu_i, \Sigma_i)$$

$$= \frac{-\frac{1}{2} |\Sigma_i|^{-3/2} |\Sigma_i| \Sigma_i^{-1}}{(2\pi)^{d/2}} \exp \left[ -\frac{1}{2} (\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \mu_i) \right] \\ + \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \mu_i) \right] \left[ \frac{1}{2} \Sigma_i^{-1} (\mathbf{x}_j - \mu_i) (\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} \right]$$

$$= -\frac{1}{2} p(\mathbf{x}_j | \mu_i, \Sigma_i) \Sigma_i^{-1} + \frac{1}{2} p(\mathbf{x}_j | \mu_i, \Sigma_i) \left[ \Sigma_i^{-1} (\mathbf{x}_j - \mu_i) (\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} \right]$$

$$= \frac{1}{2} p(\mathbf{x}_j | \mu_i, \Sigma_i) \Sigma_i^{-1} \left[ (\mathbf{x}_j - \mu_i) (\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} - I \right]$$

由于  $LL(D; \alpha, \mu, \Sigma) = \sum_{j=1}^m \ln \left\{ \sum_{i=1}^k \alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i) \right\}$



所以  $\frac{\partial}{\partial \Sigma_i} LL(D; \alpha, \mu, \Sigma) = \sum_{j=1}^m \frac{\partial}{\partial \Sigma_i} \ln \left\{ \sum_{i=1}^k \alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i) \right\} = \sum_{j=1}^m \frac{\alpha_i \frac{\partial}{\partial \Sigma_i} p(\mathbf{x}_j | \mu_i, \Sigma_i)}{\sum_{i=1}^k \alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i)}$

$$= \sum_{j=1}^m \frac{\alpha_i \frac{1}{2} p(\mathbf{x}_j | \mu_i, \Sigma_i) \Sigma_i^{-1} \left[ (\mathbf{x}_j - \mu_i)(\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} - \mathbf{I} \right]}{\sum_{i=1}^k \alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i)}$$

$$= \frac{1}{2} \sum_{j=1}^m \frac{\alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i) \Sigma_i^{-1} \left[ (\mathbf{x}_j - \mu_i)(\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} - \mathbf{I} \right]}{\sum_{i=1}^k \alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i)}$$

$$= \frac{1}{2} \sum_{j=1}^m P_M(z_j = i | \mathbf{x}_j) \Sigma_i^{-1} \left[ (\mathbf{x}_j - \mu_i)(\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} - \mathbf{I} \right]$$

即:  $\frac{\partial}{\partial \Sigma_i} LL(D; \alpha, \mu, \Sigma) = \frac{1}{2} \sum_{j=1}^m P_M(z_j = i | \mathbf{x}_j) \Sigma_i^{-1} \left[ (\mathbf{x}_j - \mu_i)(\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} - \mathbf{I} \right]$



$$\frac{\partial}{\partial \Sigma_i} LL(D; \alpha, \mu, \Sigma) = \frac{1}{2} \sum_{j=1}^m P_M(z_j = i | x_j) \Sigma_i^{-1} \left[ (\mathbf{x}_j - \mu_i)(\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} - I \right]$$

$$\text{令 } \frac{\partial}{\partial \Sigma_i} LL(D; \alpha, \mu, \Sigma) = 0$$

$$\text{则 } \sum_{j=1}^m P_M(z_j = i | x_j) \Sigma_i^{-1} (\mathbf{x}_j - \mu_i)(\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} = \sum_{j=1}^m P_M(z_j = i | x_j) \Sigma_i^{-1}$$

$$\text{即 } \Sigma_i^{-1} \left[ \sum_{j=1}^m P_M(z_j = i | x_j) (\mathbf{x}_j - \mu_i)(\mathbf{x}_j - \mu_i)^T \right] \Sigma_i^{-1} = \left[ \sum_{j=1}^m P_M(z_j = i | x_j) \right] \Sigma_i^{-1}$$

等号两侧同时左乘、右乘  $\Sigma_i$

$$\Sigma_i \Sigma_i^{-1} \left[ \sum_{j=1}^m P_M(z_j = i | x_j) (\mathbf{x}_j - \mu_i)(\mathbf{x}_j - \mu_i)^T \right] \Sigma_i^{-1} \Sigma_i = \Sigma_i \left[ \sum_{j=1}^m P_M(z_j = i | x_j) \right] \Sigma_i^{-1} \Sigma_i$$

即：

$$\sum_{j=1}^m P_M(z_j = i | x_j) (\mathbf{x}_j - \mu_i)(\mathbf{x}_j - \mu_i)^T = \Sigma_i \sum_{j=1}^m P_M(z_j = i | x_j)$$



$$\sum_{j=1}^m P_M(z_j = i | x_j) (\mathbf{x}_j - \mu_i) (\mathbf{x}_j - \mu_i)^T = \sum_i \sum_{j=1}^m P_M(z_j = i | x_j)$$

所以 
$$\hat{\Sigma}_i = \frac{\sum_{j=1}^m P_M(z_j = i | x_j) (\mathbf{x}_j - \mu_i) (\mathbf{x}_j - \mu_i)^T}{\sum_{j=1}^m P_M(z_j = i | x_j)}$$

即:

$$\hat{\Sigma}_i = \frac{\sum_{j=1}^m \gamma_{ji} (\mathbf{x}_j - \hat{\mu}_i) (\mathbf{x}_j - \hat{\mu}_i)^T}{\sum_{j=1}^m \gamma_{ji}} \quad i = 1, \dots, k$$

西瓜书209页式(9.35)



$$\frac{\partial}{\partial \alpha_i} LL(D; \alpha, \mu, \Sigma) = ?$$

固定  $\mu, \Sigma$ , 估计混合系数  $\alpha$ , 有:

$$\begin{cases} \max_{\alpha} \sum_{j=1}^m \ln \left\{ \sum_{i=1}^k \alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i) \right\} \\ \text{s.t. } \alpha_i \geq 0 \quad \sum_{i=1}^k \alpha_i = 1 \end{cases}$$

构造 *Lagrange* 辅助函数

$$H(\alpha, \lambda) = \sum_{j=1}^m \ln \left\{ \sum_{i=1}^k \alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i) \right\} + \lambda \left[ \sum_{i=1}^k \alpha_i - 1 \right]$$

$$\text{令 } \frac{H(\alpha, \lambda)}{\partial \alpha_i} = \sum_{j=1}^m \frac{p(\mathbf{x}_j | \mu_i, \Sigma_i)}{\sum_{i=1}^k \alpha_i p(\mathbf{x}_j | \mu_i, \Sigma_i)} + \lambda = 0 \quad i = 1, \dots, k$$

由贝叶斯公式知

$$\sum_{j=1}^m \frac{P_M(z_j = i | \mathbf{x}_j)}{\alpha_i} + \lambda = 0 \quad i = 1, \dots, k$$



$$\sum_{j=1}^m \frac{P_M(z_j = i | x_j)}{\alpha_i} + \lambda = 0 \quad i = 1, \dots, k$$

则有  $\sum_{j=1}^m P_M(z_j = i | x_j) = -\alpha_i \lambda \quad i = 1, \dots, k$

从而有  $\sum_{i=1}^k \left[ \sum_{j=1}^m P_M(z_j = i | x_j) \right] = \sum_{i=1}^k [-\alpha_i \lambda]$

$$\sum_{j=1}^m \sum_{i=1}^k P_M(z_j = i | x_j) = -\lambda \sum_{i=1}^k \alpha_i$$

由于  $\sum_{i=1}^k P_M(z_j = i | x_j) = 1, \sum_{i=1}^k \alpha_i = 1$  所以  $m = -\lambda$



$$\sum_{j=1}^m P_M(z_j = i | x_j) = -\alpha_i \lambda \quad i = 1, \dots, k$$

$$m = -\lambda$$

$$\sum_{j=1}^m P_M(z_j = i | x_j) = \alpha_i m \quad i = 1, \dots, k$$

$$\hat{\alpha}_i = \frac{1}{m} \sum_{j=1}^m P_M(z_j = i | x_j) = \frac{1}{m} \sum_{j=1}^m \gamma_{ji}$$

即

$$i = 1, \dots, k$$

西瓜书209页式(9.38)

# GMM模型参数估计的交叉迭代



河北师范大学软件学院  
Software College of Hebei Normal University

(1) 初始化  $\hat{\mu}_i(0), \hat{\Sigma}_i(0), \hat{\alpha}_i(0), i = 1, \dots, k$

(2) 交叉迭代

$$E-STEP \quad \gamma_{ji}(n+1) = \frac{\alpha_i(n) p(x_j | \mu_i(n), \Sigma_i(n))}{\sum_{l=1}^k \alpha_l(n) p(x_j | \mu_l(n), \Sigma_l(n))} \quad \begin{cases} i = 1, \dots, k \\ j = 1, \dots, m \end{cases}$$

$$\gamma_{ji} = P_M(z_j = i | x_j)$$

$$M-STEP \quad \left\{ \begin{aligned} \hat{\mu}_i(n+1) &= \frac{\sum_{j=1}^m \gamma_{ji}(n+1) x_j}{\sum_{j=1}^m \gamma_{ji}(n+1)} & i = 1, \dots, k \\ \hat{\Sigma}_i(n+1) &= \frac{\sum_{j=1}^m \gamma_{ji}(n+1) (x_j - \hat{\mu}_i(n+1)) (x_j - \hat{\mu}_i(n+1))^T}{\sum_{j=1}^m \gamma_{ji}(n+1)} & i = 1, \dots, k \\ \hat{\alpha}_i(n+1) &= \frac{1}{m} \sum_{j=1}^m \gamma_{ji}(n+1) & i = 1, \dots, k \end{aligned} \right.$$



输入: 样本集  $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ ;  
高斯混合成分个数  $k$ .

## 如何初始化GMM模型的参数?

过程:

1: 初始化高斯混合分布的模型参数  $\{(\alpha_i, \mu_i, \Sigma_i) \mid 1 \leq i \leq k\}$

2: repeat

3: for  $j = 1, 2, \dots, m$  do

4: 根据式(9.30)计算  $\mathbf{x}_j$  由各混合成分生成的后验概率, 即  
 $\gamma_{ji} = p_{\mathcal{M}}(z_j = i \mid \mathbf{x}_j) \quad (1 \leq i \leq k)$

5: end for

6: for  $i = 1, 2, \dots, k$  do

7: 计算新均值向量:  $\mu'_i = \frac{\sum_{j=1}^m \gamma_{ji} \mathbf{x}_j}{\sum_{j=1}^m \gamma_{ji}};$

8: 计算新协方差矩阵:  $\Sigma'_i = \frac{\sum_{j=1}^m \gamma_{ji} (\mathbf{x}_j - \mu'_i)(\mathbf{x}_j - \mu'_i)^T}{\sum_{j=1}^m \gamma_{ji}};$

9: 计算新混合系数:  $\alpha'_i = \frac{\sum_{j=1}^m \gamma_{ji}}{m};$

10: end for

11: 将模型参数  $\{(\alpha_i, \mu_i, \Sigma_i) \mid 1 \leq i \leq k\}$  更新为  $\{(\alpha'_i, \mu'_i, \Sigma'_i) \mid 1 \leq i \leq k\}$

12: until 满足停止条件

13:  $C_i = \emptyset \quad (1 \leq i \leq k)$

14: for  $j = 1, 2, \dots, m$  do

15: 根据式(9.31)确定  $\mathbf{x}_j$  的簇标记  $\lambda_j$ ;

16: 将  $\mathbf{x}_j$  划入相应的簇:  $C_{\lambda_j} = C_{\lambda_j} \cup \{\mathbf{x}_j\}$

17: end for

输出: 簇划分  $C = \{C_1, C_2, \dots, C_k\}$

PART1. 基于 EM  
算法的 GMM 模  
型的学习

PART2. 根据  
GMM 模型生成  
最终划分结果.



1. 聚类的引入

2. 动态聚类

3. 基于模型的聚类      高斯混合聚类

4. 密度聚类 (density-based clustering)

DBSCAN

*Density - Based Spatial Clustering of Applications with Noise*

6. 层次聚类

7. 其它聚类算法



- DBSCAN 算法是一种基于高密度连通区域的、基于密度的聚类算法
- 能够将具有足够高密度的区域划分为簇
- 能够在具有噪声的数据中发现任意形状的簇
- 能够发现异常点

# 1. 有关概念 [给定样本集 $D = \{x_1, \dots, x_m\}$ ]



河北师范大学软件学院  
Software College of Hebei Normal University

## [1] 两个全局邻域参数 ( $\varepsilon, MinPts$ )

$\varepsilon$ --邻域最大半径

$MinPts$ --给定样本的  $\varepsilon$ -邻域内最小样本数.

## [2] $\varepsilon$ -邻域

对于  $\forall x_j \in D$ ,  $x_j$  的  $\varepsilon$ -邻域为  $N_\varepsilon(x_j) = \{x_i \in D \mid dist(x_i, x_j) \leq \varepsilon\}$

## [3] 核心对象 (core object)

若  $|N_\varepsilon(x_j)| \geq MinPts$ , 则称  $x_j$  为一个核心对象.

## [4] 密度直达 (directly density-reachable)

若  $x_j \in N_\varepsilon(x_i)$ , 并且  $x_i$  为一个核心对象, 则称  $x_j$  为由  $x_i$  密度直达.

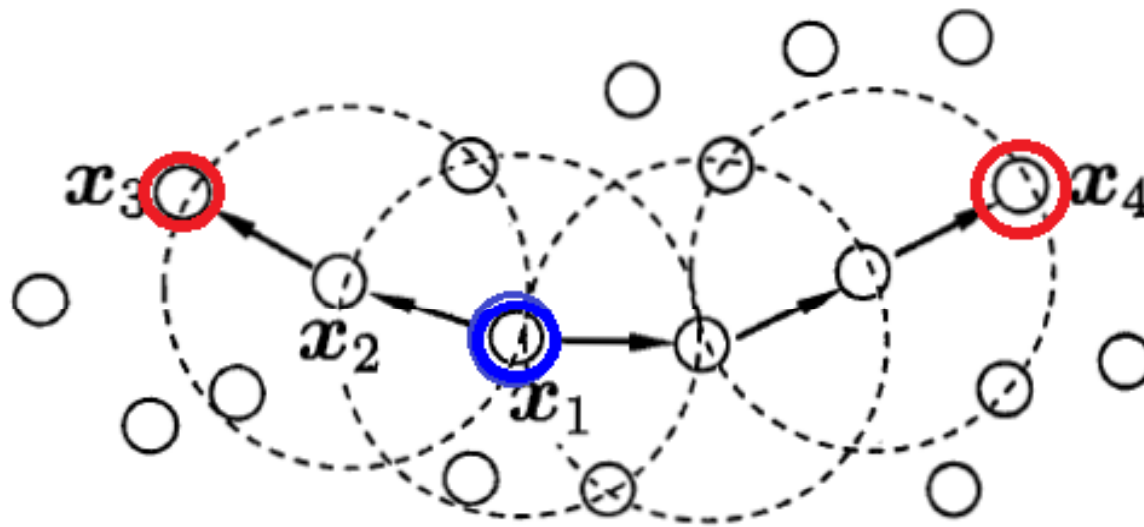
## [5] 密度可达 (density-reachable)

对于  $x_i, x_j$ , 若存在样本序列  $p_1, p_2, \dots, p_n$ , 其中  $p_1 = x_i, p_n = x_j$ , 且  $p_{i+1}$  由  $p_i$  密度直达, 则称  $x_j$  由  $x_i$  密度可达.

[6] **密度相连** (*density-connected*)

对于 $x_i$ 与 $x_j$ , 若存在样本 $x_k$ , 使得 $x_i$ 与 $x_j$ 均由 $x_k$  **密度可达**, 则称 $x_i$ 与 $x_j$  **密度相连**.

邻域参数( $\varepsilon, MinPts$ ),  $\varepsilon$ -邻域, **核心对象**, **密度直达**, **密度可达**, **密度相连**



DBSCAN 定义的基本概念 ( $MinPts = 3$ ): 虚线显示出  $\varepsilon$ -邻域,  $x_1$  是核心对象,  $x_2$  由  $x_1$  密度直达,  $x_3$  由  $x_1$  密度可达,  $x_3$  与  $x_4$  密度相连.

### [7] 边界对象 (*border object*)

若  $N_\varepsilon(\mathbf{x}_i) \geq MinPts$ ,  $\mathbf{x}_j \in N_\varepsilon(\mathbf{x}_i)$ , 并且  $N_\varepsilon(\mathbf{x}_j) < MinPts$ ,

则称  $\mathbf{x}_j$  为一个边界对象.

边界对象位于聚类簇的边界处.

### [8] 噪声对象 (*noise object*)

核心对象、边界对象以外的其它样本.

### [9] 簇

由 **密度可达** 关系导出的 **最大密度相连样本集**。

给定邻域参数  $(\varepsilon, MinPts)$ , 簇  $C \subseteq D$  为满足以下性质的非空样本子集:

**连接性** (*connectivity*):  $\mathbf{x}_i, \mathbf{x}_j \in C \Rightarrow \mathbf{x}_i$  与  $\mathbf{x}_j$  **密度相连**.

**最大性** (*maximality*):  $\mathbf{x}_i \in C$ , 并且  $\mathbf{x}_j$  由  $\mathbf{x}_i$  **密度可达**  $\Rightarrow \mathbf{x}_j \in C$ .

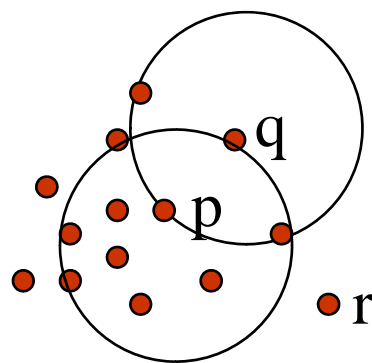
## 样本集 $D$ 的主要组成

**核心对象** (聚类簇的种子点)  
**边界对象** (可由核心对象密度可达的非核心对象)  
**噪声对象** (异常点, 孤立点)

} 聚类簇

高密度区

低密度区



MinPts = 5

Eps = 1 cm

## 2. DBSCAN算法纲要



### 特点

- 无需指定聚类簇的数目;
- 只需输入两个全局参数(**全局参数取值如何选?**).

### 启发

- 识别样本集内的核心对象;
- 任选一个核心对象, 作为一个聚类簇的种子点;
- 获取该种子点及其所有密度可达样本, 构成一个聚类簇.





- 遍历整个样本集 $D$ ，确定核心对象、边界对象、及噪声对象
- 将所有噪声对象标记为异常。
- 考查每个核心对象：核心对象彼此可密度直达，则应划分至相同的聚类簇。
- 所有边界对象：若它们是密度相连的，应划分至与其核心对象一致的聚类簇内。

## 2. 算法描述

STEP1. 识别给定样本集  $D$  的所有核心对象。

得到核心对象集合  $\Omega$

STEP2. 初始化聚类簇的数目为0；初始化未被访问的样本集为整个数据集  $D$ 。

STEP3. 重复如下过程，生成一系列聚类簇，直到核心对象集合为空。

➤ 从核心对象集合中，任选1核心对象，作为聚类簇的一个种子点，找出其密度可达的所有样本，构成1个聚类簇。

➤ 更新核心对象集合；

➤ 更新未访问的样本集合

STEP4. 输出所有聚类簇。

输入：样本集  $D = \{x_1, x_2, \dots, x_m\}$ ；  
邻域参数  $(\epsilon, MinPts)$ 。

过程：

```
1: 初始化核心对象集合:  $\Omega = \emptyset$ 
2: for  $j = 1, 2, \dots, m$  do
3:   确定样本  $x_j$  的  $\epsilon$ -邻域  $N_\epsilon(x_j)$ ;
4:   if  $|N_\epsilon(x_j)| \geq MinPts$  then
5:     将样本  $x_j$  加入核心对象集合:  $\Omega = \Omega \cup \{x_j\}$ 
6:   end if
7: end for

8: 初始化聚类簇数:  $k = 0$ 
9: 初始化未访问样本集合:  $\Gamma = D$ 

10: while  $\Omega \neq \emptyset$  do
11:   记录当前未访问样本集合:  $\Gamma_{old} = \Gamma$ ;
12:   随机选取一个核心对象  $o \in \Omega$ , 初始化队列  $Q = \langle o \rangle$ ;
13:    $\Gamma = \Gamma \setminus \{o\}$ ; 更新未访问的样本集合
14:   while  $Q \neq \emptyset$  do
15:     取出队列  $Q$  中的首个样本  $q$ ;
16:     if  $|N_\epsilon(q)| \geq MinPts$  then
17:       令  $\Delta = N_\epsilon(q) \cap \Gamma$ ;
18:       将  $\Delta$  中的样本加入队列  $Q$ ;
19:        $\Gamma = \Gamma \setminus \Delta$ ; 更新未访问的样本集合
20:     end if
21:   end while
22:    $k = k + 1$ , 生成聚类簇  $C_k = \Gamma_{old} \setminus \Gamma$ ;
23:    $\Omega = \Omega \setminus C_k$  更新核心对象集合
24: end while
```

输出：簇划分  $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$



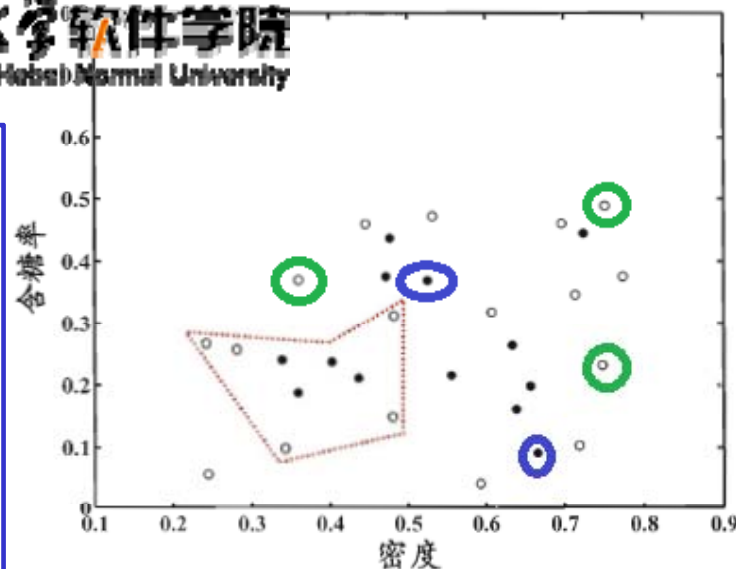
## DBSCAN算法的 聚类结果

核心对象  
非核心对象

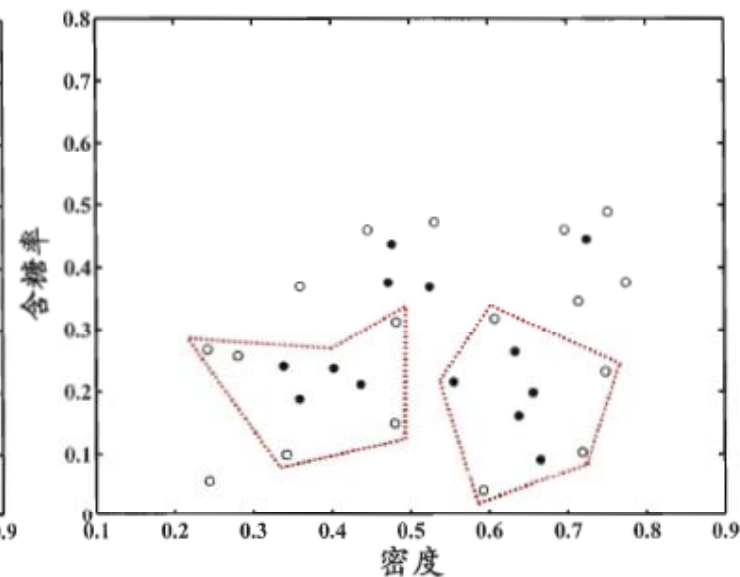
密度直达  
密度可达  
密度相连

聚类簇

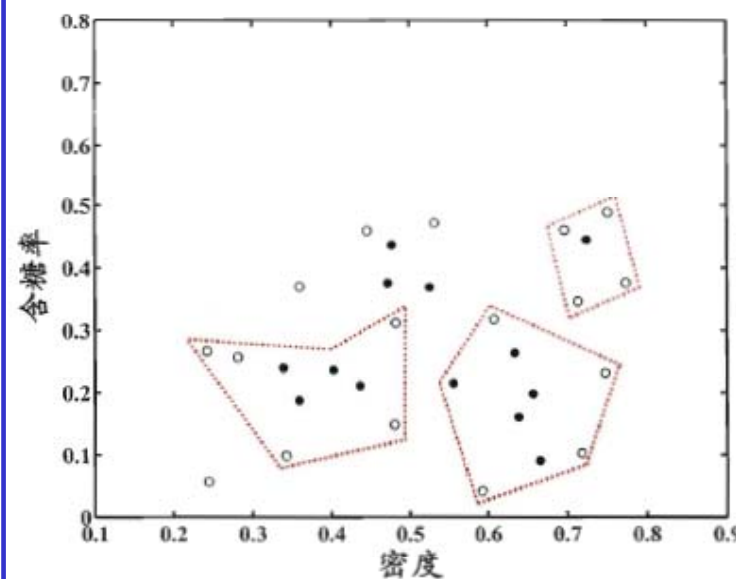
噪声或异常点



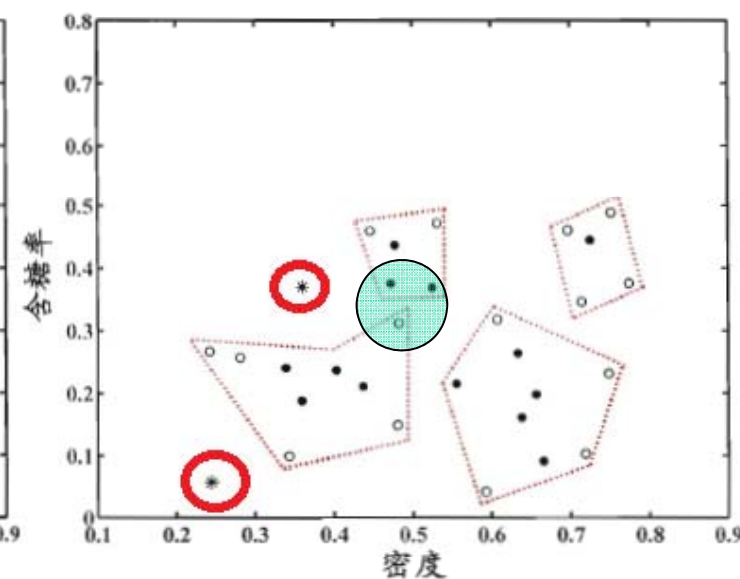
(a) 生成聚类簇  $C_1$



(b) 生成聚类簇  $C_2$

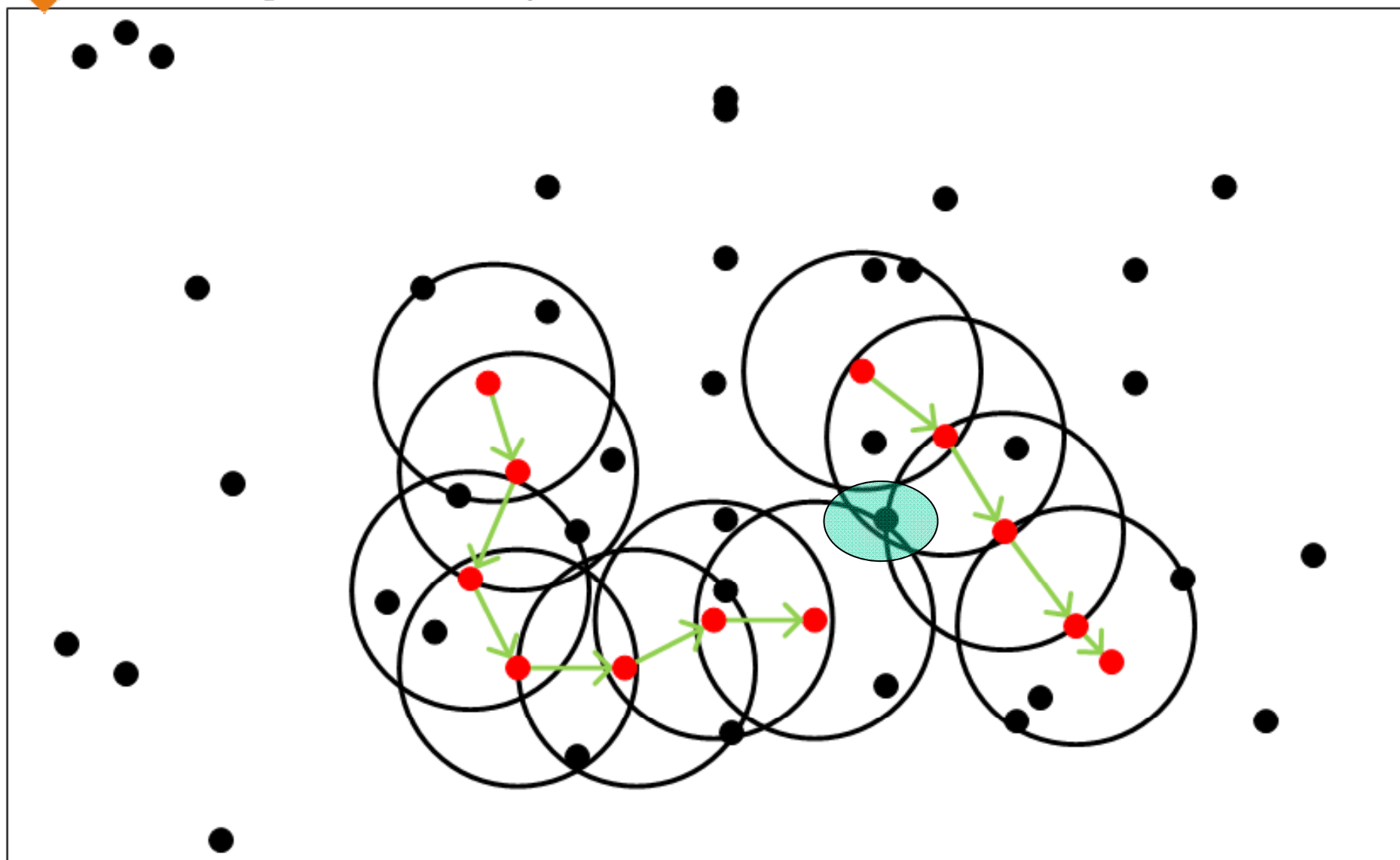


(c) 生成聚类簇  $C_3$



(d) 生成聚类簇  $C_4$







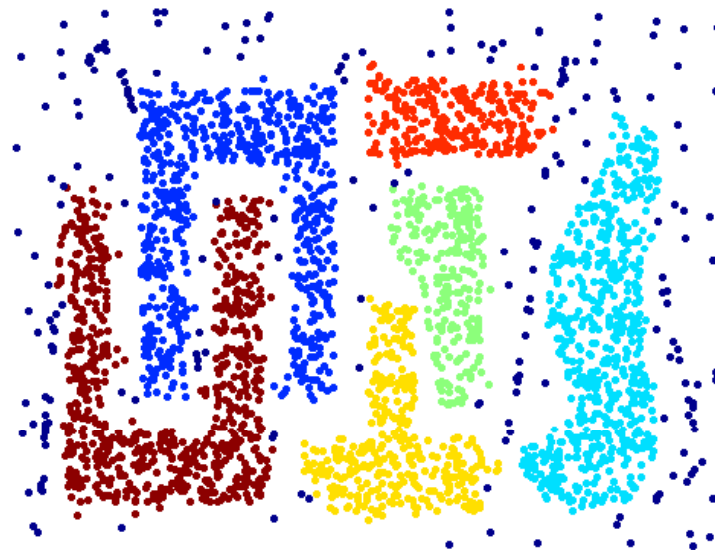
河北师范大学软件学院  
Software College of Hebei Normal University

# DBSCAN 聚类算法的细节

# DBSCAN运行效果好的时候



Original Points



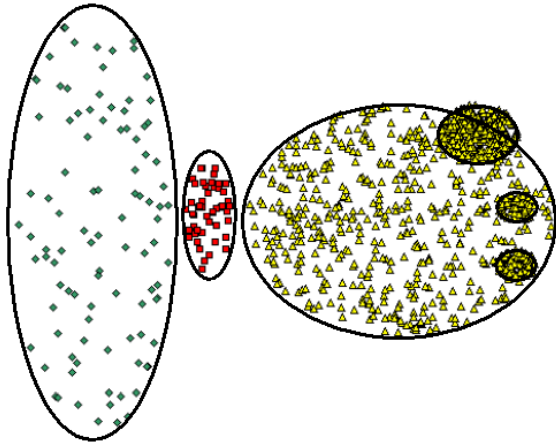
Clusters

- 对噪声不敏感
- 可处理不同形状和大小的数据



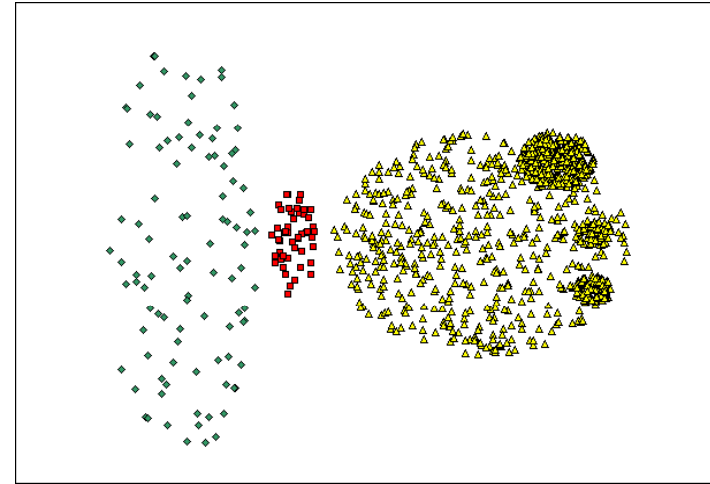


# DBSCAN运行不好的效果

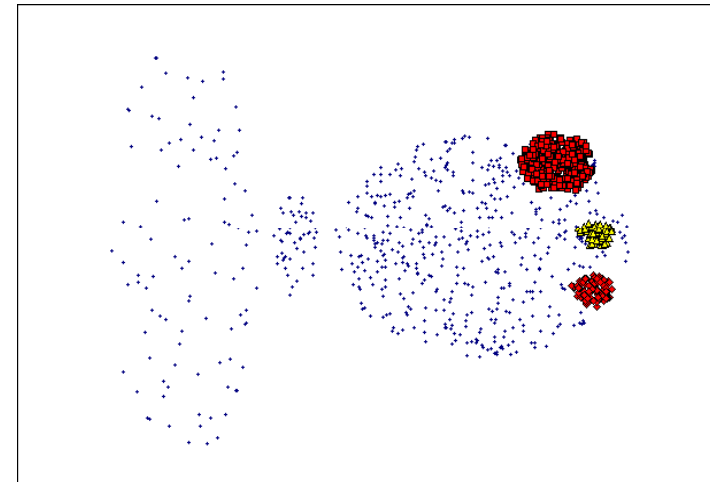


Original Points

- 密度变化的数据
- 高维数据



(MinPts=4, Eps=9.75).



(MinPts=4, Eps=9.92)



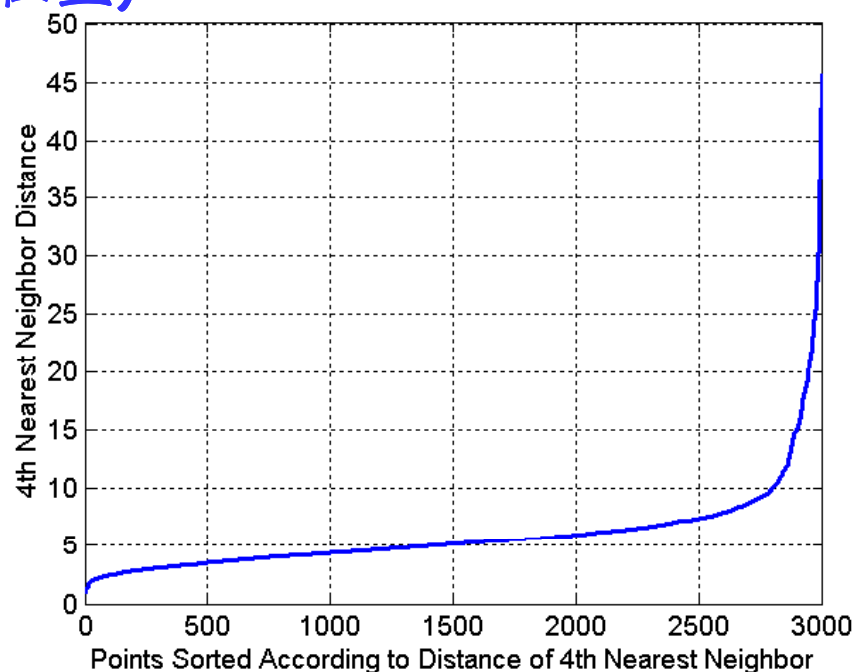


## 如何适当选取EPS和MinPts

### --基于k-距离

- 位于同一聚类簇的所有样本点,它们到其第k个最近邻的距离应大致一样
- 噪声点到其第k个最近邻的距离比较远
- 获取每个样本到其第k个最近邻的距离,从小到大升序排列得到k-距离变化曲线图
- 找到曲线的拐点(变化剧烈的位置)
- 然后:  
Eps即为变化剧烈位置对应的K-距离.

MinPts取k





# DBSCAN算法的优缺点

- 优点

基于密度定义，相对抗噪音；  
能处理任意形状和大小的簇

- 缺点

密度分布不均匀、或密度变化较大的簇时，会有麻烦；  
邻域半径小，数量多的小簇，核心点数量减小  
邻域半径大，本不属于同一簇的样本会聚至相同簇

对于高维问题，密度定义是个比较麻烦的问题



1. 聚类的引入

2. 动态聚类

3. 基于模型的聚类      高斯混合聚类

4. 密度聚类 (density-based clustering) DBSCAN

5. 层次聚类

hierarchical clustering

也称系统聚类、分级聚类

6. 其它聚类算法

# 1. 层次聚类的引入

## (1) 问题描述

数据集  $D = \{x_1, \dots, x_m\} \Rightarrow$  划分为合理的  $k$  个聚类簇。

$$1 \leq k \leq m$$

$k$  值的极端情况  $\begin{cases} \text{最多 } m \text{ 簇, 每簇包含一个样本} \\ \text{最少 1 簇, 所有样本同属一簇} \end{cases}$

## (2) 层次聚类的实现方式

样本数据的递归聚类：聚类过程中逐级考察簇间的距离，  
以此决定类别数。



### 合并式(聚合式)聚类(*agglomerative clustering*)

从聚类簇数目最多开始，**自底向上**，逐级合并距离最近的两个聚类簇，聚类簇的数目逐渐减少，直至与预设的聚类簇数目一致。经常使用。

如：**AGNES**算法(*AGglomerative NESting*)

### 分裂式聚类(*divisive clustering*)

从聚类簇数目最少开始，**自顶向下**，逐级分裂每级中最松散的聚类簇，聚类簇的数目逐级增加，直至与预设聚类簇数目一致。

如：**TSVQ**，可用于码本(*codebook*)的生成

### (3) 层次聚类的几个关键问题



河北师范大学软件学院  
Software College of Hebei Normal University

#### 相似性(相异性)度量

聚类簇之间;  
簇内样本之间

#### 聚类簇合并 / 分裂停止的条件

距离阈值;  
预设的聚类簇数目

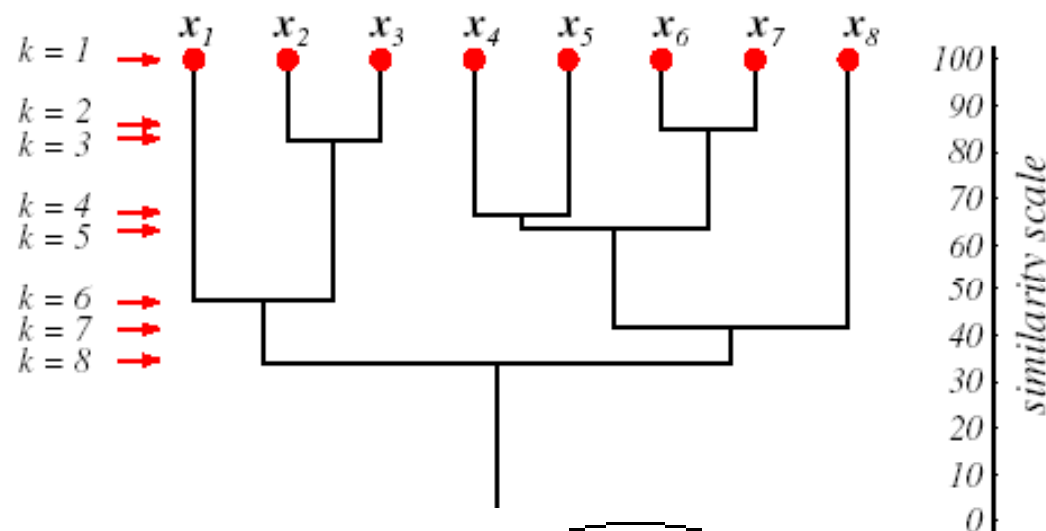
#### 计算复杂度

#### (4) 聚合式系统聚类结果的表示

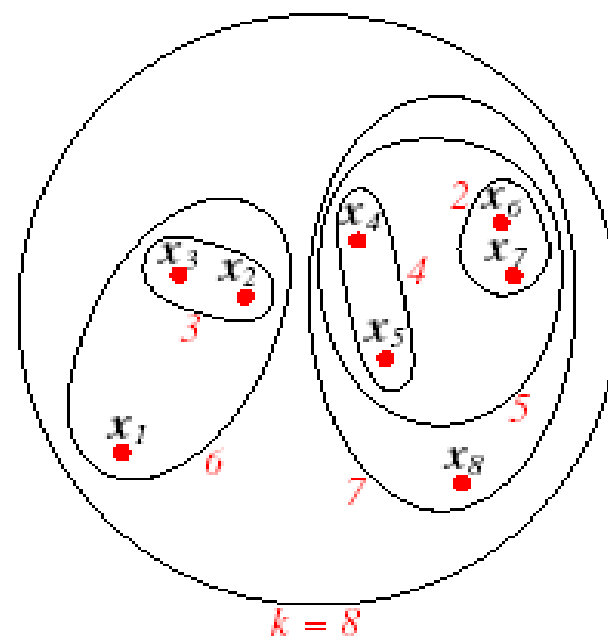


河北师范大学软件学院  
Software College of Hebei Normal University

##### A. 树状图(*dendrogram*), 聚类树



可定量表示聚类簇间的相似性(或相异性), 簇间相似性标尺随着层数增加, 逐渐减小。



##### B. 维恩图(*venn diagram*), 集合图

只能定性表示类间相似性

## 2. 聚合式系统聚类算法



给定样本集  $D = \{x_1, \dots, x_m\}$ , 选定簇间距离度量  $d(C_i, C_j)$ , 将样本  $D$  聚成  $k$  簇.

不同的距离度量, 可产生不同的聚类结果

### STEP1 初始化.

A. 指定最终类别数  $k$ ;

B. 初始划分  $\begin{cases} \text{每个样本自成1簇} & C_i \leftarrow \{x_i\} \quad i = 1, \dots, m \\ \text{初始类别数} & q \leftarrow m \end{cases}$

STEP2 合并. 重复以下工作, 直到  $q = k$ .

A. 寻找最相似 (距离最小) 的两簇  $C_i, C_j (i < j)$

$$d(C_i, C_j) = \min_{\substack{\forall l, m \in \{1, \dots, q\} \\ \text{并且 } l < m}} d(C_l, C_m)$$

B. 记录最小距离, 合并两类  $C_i, C_j$ :  $q \leftarrow q - 1$

STEP3 输出  $k$  个聚类簇.



# AGNES算法描述

STEP1.初始化聚类簇

STEP2.初始化距离矩阵

STEP3.初始化聚类簇的数目

STEP4.逐级合并最相似的两簇，直到满足规定的聚类簇数目：

- (1) 合并最相似的两簇；
- (2) 更新相关簇的序号；
- (3) 更新相关距离矩阵；
- (4) 更新簇的数目

STEP5.输出结果。

输入：样本集  $D = \{x_1, x_2, \dots, x_m\}$ ;  
聚类簇距离度量函数  $d$ ;  
聚类簇数  $k$ .

过程：

```
1: for  $j = 1, 2, \dots, m$  do
2:    $C_j = \{x_j\}$ 
3: end for
4: for  $i = 1, 2, \dots, m$  do
5:   for  $j = 1, 2, \dots, m$  do
6:      $M(i, j) = d(C_i, C_j)$ ;
7:      $M(j, i) = M(i, j)$ 
8:   end for
9: end for
10: 设置当前聚类簇个数:  $q = m$ 
11: while  $q > k$  do
12:   找出距离最近的两个聚类簇  $C_{i^*}$  和  $C_{j^*}$ ;
13:   合并  $C_{i^*}$  和  $C_{j^*}$ :  $C_{i^*} = C_{i^*} \cup C_{j^*}$ ;
14:   for  $j = j^* + 1, j^* + 2, \dots, q$  do
15:     将聚类簇  $C_j$  重编号为  $C_{j-1}$ 
16:   end for
17:   删除距离矩阵  $M$  的第  $j^*$  行与第  $j^*$  列;
18:   for  $j = 1, 2, \dots, q - 1$  do
19:      $M(i^*, j) = d(C_{i^*}, C_j)$ ;
20:      $M(j, i^*) = M(i^*, j)$ 
21:   end for
22:    $q = q - 1$ 
23: end while
```

输出：簇划分  $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$

### 3. 聚类簇 $C_i, C_j$ 之间的连接 (*linkage*)

#### (1) 最小距离

$$d_{\min}(C_i, C_j) = \min_{\substack{x \in C_i \\ z \in C_j}} \text{dist}(x, z)$$

相应聚类算法就是“单链接”算法 (*single linkage* 或 *single-link*)

采用该距离度量 **聚类簇  $C_i, C_j$  之间** 的相异程度，进行聚类，就是产生 **最小生成树** (*minimal spanning tree*)

#### 缺陷

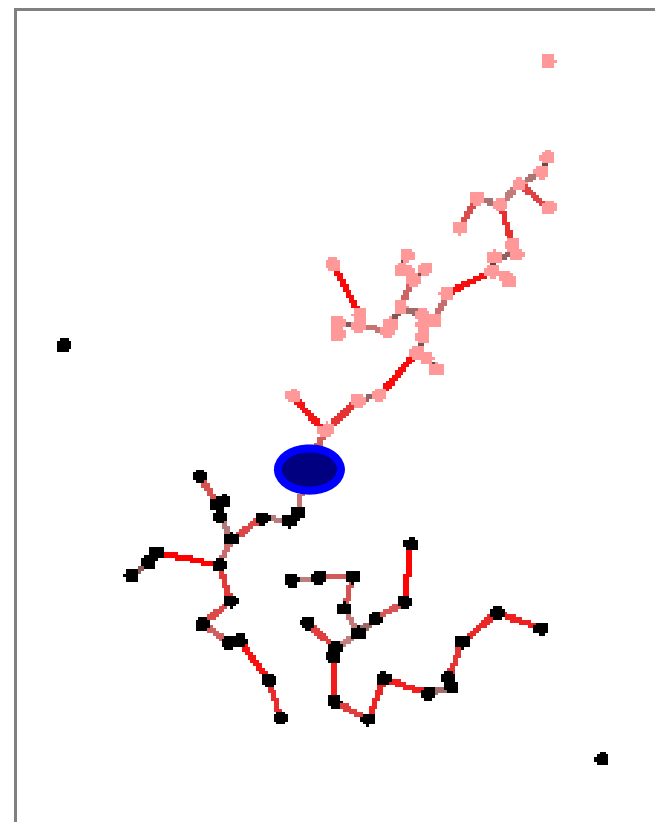
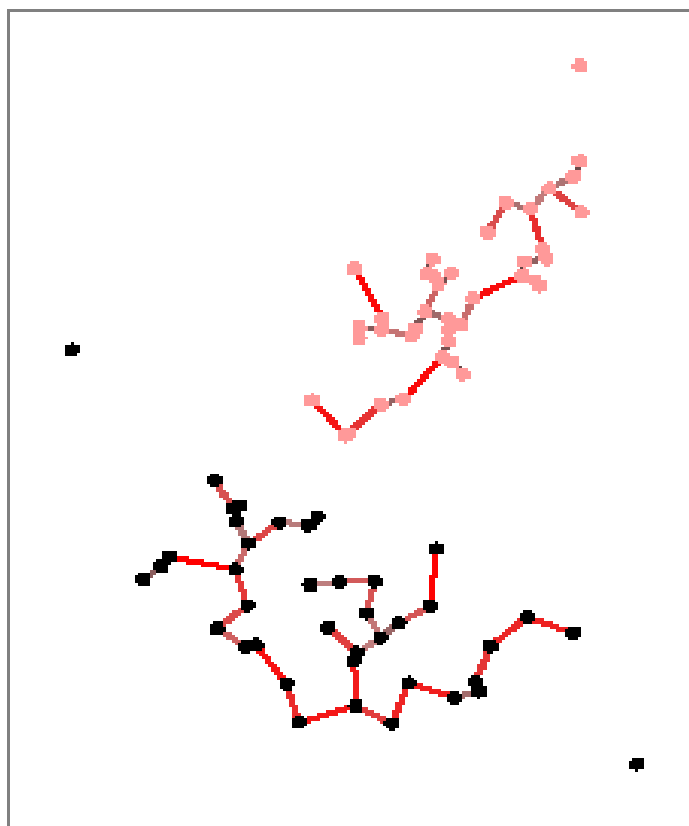
链接效应，产生细长的聚类；

聚类结果对噪声或数据点波动敏感。

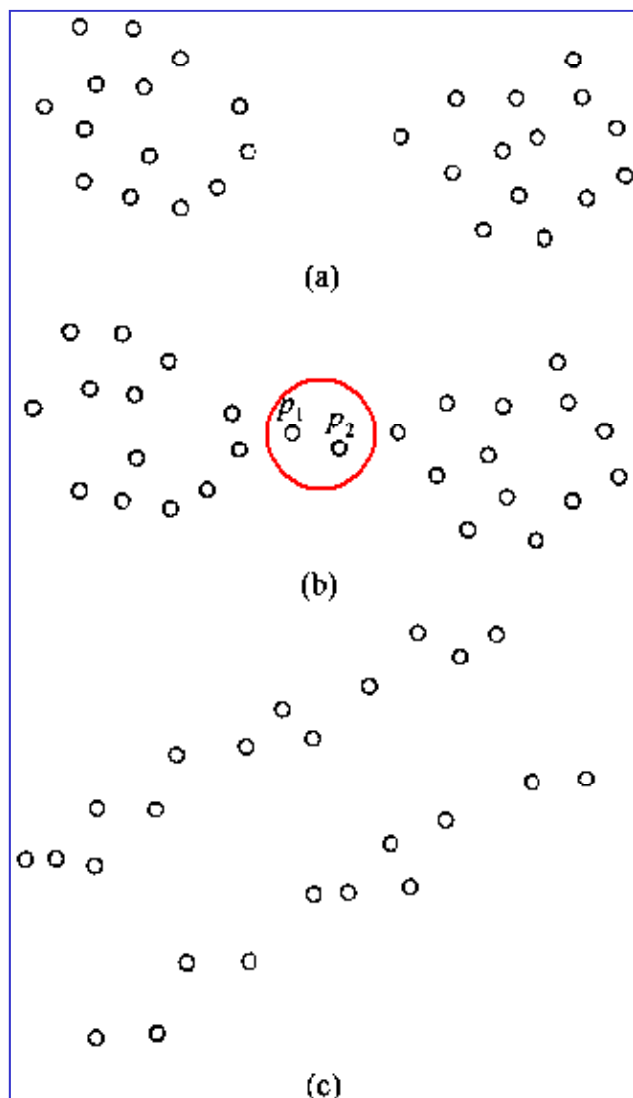
## 图. 二维高斯样本最小距离法层次聚类

左图 无干扰点

右图 有干扰点



## 最小距离法层次聚类 (*single linkage* 或 *single-link*)



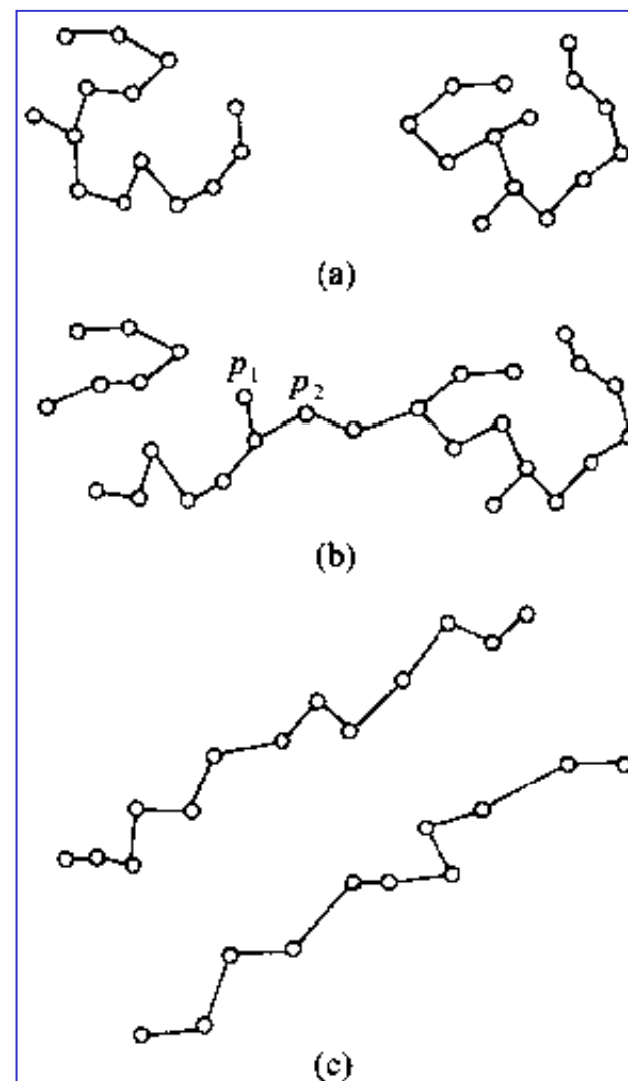
左图：

三种数据分布；

右图：

最近距离法的聚  
类结果

注意：a,b的区别



## (2) 最大距离

$$d_{\max}(C_i, C_j) = \max_{x \in C_i, z \in C_j} \text{dist}(x, z)$$

### 最远邻算法

*(the farthest - neighbor clustering algorithm)*

### 全连接算法

*(complete - linkage algorithm)*

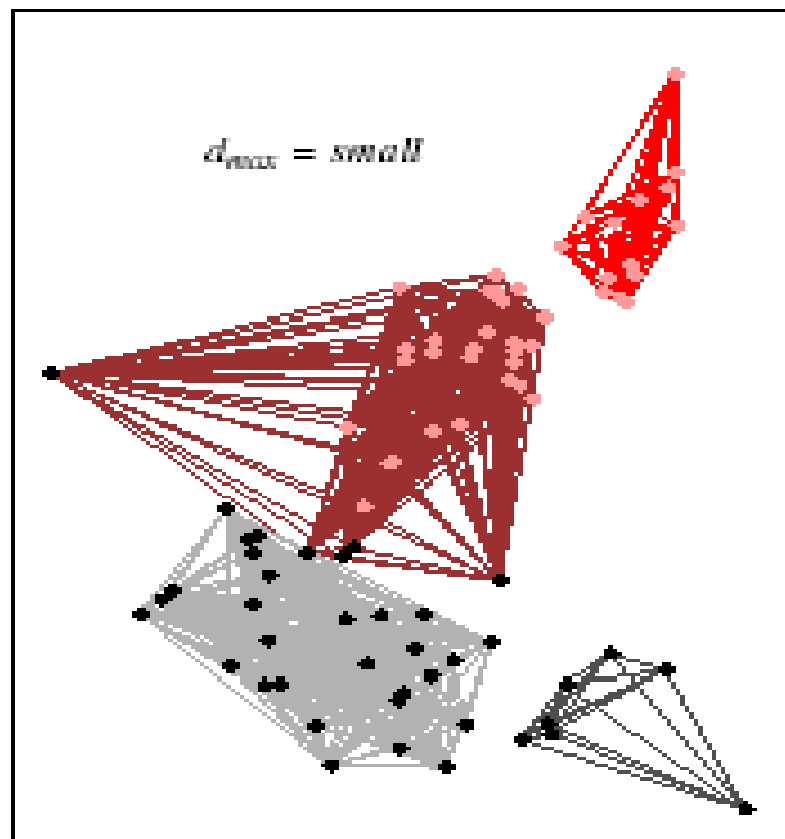
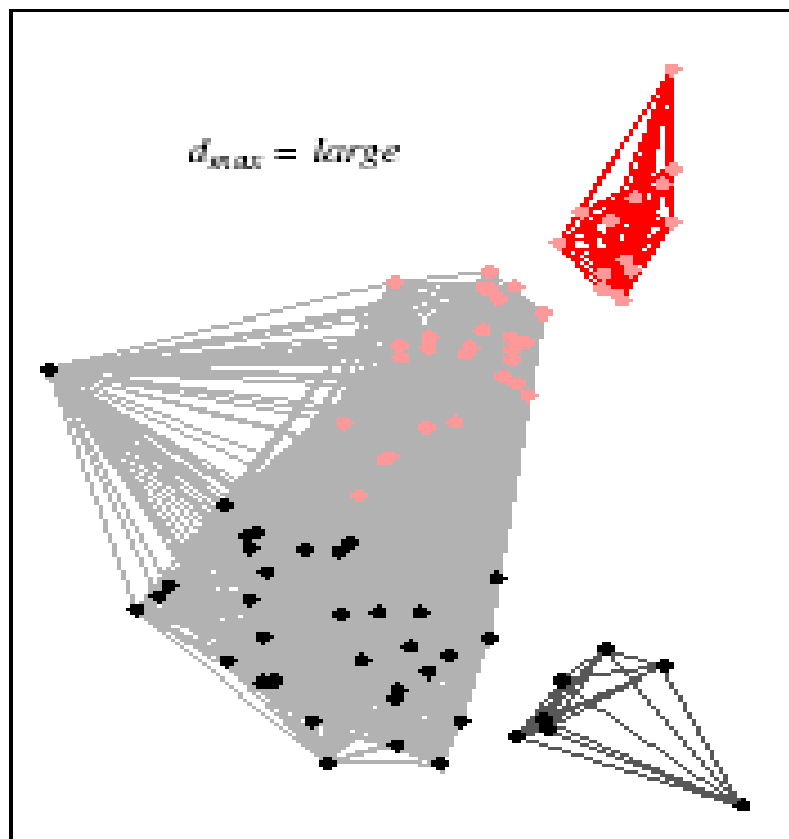
若  $d_{\max}(C_i, C_j) \geq d_0$ , 则聚类过程结束。

特点 { 防止两个密集点集通过某个路径聚为一簇的可能  
不能检测出具有长条形状的聚类  
适合紧密、体积相近的类别划分  
聚类结果对个别远离点敏感

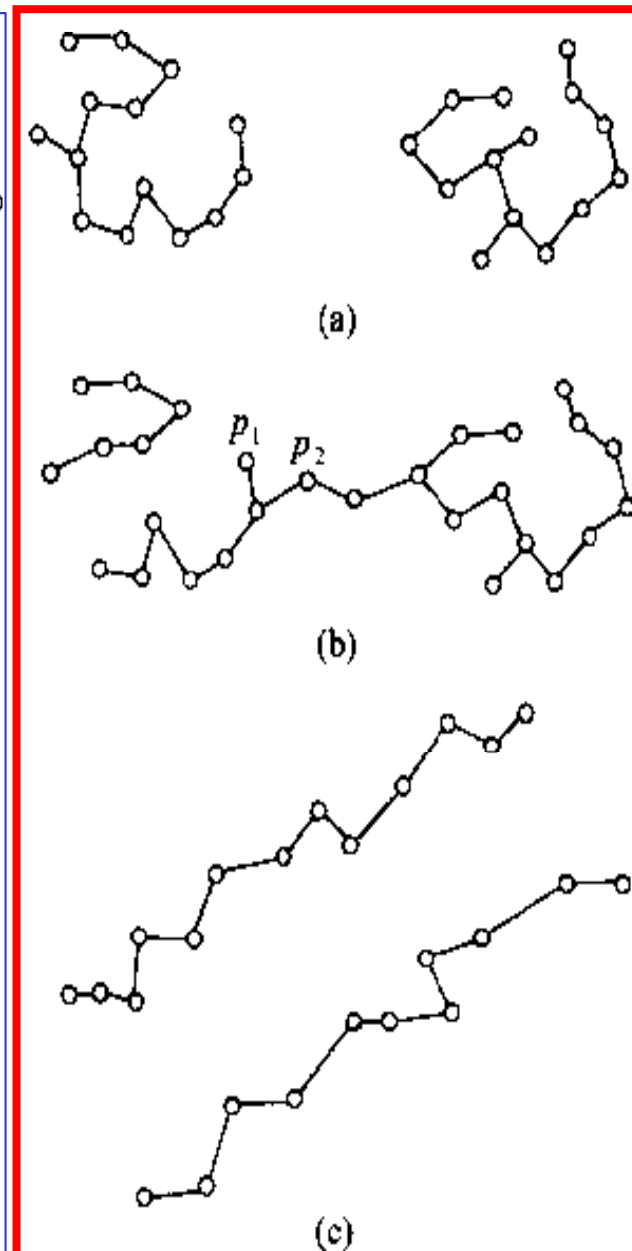
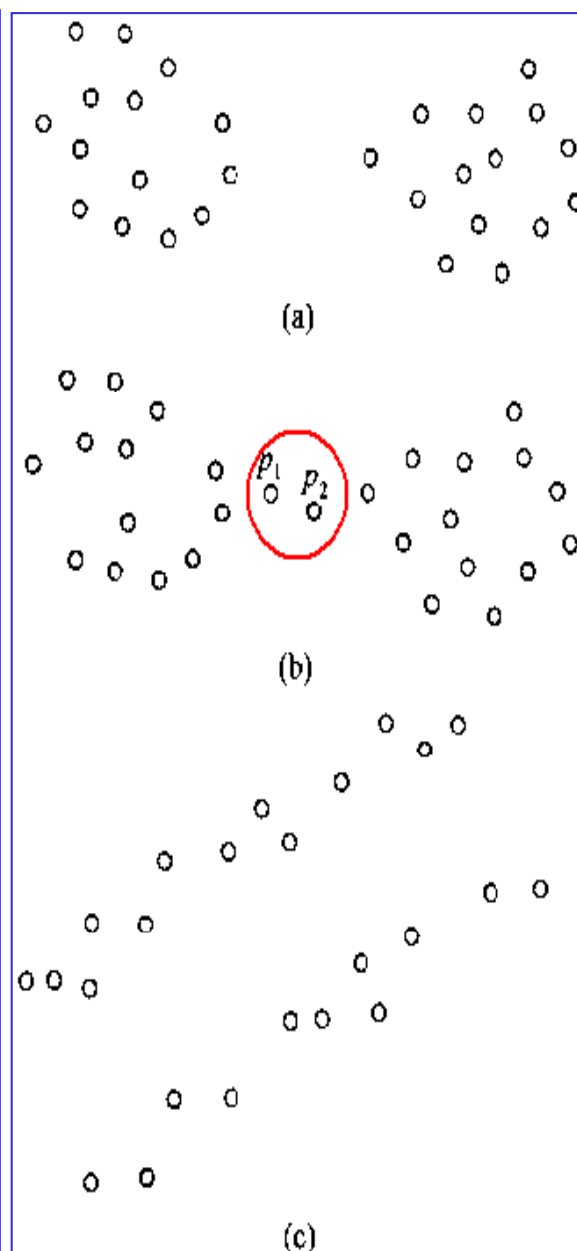
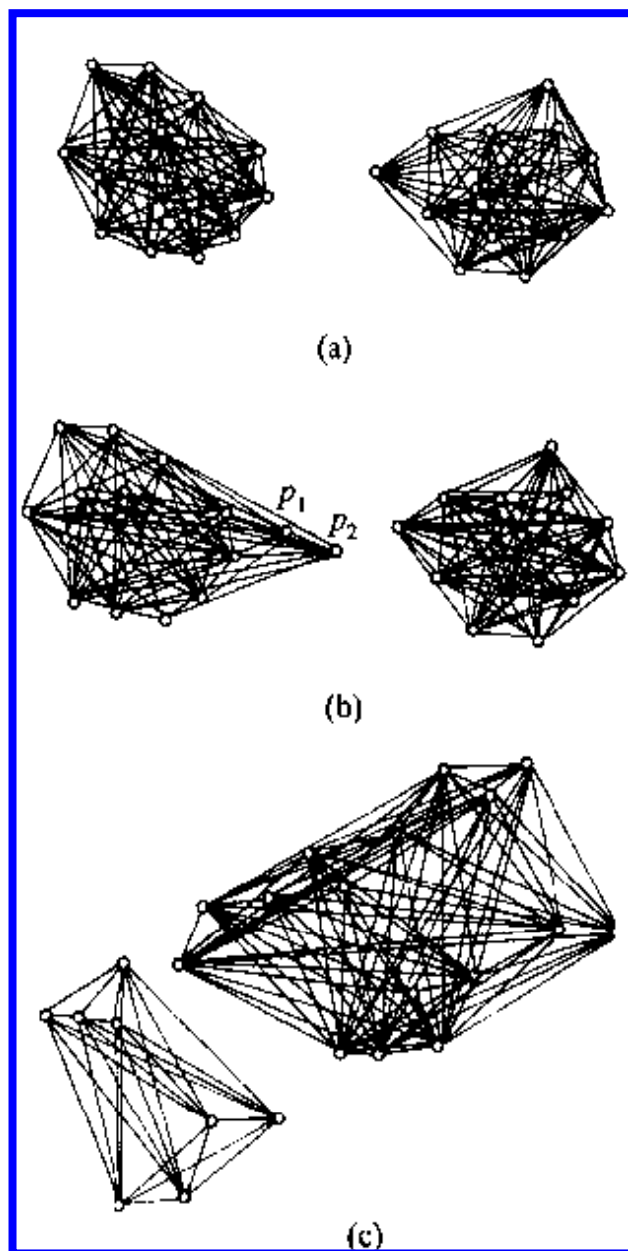


河北师范大学软件学院  
Software College of Hebei Normal University

最大距离阈值 $d_0$ 越大，聚类簇的数目越小。



# 最大距离法与最小距离法聚类结果比较



### (3) 平均距离

$$d_{avg}(C_i, C_j) = \frac{1}{|C_i||C_j|} \sum_{x \in C_i} \sum_{z \in C_j} dist(x, z)$$

相应的聚类算法称为“均链接”算法  
(*average-link, average linkage*).

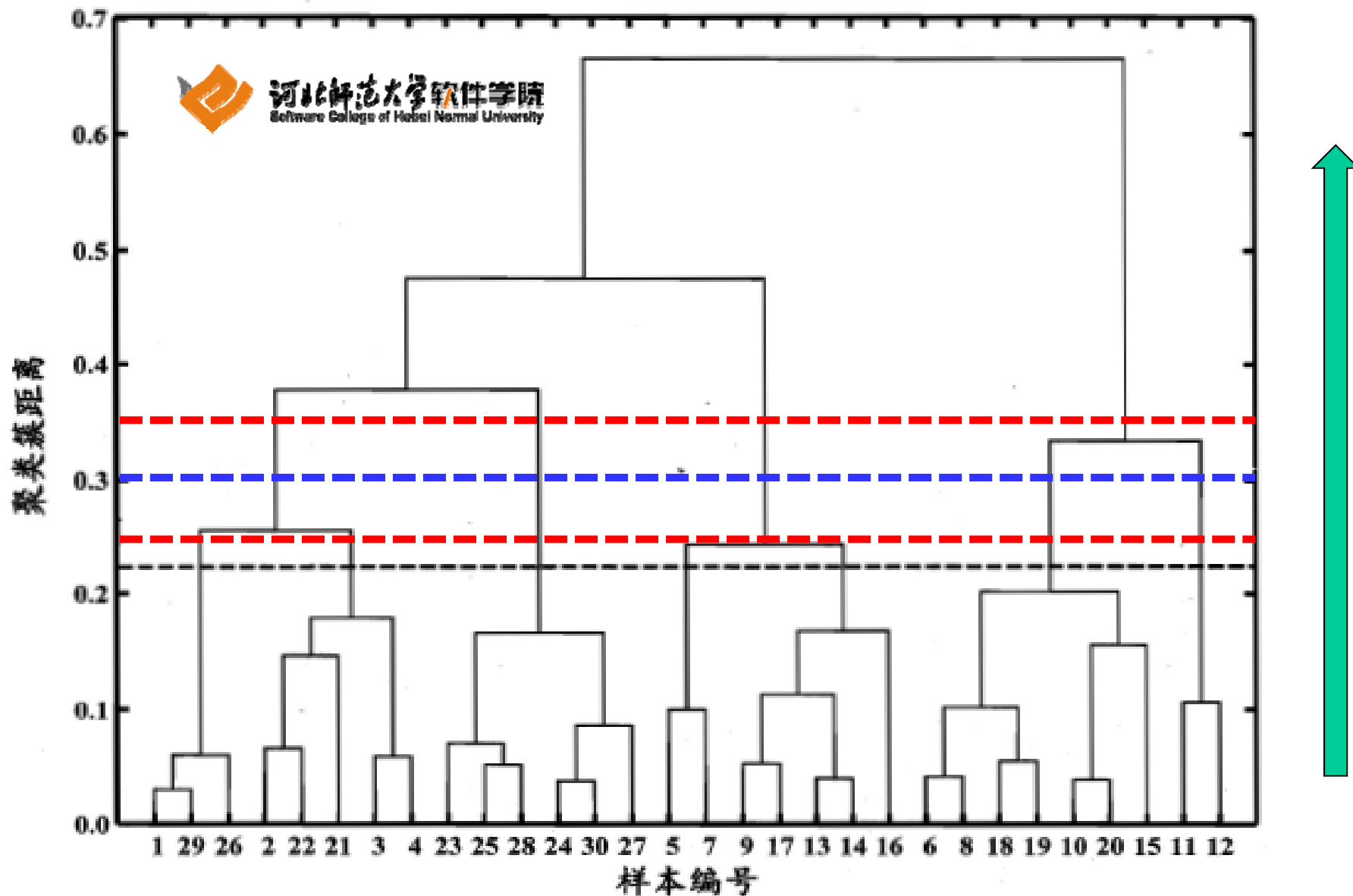
是关于前两种聚类算法的折中；计算简单。



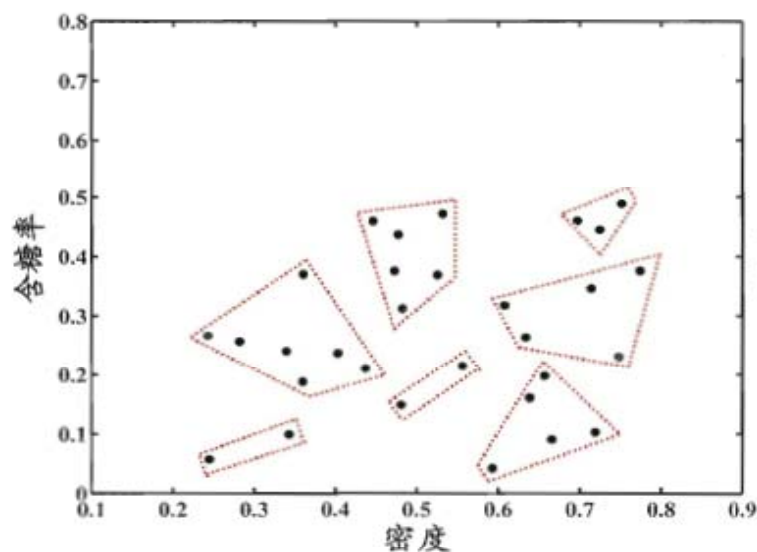


## 西瓜数据集4.0

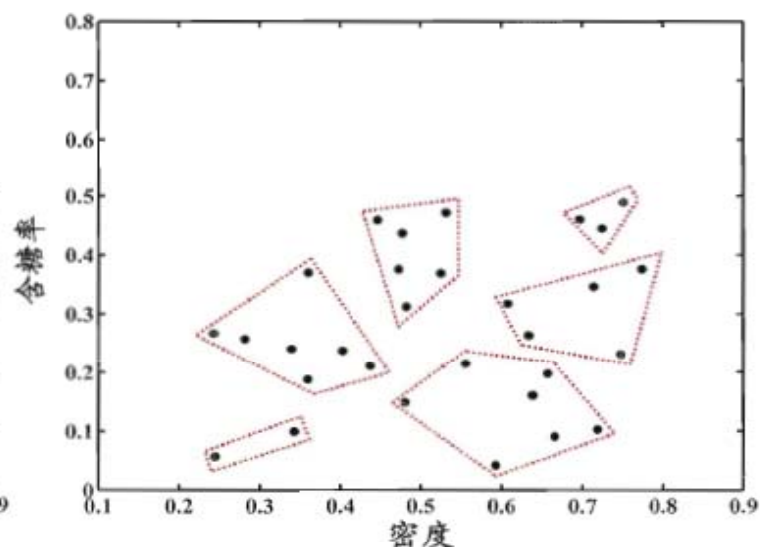
编号	密度	含糖率	编号	密度	含糖率	编号	密度	含糖率
1	0.697	0.460	11	0.245	0.057	21	0.748	0.232
2	0.774	0.376	12	0.343	0.099	22	0.714	0.346
3	0.634	0.264	13	0.639	0.161	23	0.483	0.312
4	0.608	0.318	14	0.657	0.198	24	0.478	0.437
5	0.556	0.215	15	0.360	0.370	25	0.525	0.369
6	0.403	0.237	16	0.593	0.042	26	0.751	0.489
7	0.481	0.149	17	0.719	0.103	27	0.532	0.472
8	0.437	0.211	18	0.359	0.188	28	0.473	0.376
9	0.666	0.091	19	0.339	0.241	29	0.725	0.445
10	0.243	0.267	20	0.282	0.257	30	0.446	0.459



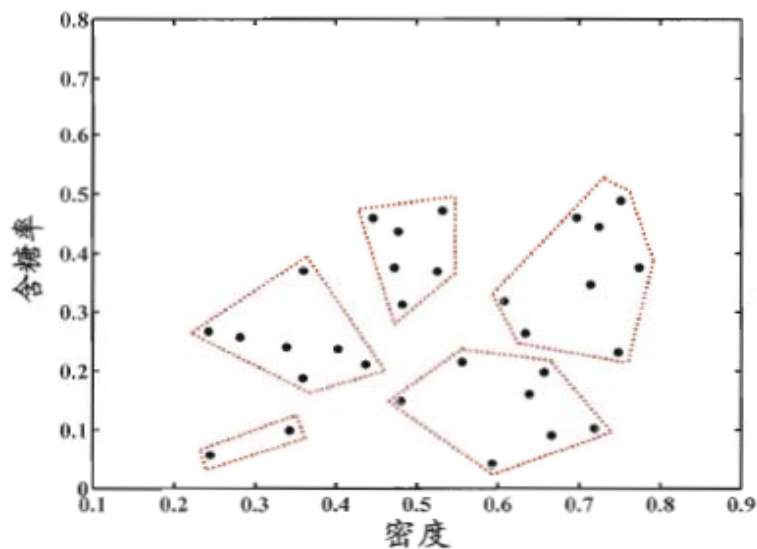
基于西瓜数据集4.0，采用**聚合式层次聚类**，得到的聚类树状图，簇之间的相异性度量采用**最大距离法**。**逐步提升分割层，聚类数目逐渐减小。**



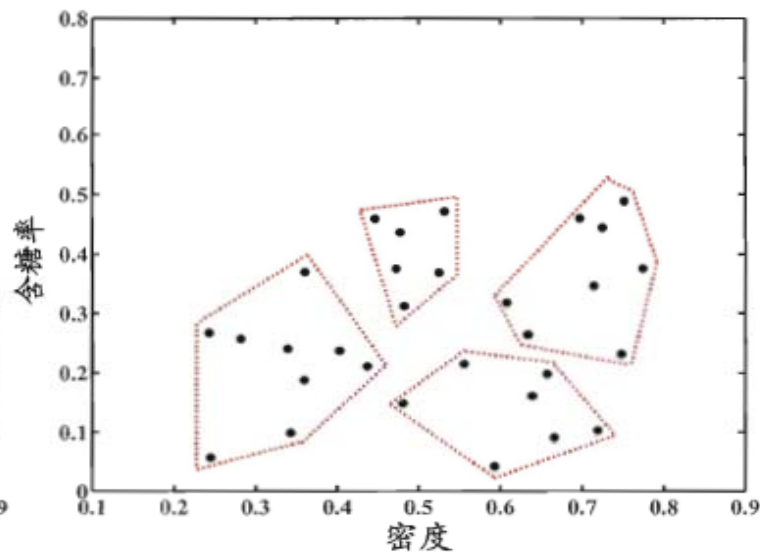
(a) 聚类簇数  $k = 7$



(b) 聚类簇数  $k = 6$



(c) 聚类簇数  $k = 5$



(d) 聚类簇数  $k = 4$

## 思考题

1. 什么是聚类？什么是分类？  
请给出二者的区别与联系。
2. 若采用不同模型对给定的数据集D进行划分.请给出不同聚类算法的实现步骤；影响各聚类算法的因素有哪些？
  - (1) k-均值聚类
  - (2) 高斯混合聚类
  - (3) 层次聚类
  - (4) DBSCAN聚类
3. 上述四种聚类模型的适用场合。

