



机器学习导论

第10章 聚类

谢茂强

南开大学软件学院

本章部分内容来自于吴恩达Coursera网课

目录

07. 聚类可视化

聚类：一种无监督学习算法

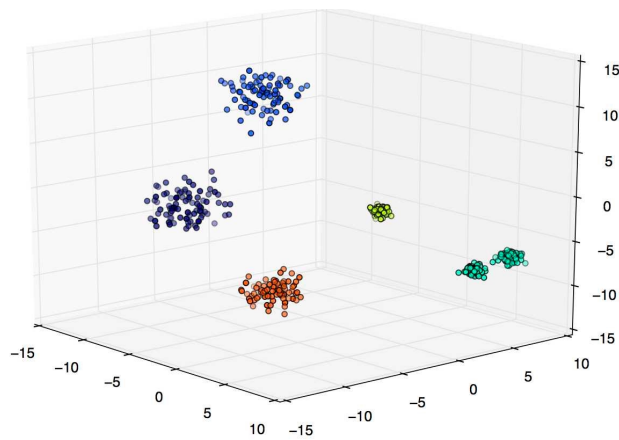
聚类目标： 对数据集进行聚合,使形成一些簇（类），用以了解数据（分布）的内在性质及规律

聚类原则： 物以类聚

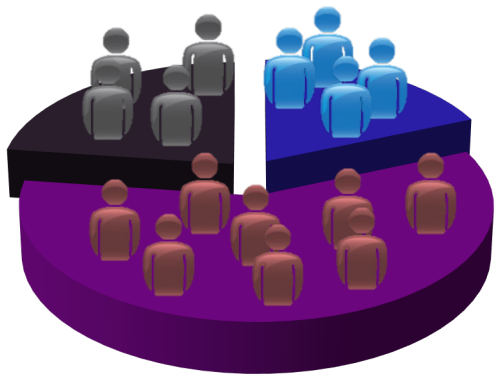
聚类结果： 相似的归为一类

地位： 了解数据分布的重要工具，是后续分析的重要基础

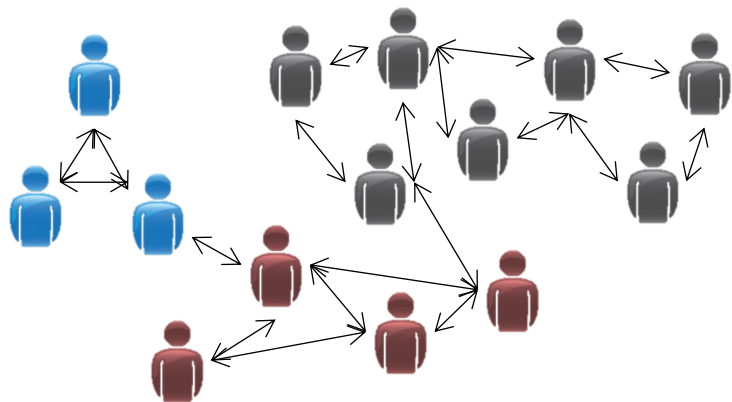
示例：
根据客户的特征进行类别划分



聚类的应用



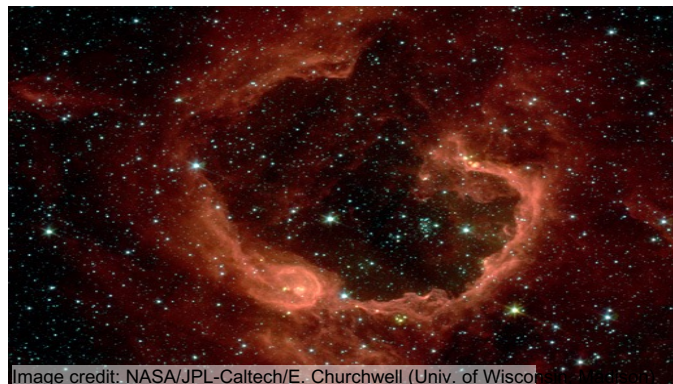
Market segmentation



Social network analysis



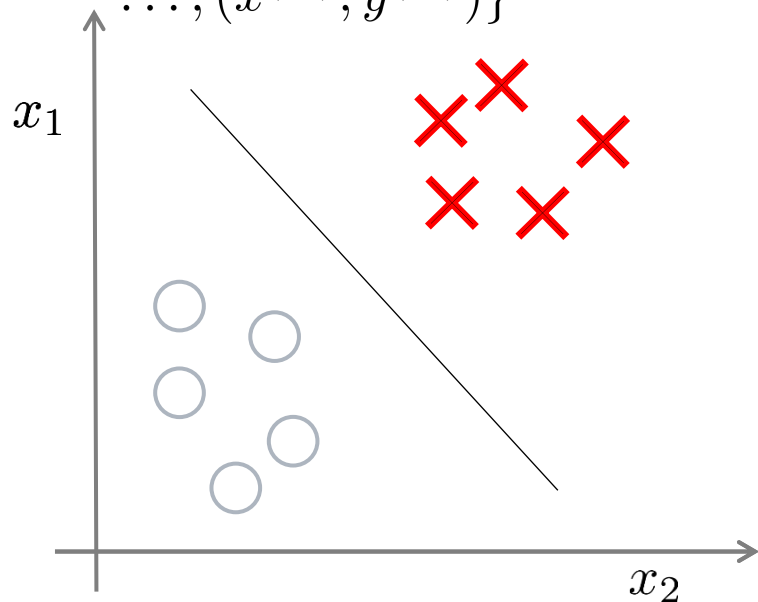
Organize computing clusters



Astronomical data analysis

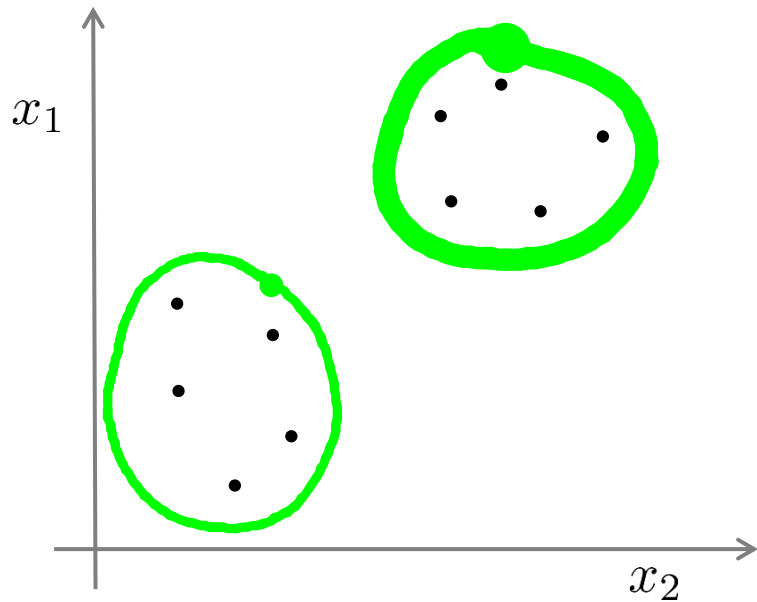
监督学习:

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(m)}, y^{(m)})\}$$



非监督学习:

$$\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$$



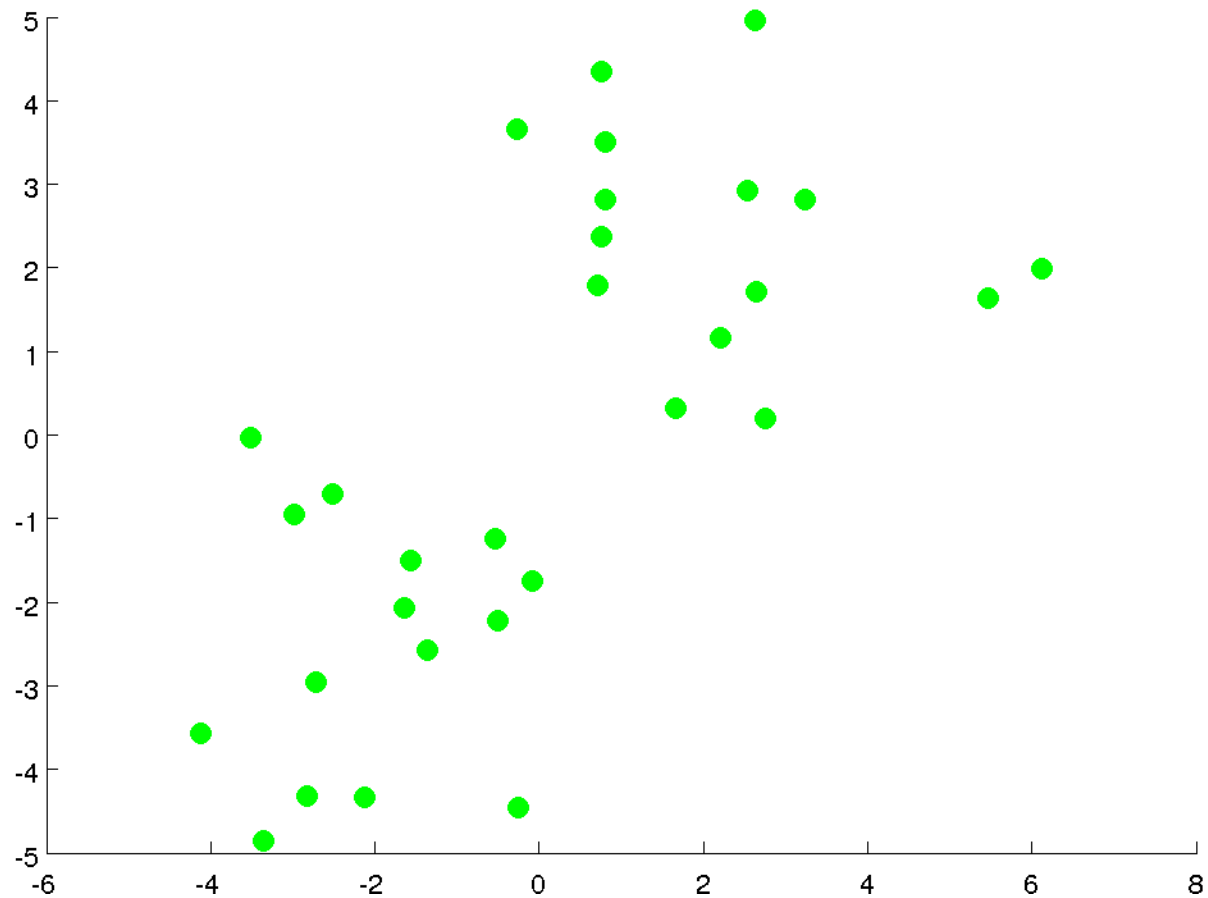
1. 聚类的输入输出定义

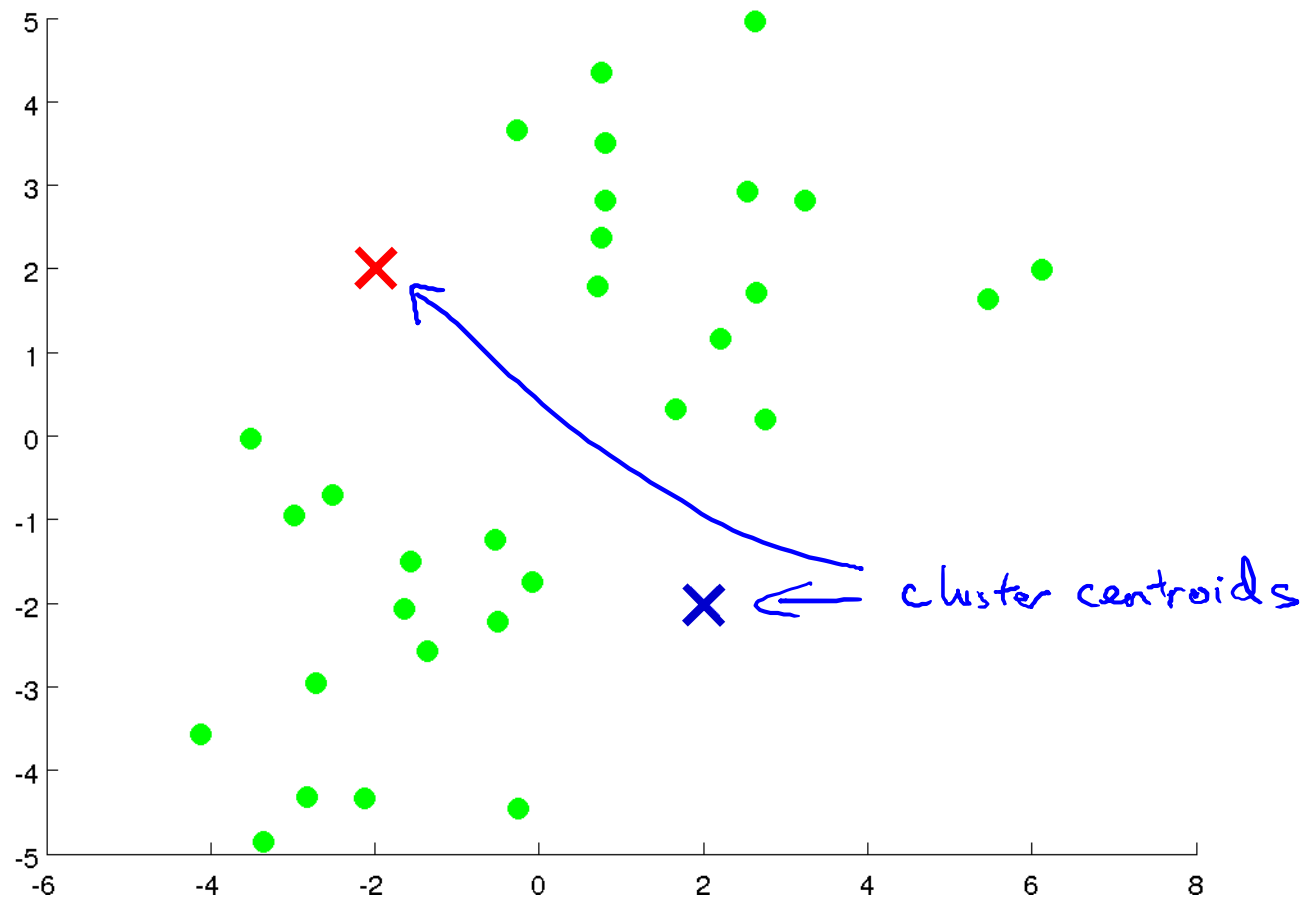
输入: $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ $\mathbf{x}_i = (x_{i1}; x_{i2}; \dots; x_{in})$

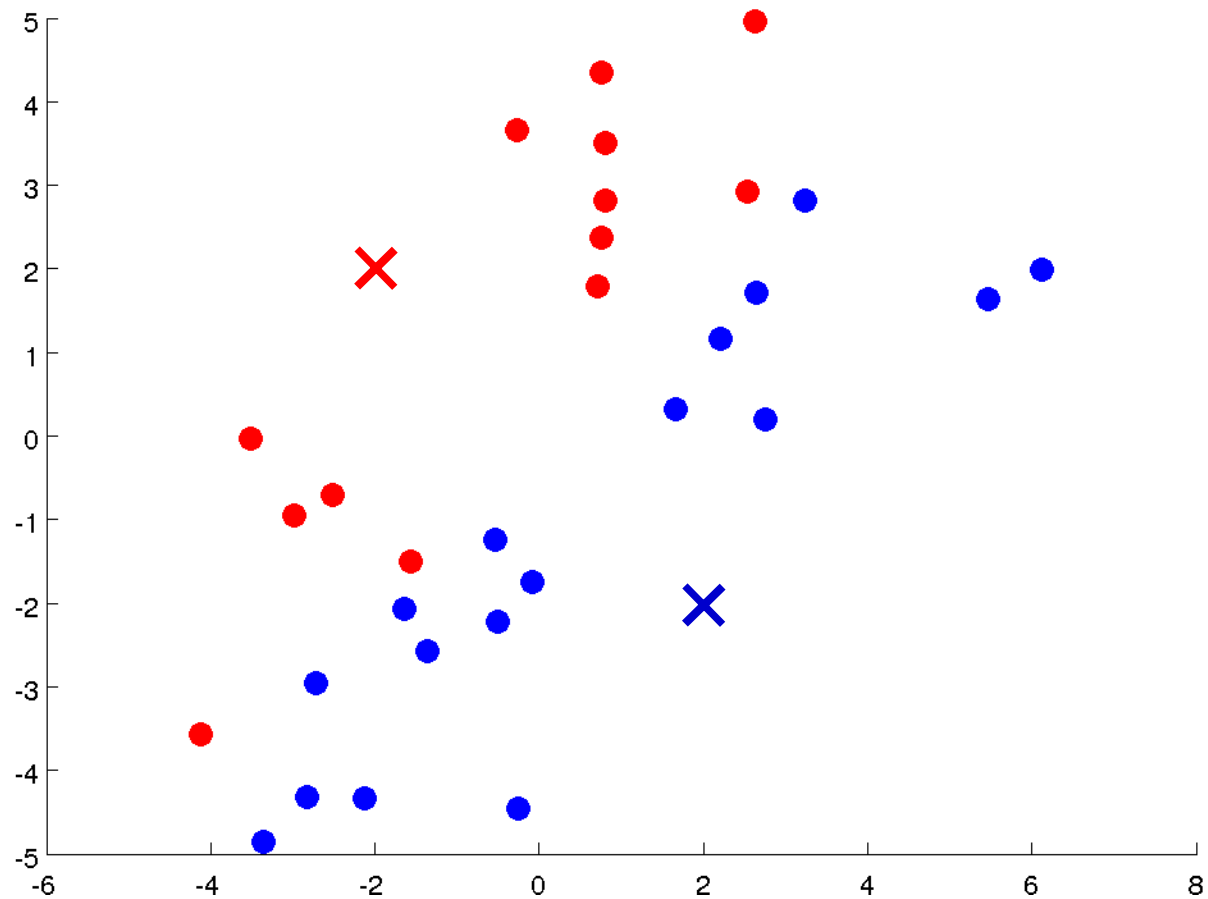
输出: $\mathbf{x} \in C_l$ 簇集合: $\{C_l \mid l = 1, 2, \dots, k\}$

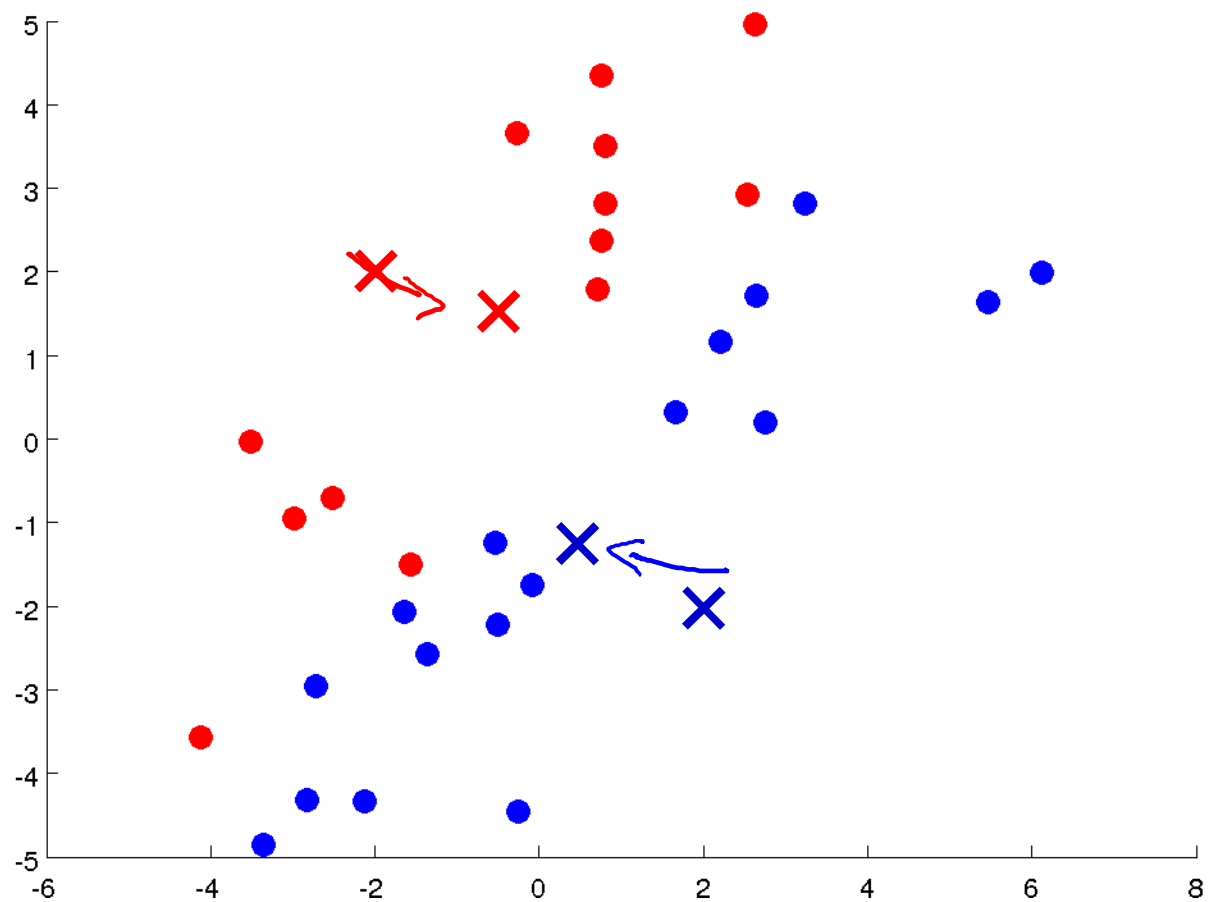
$$C_{l'} \cap_{l' \neq l} C_l = \emptyset \qquad D = \bigcup_{l=1}^k C_l$$

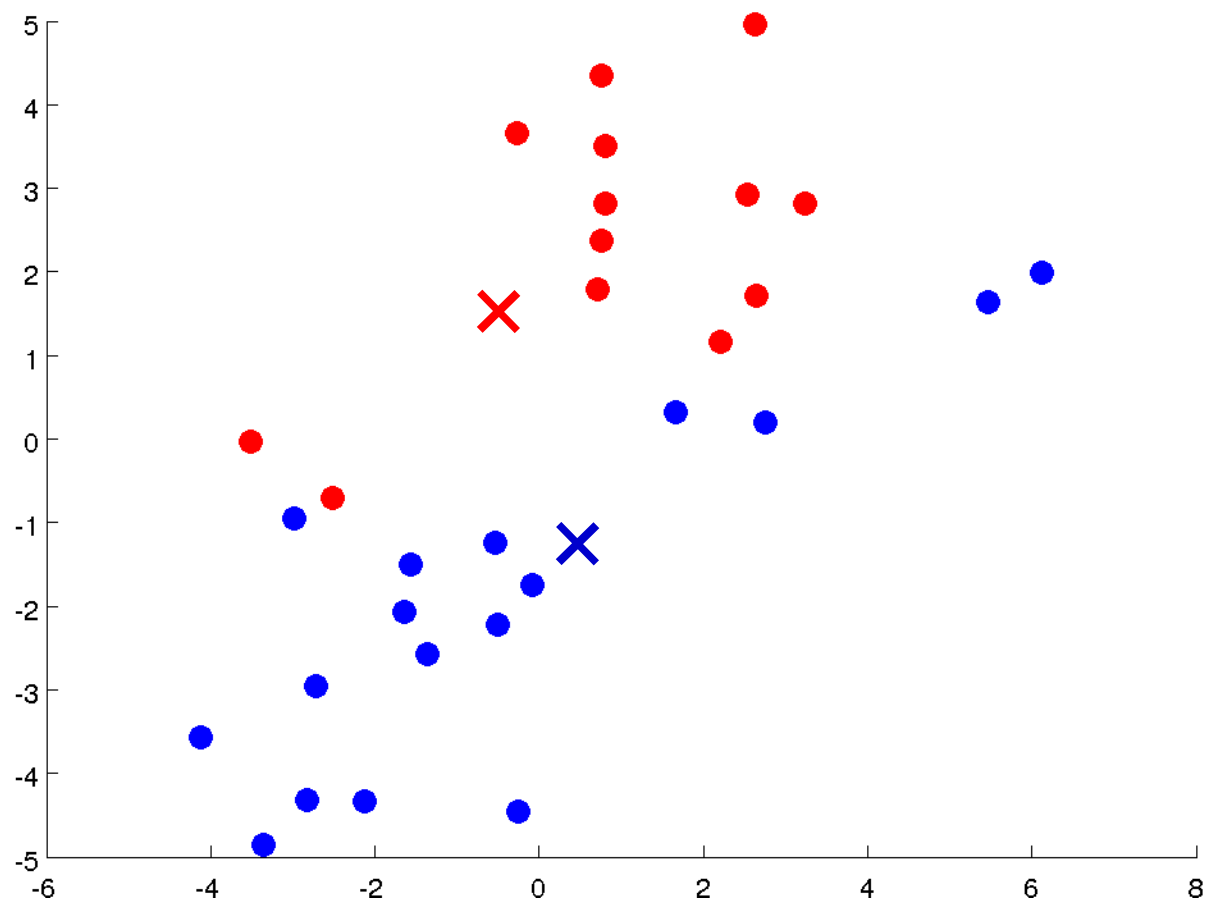
聚类中一种算法思路的示意（下页开始）

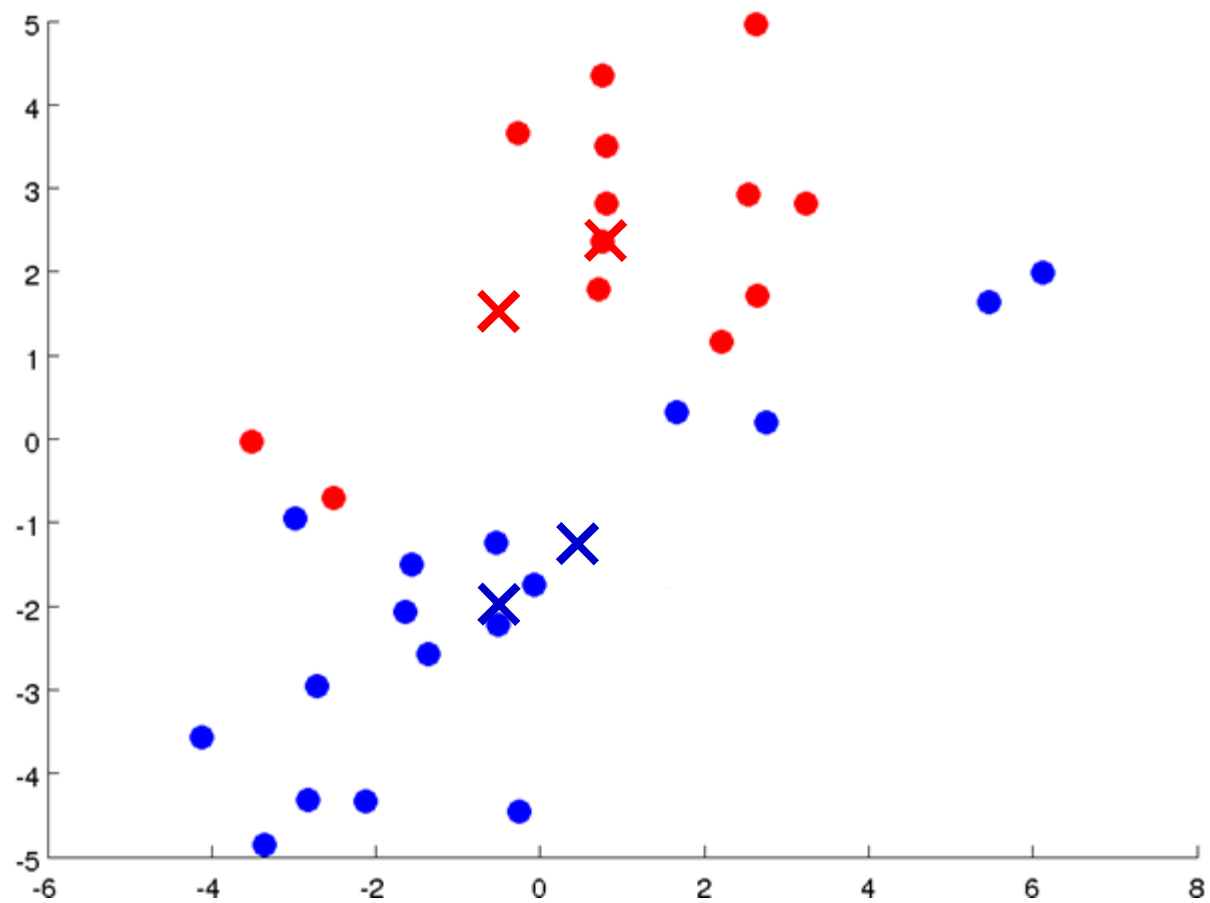


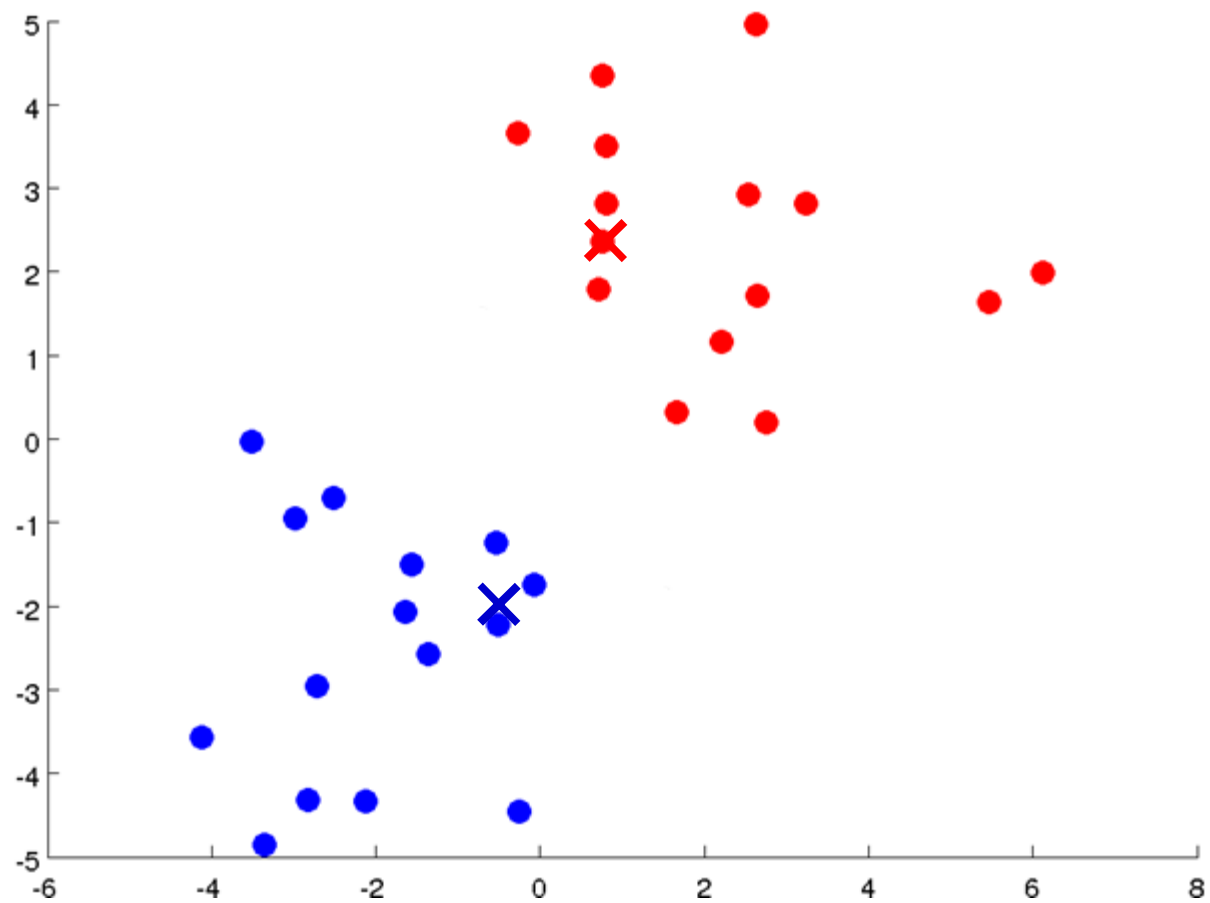


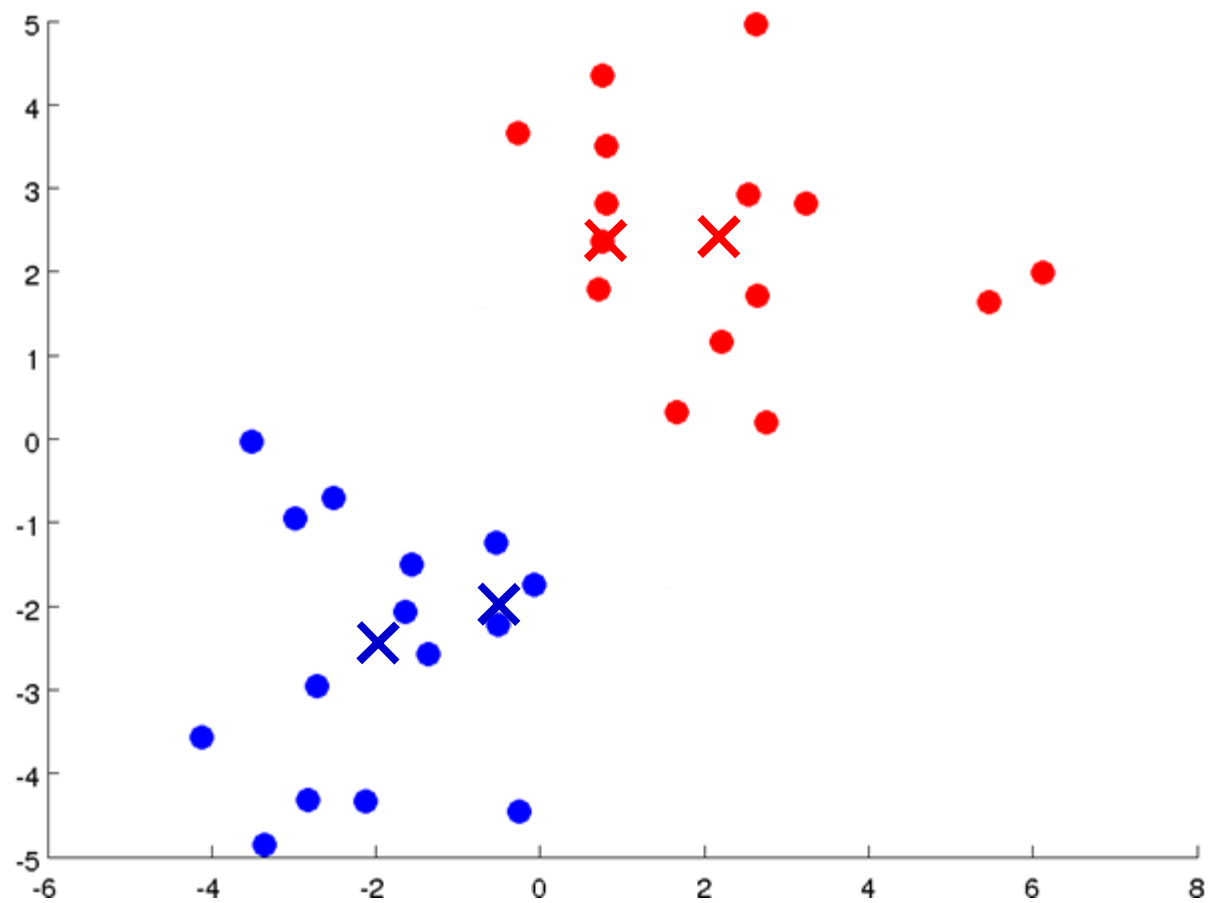


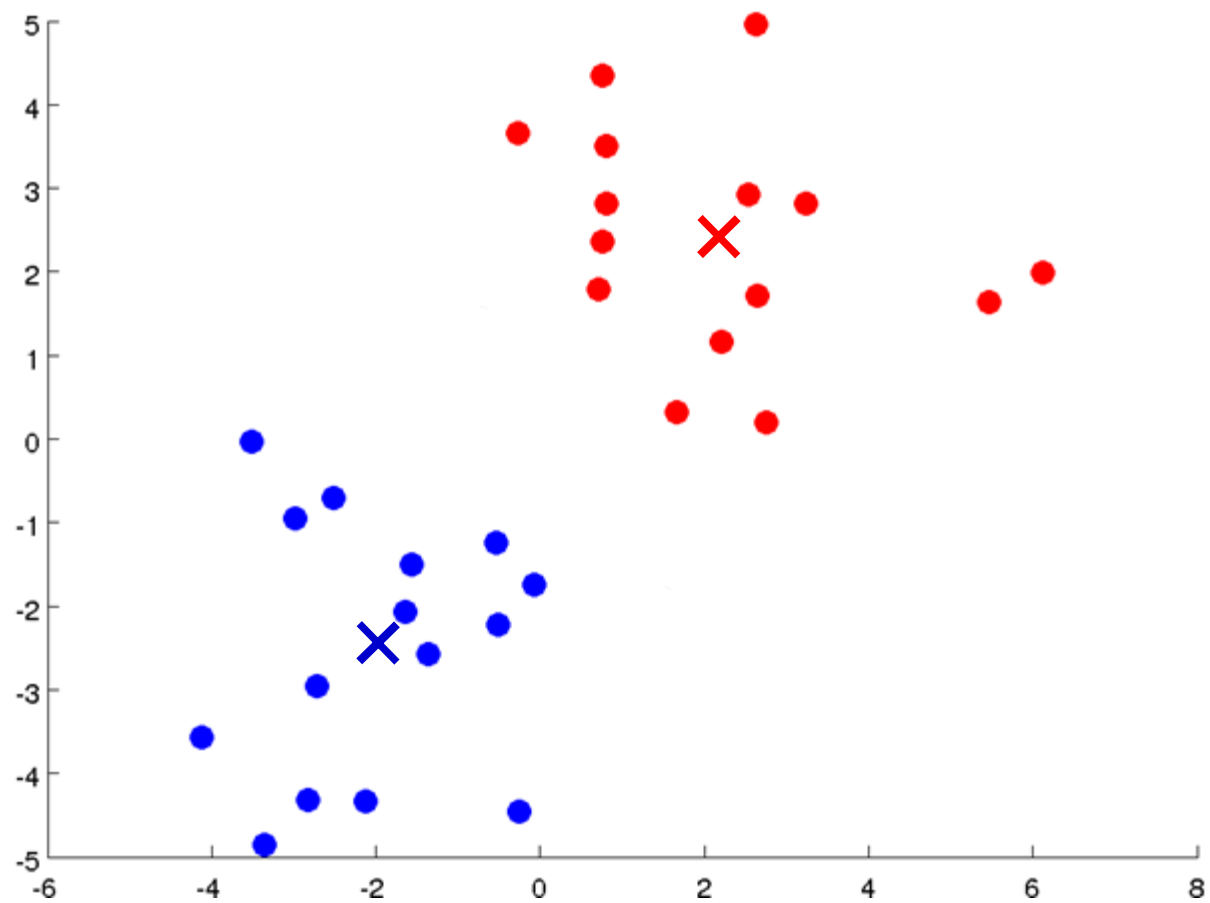












2. 聚类的性能指标

- 聚类性能指标（有效性指标，validity index），用来评价结果，进而作为优化目标。（物以类聚）

外部指标（External Index）

与某个“参考模型”进行比较

比如专家标注的簇标签

内部指标（Internal Index）

簇内相似度（intra-cluster sim）高

簇间相似度（inter-cluster sim）低

类内距、方差 ...

2.1 聚类性能的外部指标 (从“样本对”归属的角度)

对数据集 $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$, 假定通过聚类给出的簇划分为 $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$, 参考模型给出的簇划分为 $\mathcal{C}^* = \{C_1^*, C_2^*, \dots, C_s^*\}$. 相应地, 令 λ 与 λ^* 分别表示与 \mathcal{C} 和 \mathcal{C}^* 对应的簇标记向量. 我们将样本两两配对考虑, 定义

$$a = |SS|, \quad SS = \{(\mathbf{x}_i, \mathbf{x}_j) \mid \lambda_i = \lambda_j, \lambda_i^* = \lambda_j^*, i < j\}, \quad (9.1)$$

$$b = |SD|, \quad SD = \{(\mathbf{x}_i, \mathbf{x}_j) \mid \lambda_i = \lambda_j, \lambda_i^* \neq \lambda_j^*, i < j\}, \quad (9.2)$$

$$c = |DS|, \quad DS = \{(\mathbf{x}_i, \mathbf{x}_j) \mid \lambda_i \neq \lambda_j, \lambda_i^* = \lambda_j^*, i < j\}, \quad (9.3)$$

$$d = |DD|, \quad DD = \{(\mathbf{x}_i, \mathbf{x}_j) \mid \lambda_i \neq \lambda_j, \lambda_i^* \neq \lambda_j^*, i < j\}, \quad (9.4)$$

其中集合 SS 包含了在 \mathcal{C} 中隶属于相同簇且在 \mathcal{C}^* 中也隶属于相同簇的样本对, 集合 SD 包含了在 \mathcal{C} 中隶属于相同簇但在 \mathcal{C}^* 中隶属于不同簇的样本对, ……由于每个样本对 $(\mathbf{x}_i, \mathbf{x}_j)$ ($i < j$) 仅能出现在一个集合中, 因此有 $a + b + c + d = m(m-1)/2$ 成立.

2.1 聚类性能的外部指标

- Jaccard系数 (Jaccard Index, Jaccard Similarity Coefficient)

$$JC = \frac{a}{a + b + c}$$

- FM指数 (Fowlkers & Mallows Index, FMI)

$$FMI = \sqrt{\frac{a}{a + b} \cdot \frac{a}{a + c}}$$

- Rand指数 (Rand Index, RI)

$$RI = \frac{2(a + d)}{m(m - 1)}$$

| | | Groundtruth | |
|------------|----|-------------|-------|
| | | 同簇 | 异簇 |
| Prediction | 同簇 | SS(a) | SD(b) |
| | 异簇 | DS(c) | DD(d) |

2.2 聚类性能的内部指标

- 簇C内样本间平均距离

$$\text{avg}(C) = \frac{2}{|C|(|C| - 1)} \sum_{1 \leq i < j \leq |C|} \text{dist}(\mathbf{x}_i, \mathbf{x}_j) , \quad (9.8)$$

- 簇C内样本间最远距离

$$\text{diam}(C) = \max_{1 \leq i < j \leq |C|} \text{dist}(\mathbf{x}_i, \mathbf{x}_j) , \quad (9.9)$$

2.2 聚类性能的内部指标

- 簇间距离

$$d_{\min}(C_i, C_j) = \min_{\mathbf{x}_i \in C_i, \mathbf{x}_j \in C_j} \text{dist}(\mathbf{x}_i, \mathbf{x}_j) , \quad (9.10)$$

$$d_{\text{cen}}(C_i, C_j) = \text{dist}(\boldsymbol{\mu}_i, \boldsymbol{\mu}_j) , \quad (9.11)$$

- DB指数(Davies-Bouldin Index, DBI)

$$\text{DBI} = \frac{1}{k} \sum_{i=1}^k \max_{j \neq i} \left(\frac{\text{avg}(C_i) + \text{avg}(C_j)}{d_{\text{cen}}(\boldsymbol{\mu}_i, \boldsymbol{\mu}_j)} \right) . \quad (9.12)$$

DBI越小聚类性能越好

2.2 聚类性能的内部指标

- 簇间距离

$$d_{\min}(C_i, C_j) = \min_{\mathbf{x}_i \in C_i, \mathbf{x}_j \in C_j} \text{dist}(\mathbf{x}_i, \mathbf{x}_j) , \quad (9.10)$$

$$d_{\text{cen}}(C_i, C_j) = \text{dist}(\boldsymbol{\mu}_i, \boldsymbol{\mu}_j) , \quad (9.11)$$

- Dunn指数(Dunn Index, DI)

$$\text{DI} = \min_{1 \leq i \leq k} \left\{ \min_{j \neq i} \left(\frac{d_{\min}(C_i, C_j)}{\max_{1 \leq l \leq k} \text{diam}(C_l)} \right) \right\} . \quad (9.13)$$

DI越大聚类性能越好

3. 距离度量

- 非负性: $\text{dist}(\mathbf{x}_i, \mathbf{x}_j) \geq 0$
- 同一性: $\text{dist}(\mathbf{x}_i, \mathbf{x}_j) = 0$ 当且仅当 $\mathbf{x}_i = \mathbf{x}_j$
- 对称性: $\text{dist}(\mathbf{x}_i, \mathbf{x}_j) = \text{dist}(\mathbf{x}_j, \mathbf{x}_i)$
- 三角不等式: $\text{dist}(\mathbf{x}_i, \mathbf{x}_j) \leq \text{dist}(\mathbf{x}_i, \mathbf{x}_k) + \text{dist}(\mathbf{x}_k, \mathbf{x}_j)$

3. 距离度量

- 明氏距离(Minkowski distance, L_p 范数):

$$\text{dist}_{\text{mk}}(\mathbf{x}_i, \mathbf{x}_j) = \left(\sum_{u=1}^n |x_{iu} - x_{ju}|^p \right)^{\frac{1}{p}}$$

- 欧氏距离(Euclidean distance):

$$\text{dist}_{\text{ed}}(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_2 = \sqrt{\sum_{u=1}^n |x_{iu} - x_{ju}|^2}$$

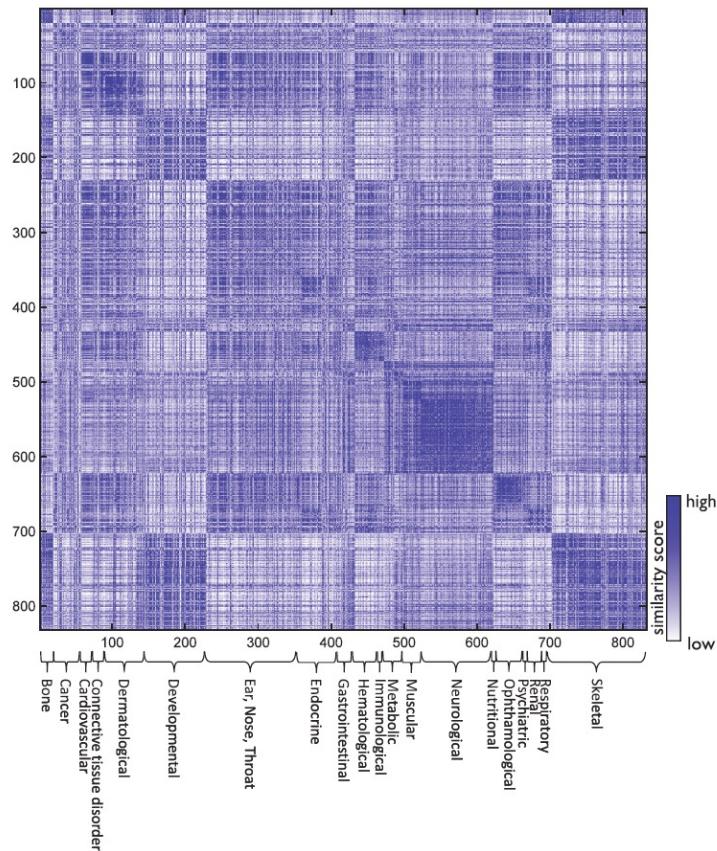
- 曼哈顿距离(Manhattan/City block distance):

$$\text{dist}_{\text{man}}(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_1 = \sum_{u=1}^n |x_{iu} - x_{ju}|$$

3. 距离和空间的关系

- 距离是定义在空间定义之上，是空间的附属
- 常见的空间问题
 - 尺度伸缩
 - 球面距离
 - 流形 (manifold)

示例：可以使用Jaccard Index 作为语义相似度度量



4. 基于划分的K-Means算法

K-means algorithm

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

Repeat {

 for $i = 1$ to m

$c^{(i)} :=$ index (from 1 to K) of cluster centroid
 closest to $x^{(i)}$

 for $k = 1$ to K

$\mu_k :=$ average (mean) of points assigned to k cluster

}



4. K-Means算法的优化目标

给定样本集 $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$, “ k 均值” (k -means) 算法针对聚类所得簇划分 $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$ 最小化平方误差

$$E = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|_2^2, \quad (9.24)$$

其中 $\boldsymbol{\mu}_i = \frac{1}{|C_i|} \sum_{\mathbf{x} \in C_i} \mathbf{x}$ 是簇 C_i 的均值向量. 直观来看, 式(9.24) 在一定程度上刻画了簇内样本围绕簇均值向量的紧密程度, E 值越小则簇内样本相似度越高.

簇中心的随机初始化

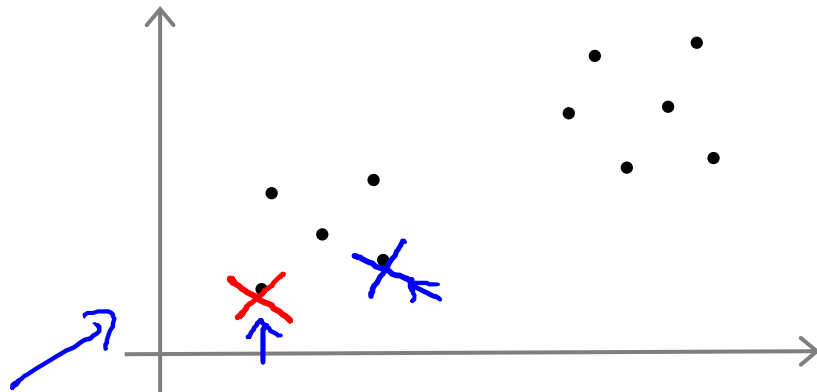
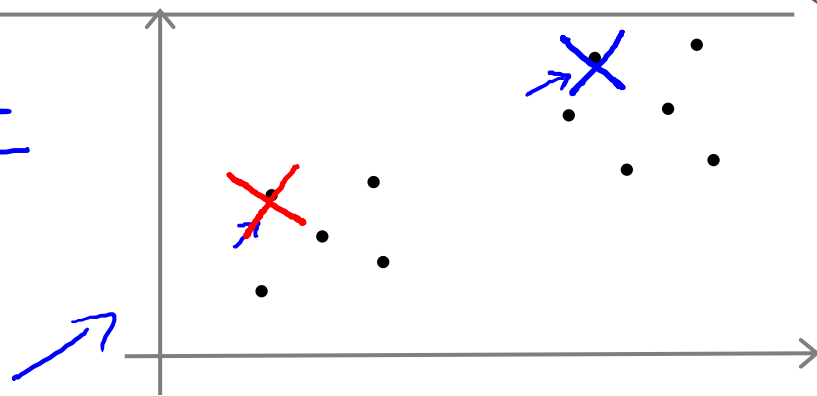
Should have $K < m$

Randomly pick K training examples.

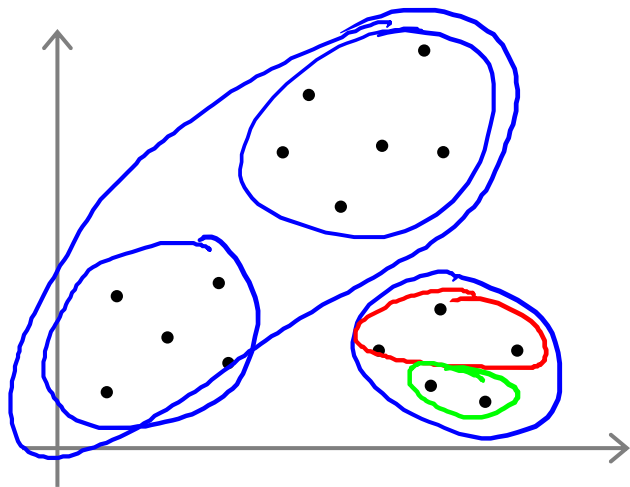
Set μ_1, \dots, μ_K equal to these K examples.

$$\begin{aligned}\mu_1 &= x^{(i)} \\ \mu_2 &= x^{(j)} \\ &\vdots\end{aligned}$$

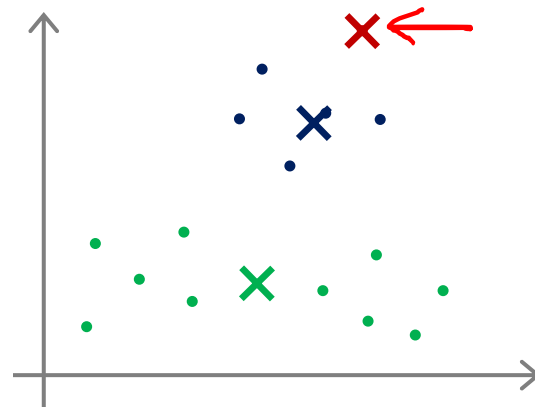
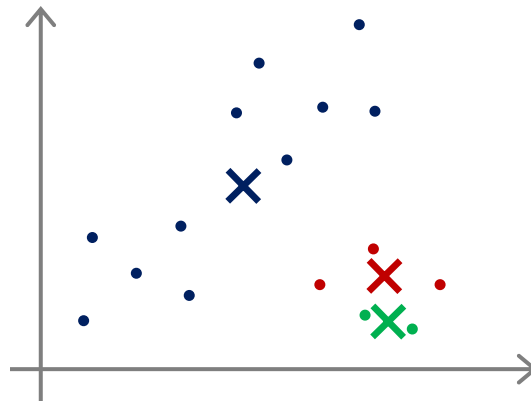
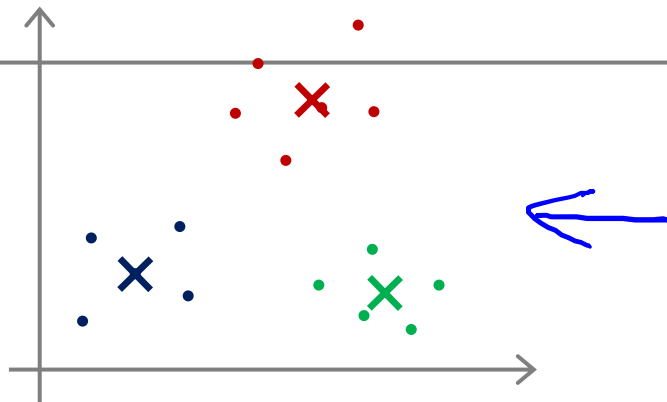
$K=2$



K-Means算法的局部收敛结果



$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$$



通过多次随机初始化解解决陷入局部极值的问题

For $i = 1$ to 100 {

Randomly initialize K-means.

Run K-means. Get $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K$.

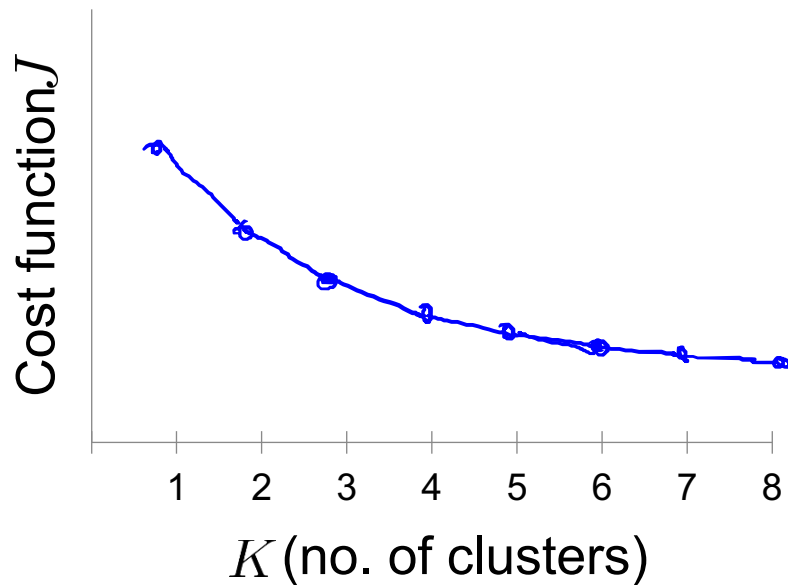
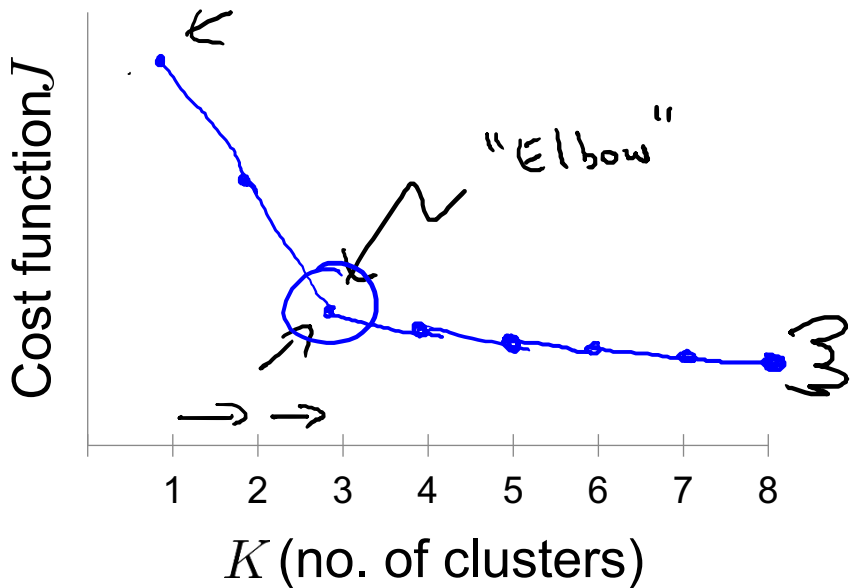
Compute cost function (distortion)

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

}

Pick clustering that gave lowest cost $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

选择合适聚类簇数K (Elbow方法)



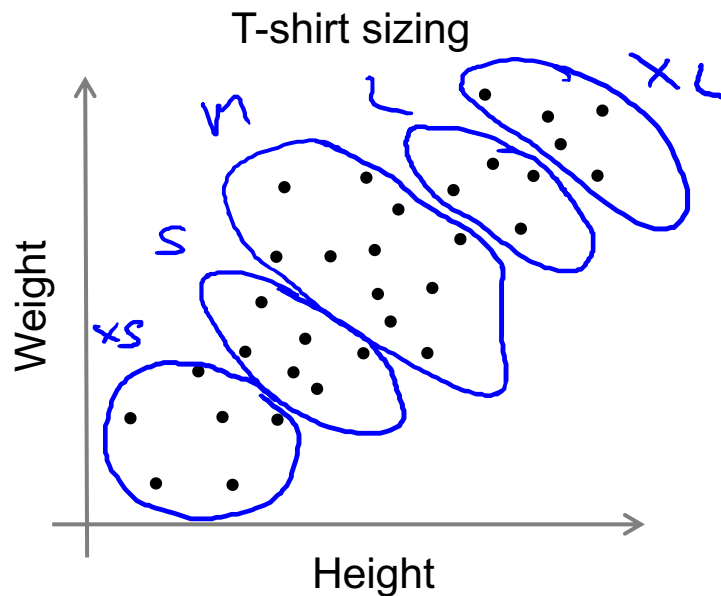
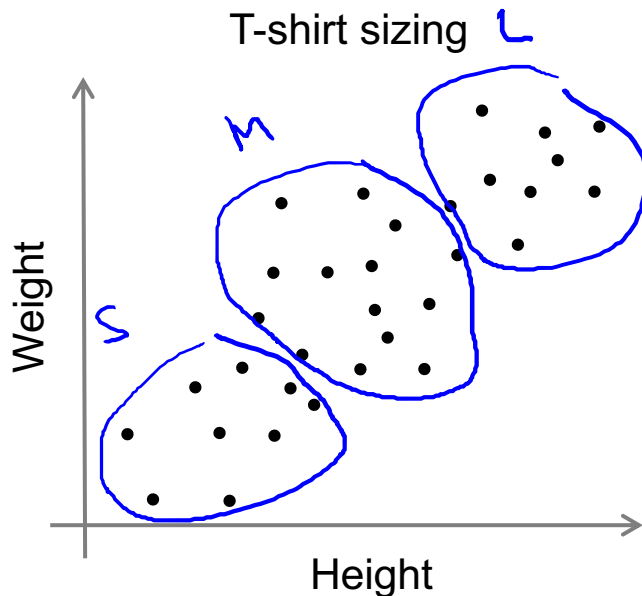
选择合适聚类簇数K

Sometimes, you're running K-means to get clusters to use for some later/downstream purpose. Evaluate K-means based on a metric for how well it performs for that later purpose.

$K=3$ S, M, L

$K=5$ XS, S, M, L, XL

E.g.



4.2 Gaussian Mixture Model

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 .

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
 $\gamma_{ji} = p_{\mathcal{M}}(z_j = i \mid \mathbf{x}_j) \ (1 \leq i \leq k)$
- 5: **end for**



- 6: **for** $i = 1, 2, \dots, k$ **do**
- 7: 计算新均值向量: $\mu'_i = \frac{\sum_{j=1}^m \gamma_{ji} x_j}{\sum_{j=1}^m \gamma_{ji}};$
- 8: 计算新协方差矩阵: $\Sigma'_i = \frac{\sum_{j=1}^m \gamma_{ji} (x_j - \mu'_i)(x_j - \mu'_i)^T}{\sum_{j=1}^m \gamma_{ji}};$
- 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: $C_i = \emptyset \ (1 \leq i \leq k)$
- 14: **for** $j = 1, 2, \dots, m$ **do**
- 15: 根据式(9.31)确定 x_j 的簇标记 λ_j ;
- 16: 将 x_j 划入相应的簇: $C_{\lambda_j} = C_{\lambda_j} \cup \{x_j\}$
- 17: **end for**
- 输出: 簇划分 $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$
-

5. Density-based Clustering

- DBSCAN(Density-Based Spatial Clustering of Applications with Noise) 连通簇簇定义(P212)

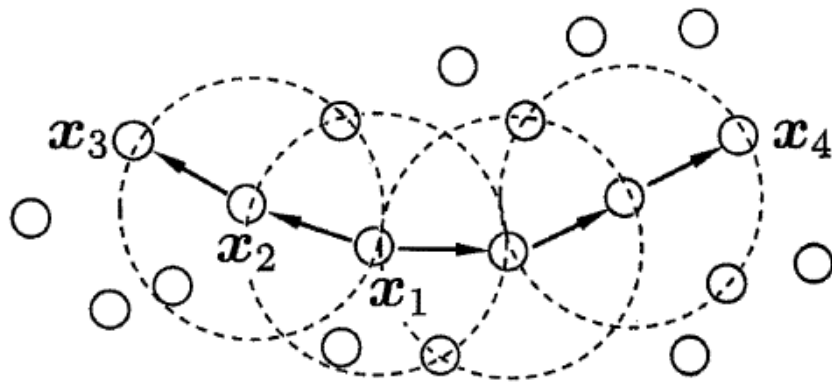


图 9.8 DBSCAN 定义的基本概念($MinPts = 3$): 虚线显示出 ϵ -邻域, x_1 是核心对象, x_2 由 x_1 密度直达, x_3 由 x_1 密度可达, x_3 与 x_4 密度相连.

5. DBSCAN

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

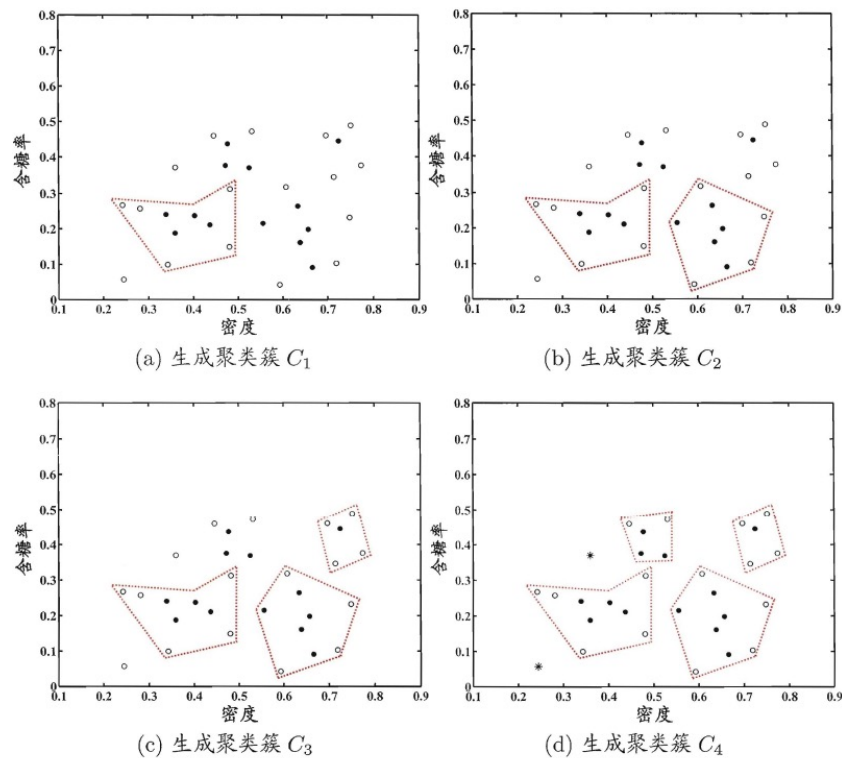


图 9.10 DBSCAN 算法($\epsilon = 0.11$, $MinPts = 5$)生成聚类簇的先后情况. 核心对象、非核心对象、噪声样本分别用“●”“○”“*”表示, 红色虚线显示出簇划分.

6. Hierarchical Clustering

- 在不同层次对数据集进行划分，从而形成树形的聚类结构
 - “自底向上” 的聚合策略 (AGNES)
 - “自顶向下” 的分拆策略
- AGNES (Agglomerative Nesting, 凝聚聚类)

6. AGNES

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  $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$ 
    
```


- 类-类距离

$$\text{最小距离: } d_{\min}(C_i, C_j) = \min_{\mathbf{x} \in C_i, \mathbf{z} \in C_j} \text{dist}(\mathbf{x}, \mathbf{z}) , \quad (9.41)$$

$$\text{最大距离: } d_{\max}(C_i, C_j) = \max_{\mathbf{x} \in C_i, \mathbf{z} \in C_j} \text{dist}(\mathbf{x}, \mathbf{z}) , \quad (9.42)$$

$$\text{平均距离: } d_{\text{avg}}(C_i, C_j) = \frac{1}{|C_i||C_j|} \sum_{\mathbf{x} \in C_i} \sum_{\mathbf{z} \in C_j} \text{dist}(\mathbf{x}, \mathbf{z}) . \quad (9.43)$$

- 对象-类距离

- 对象-对象距离

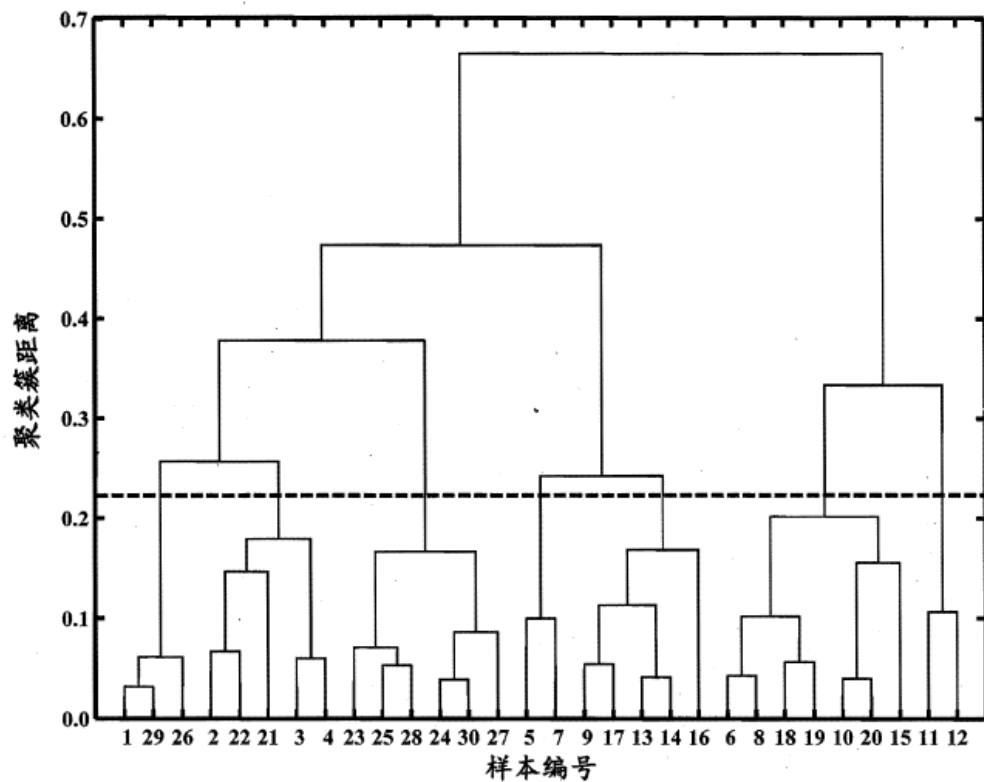
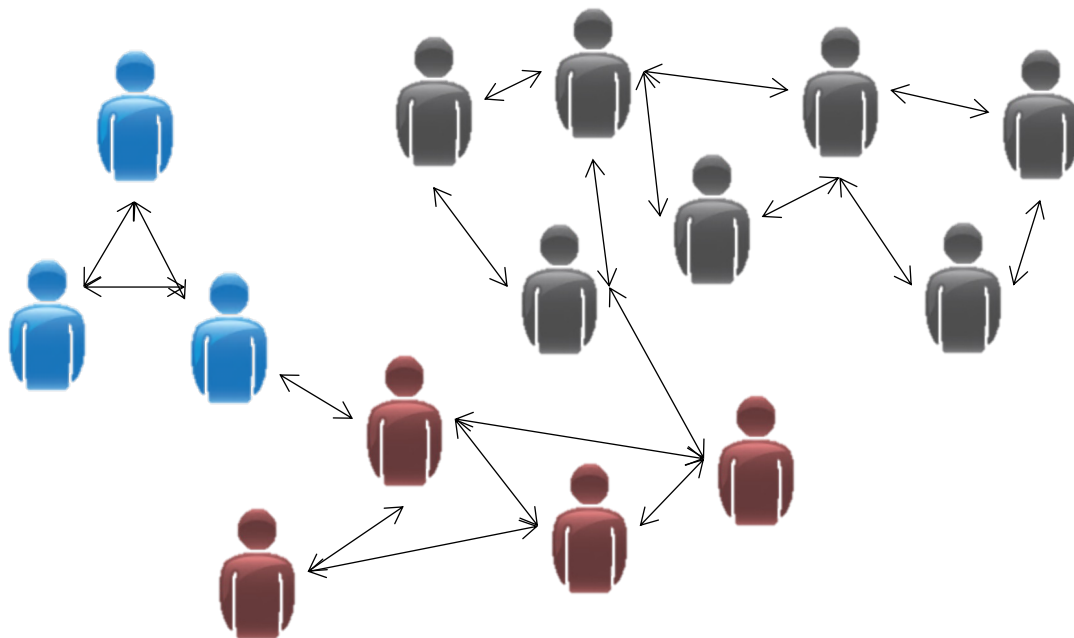


图 9.12 西瓜数据集 4.0 上 AGNES 算法生成的树状图(采用 d_{\max}). 横轴对应于样本编号, 纵轴对应于聚类簇距离.

其他：网络数据聚类算-Mincut



7. 双向层次聚类结果的可视化

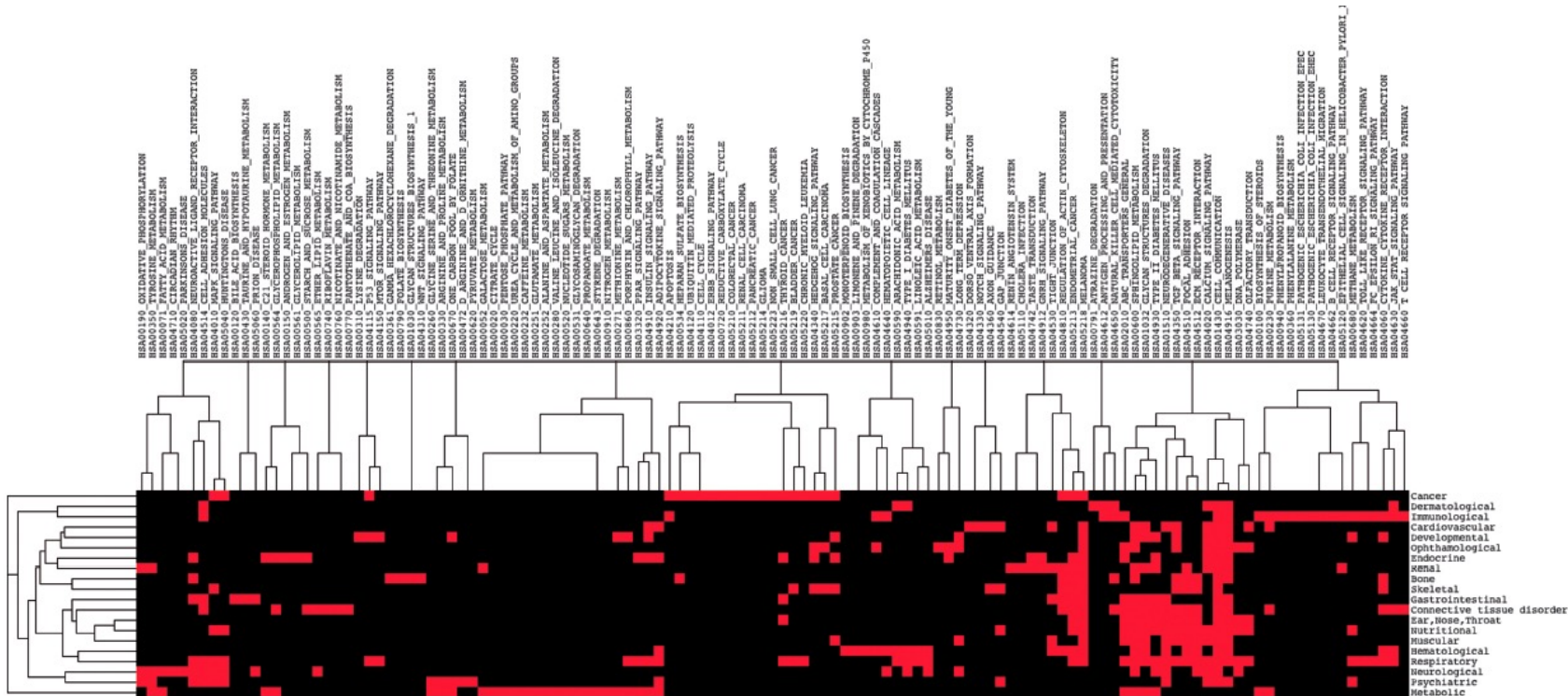
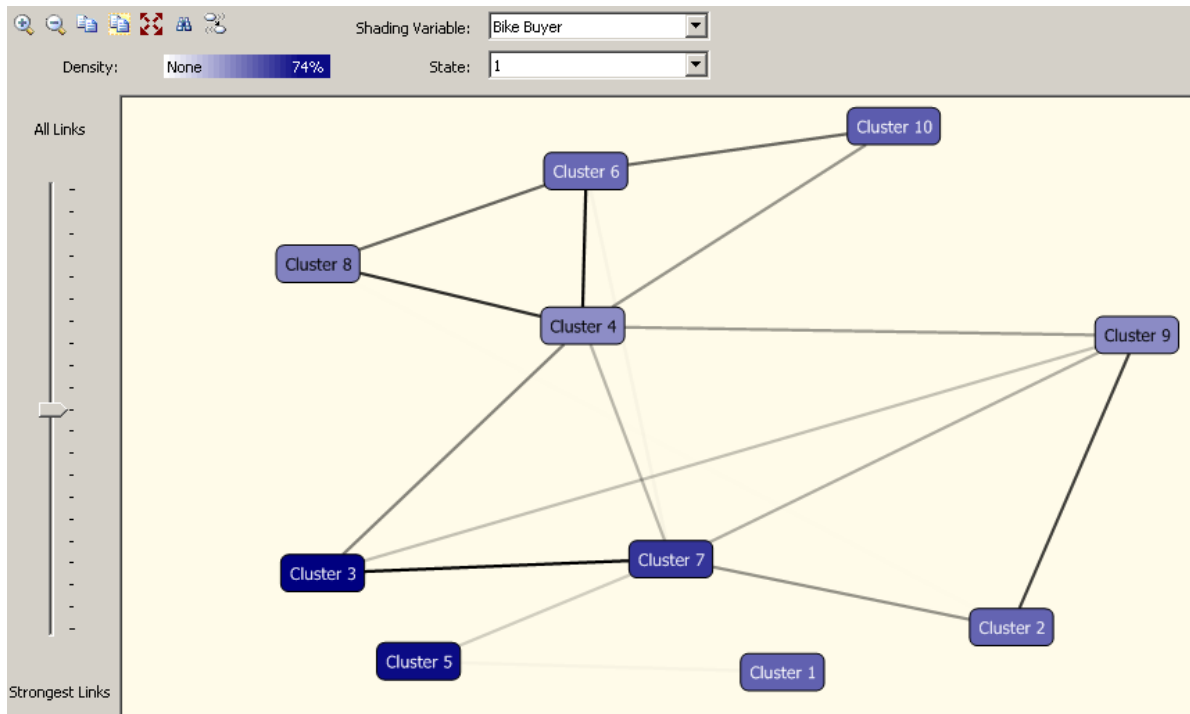


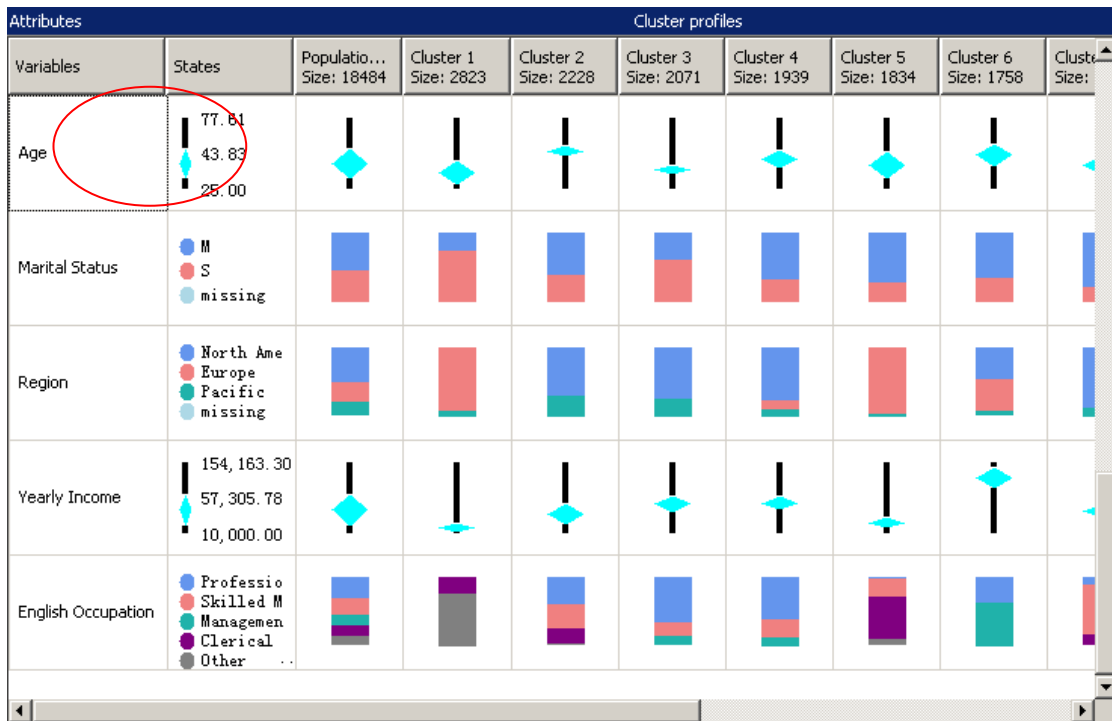
Figure 6. Predicted associations between disease classes and pathways. Each red entry represents a predicted association between 20 disease classes and 200 KEGG pathways.

7. 聚类结果可视化



- The lines between the clusters represent “closeness”
- The color of the cluster represents the frequency of the variable and state

其他：聚类结果可视化



- The distribution of an attribute's states for each cluster
- A colored bar: discrete attribute
- A diamond chart: continuous attribute























7. 聚簇特征



| Characteristics for Cluster 3 | | |
|-------------------------------|-------------------|-------------|
| Variables | Values | Probability |
| Bike Buyer | 1 | |
| Number Children At Home | 0 | |
| Region | North America | |
| Age | 36.2 - 43.8 | |
| English Occupation | Professional | |
| Marital Status | S | |
| House Owner Flag | 1 | |
| Gender | F | |
| First Name | missing | |
| English Education | Bachelors | |
| Commute Distance | 0-1 Miles | |
| Yearly Income | 57305.8 - 79082.2 | |
| Number Cars Owned | 0 | |
| Gender | M | |
| Number Cars Owned | 1 | |
| English Education | Graduate Degree | |
| House Owner Flag | 0 | |
| Last Name | missing | |
| Marital Status | M | |
| Commute Distance | 2-5 Miles | |
| Total Children | 0 | |
| Region | Pacific | |

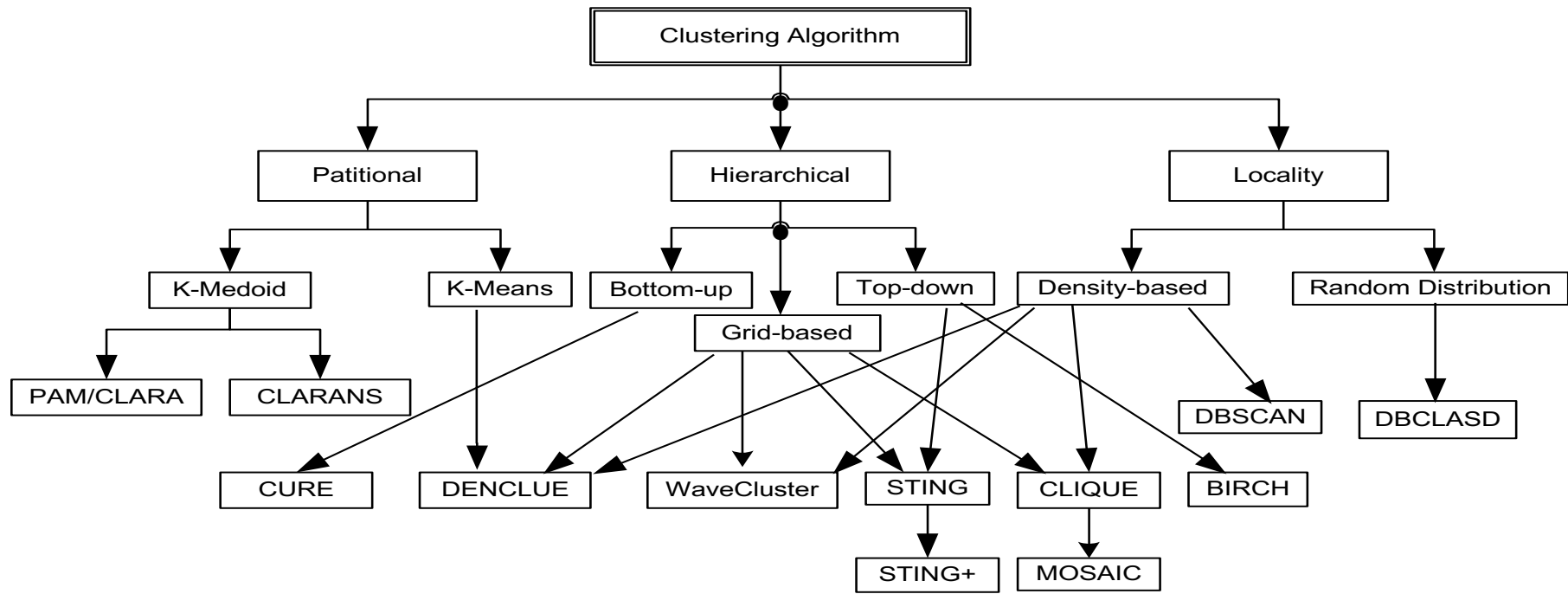
- The characteristics that make up a cluster in more detail

7. 区分簇的典型特征

| Variables | Values | Favors Cluster 3 | Favors Cluster 4 |
|-------------------------|-----------------|---|---|
| Number Children At Home | 0 |  | |
| Age | 27.9 - 43.5 |  | |
| Age | 43.5 - 95.0 | |  |
| Total Children | 0 |  | |
| English Education | Graduate Degree |  | |
| Commute Distance | 10+ Miles | |  |
| English Education | Partial College | |  |
| Bike Buyer | 0 | |  |
| Bike Buyer | 1 |  | |
| Total Children | 4 | |  |
| Number Cars Owned | 2 | |  |
| Number Cars Owned | 3 | |  |
| Number Children At Home | 1 | |  |
| Commute Distance | 0-1 Miles |  | |
| Total Children | 5 | |  |
| Total Children | 1 |  | |
| English Education | High School | |  |
| Number Cars Owned | 0 |  | |
| Number Children At Home | 4 | |  |
| House Owner Flag | 1 | |  |
| House Owner Flag | 0 |  | |
| Number Children At Home | 3 | |  |

- The characteristics that distinguish one cluster from another

常见的聚类算法之间的关系



- 可伸缩性
- 可以处理噪声
- 对输入次序不敏感
- 能够处理高维数据
- 可解释和可用性
- 处理不同类型属性的能力
- 可以发现任意形状的簇
- 用于决定输入参数的领域知识最小化

1. 聚类简介
2. 聚类性能评价指标
3. 空间和距离（相似度）度量
4. 基于划分的聚类 (K-Means, GMM)
5. 基于密度的聚类 (DBSCAN)
6. 层次聚类方法 (AGNES)
7. 聚类可视化