

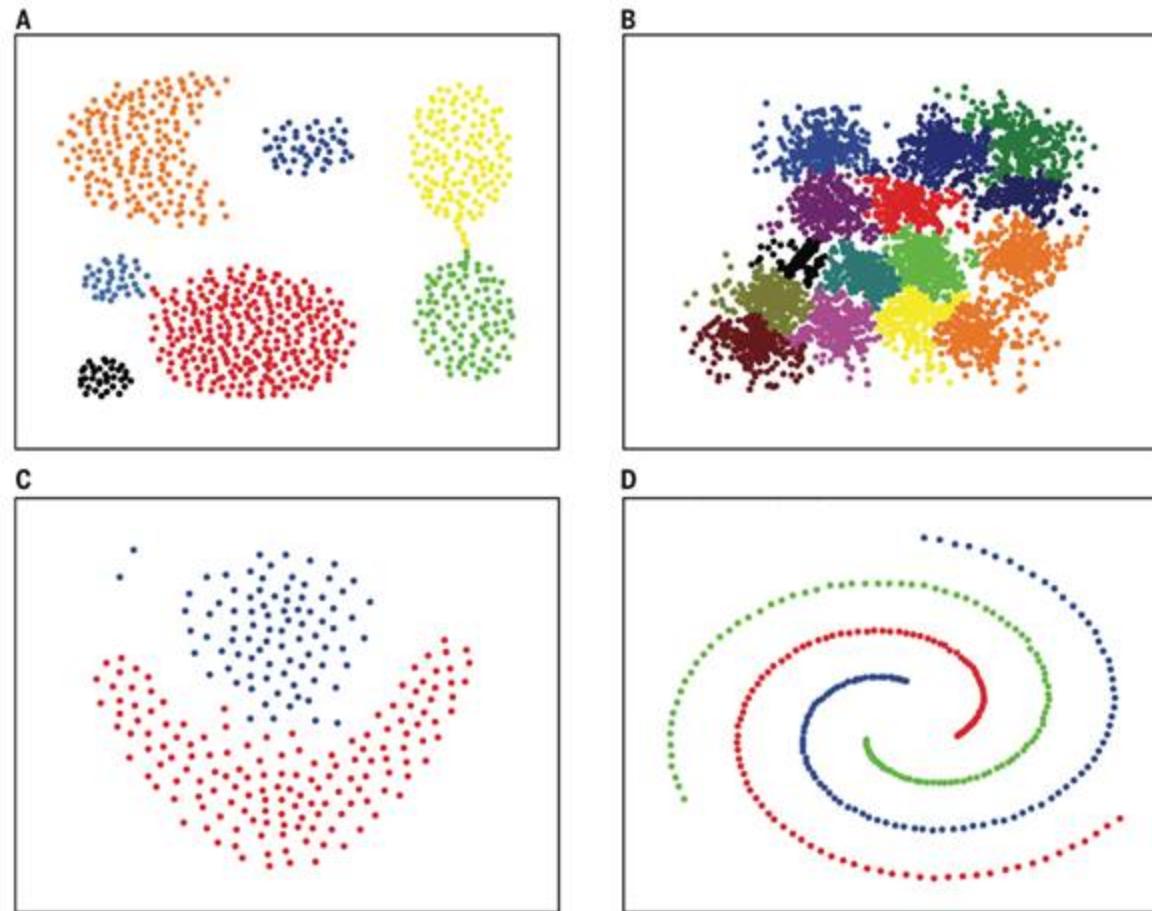
聚类分析-丁兆云

聚类分析-丁兆云

物以类聚、人以群分



聚类效果



Clustering by fast search and find of density peak. Alex Rodriguez, Alessandro Laio

学习目标

- ◆ 描述聚类分析的概念，区分聚类和分类
- ◆ 应用基本的聚类方法解决实际应用问题
- ◆ 了解聚类评估的任务有哪些

主要内容

- ◆ 1. 聚类分析概述
- ◆ 2. 基本聚类方法
- ◆ 3. 聚类评估

1. 聚类分析概述

- ◆ 什么是聚类?
 - 是把数据对象集合按照相似性划分成多个子集的过程。
 - 每个子集是一个簇 (cluster) , 使得簇中的对象彼此相似, 但与其他簇中的对象不相似。
- ◆ 聚类是无监督学习, 因为给的数据没有类标号信息

分类 vs. 聚类

- ◆ 分类
 - 有监督学习
 - 通过有标签样本学习分类器
- ◆ 聚类
 - 无监督学习
 - 通过观察学习，将数据分割成多个簇

聚类的应用

- ◆ 商业领域
 - 聚类分析被用来发现不同的客户群，并且通过购买模式刻画不同的客户群的特征。
- ◆ 电子商务
 - 聚类出具有相似浏览行为的客户，并分析客户的共同特征，可以更好的帮助电子商务的用户了解自己的客户，向客户提供更合适的服务。
- ◆ 舆情监控
 - 发现热点主题、话题、事件等
 - 发现未知异常

- ◆ **可扩展性**

- 许多聚类算法在小于几百个数据对象的小数据集合上工作得很好
- 而一个大规模数据库可能包含几百万个对象

- ◆ **处理不同数据类型的能力**

- 许多聚类算法专门用于数值类型的数据
- 而实际应用涉及不同的数据类型，如二元的、分类的、图像的等

- ◆ **发现任意形状的能力**

- 基于距离的聚类算法往往发现的是球形的聚类
- 而现实的聚类是任意形状的

- ◆ **用于决定输入参数的领域知识最小化**

- 聚类结果对于输入参数十分敏感
- 而参数很难决定，聚类的质量也很难控制

- ◆ **处理噪声数据的能力**

- 很多数据库都包含了孤立点，缺失或错误的数据
- 而一些聚类算法对于这样的数据敏感，可能导致低质量的聚类结果

- ◆ **对输入数据的顺序不敏感和增量聚类**

- 同一个数据集合，以不同的次序提交给同一个算法，应该产生相似的结果
- 能将新加入的数据合并到已有聚类中

- ◆ **高维度**

- 许多聚类算法擅长处理低维数据，可能只涉及两到三维
- 而数据库或者数据仓库可能包含若干维或属性

2. 基本聚类方法

- ◆ 2.1 划分方法
- ◆ 2.2 层次方法
- ◆ 2.3 基于密度的方法

2.1 划分方法

- ◆ 划分方法：将有 n 个对象的数据集D划分成 k 个簇，并且 $k \leq n$ ，满足如下的要求：
 - 每个簇至少包含一个对象
 - 每个对象属于且仅属于一个簇
- ◆ 基本思想
 - 首先创建一个初始 k 划分(k 为要构造的划分数)
 - 然后不断迭代地计算各个簇的聚类中心并依新的聚类中心调整聚类情况，直至收敛
- ◆ 目标
 - 同一个簇中的对象之间尽可能“接近”或相关
 - 不同簇中的对象之间尽可能“远离”或不同

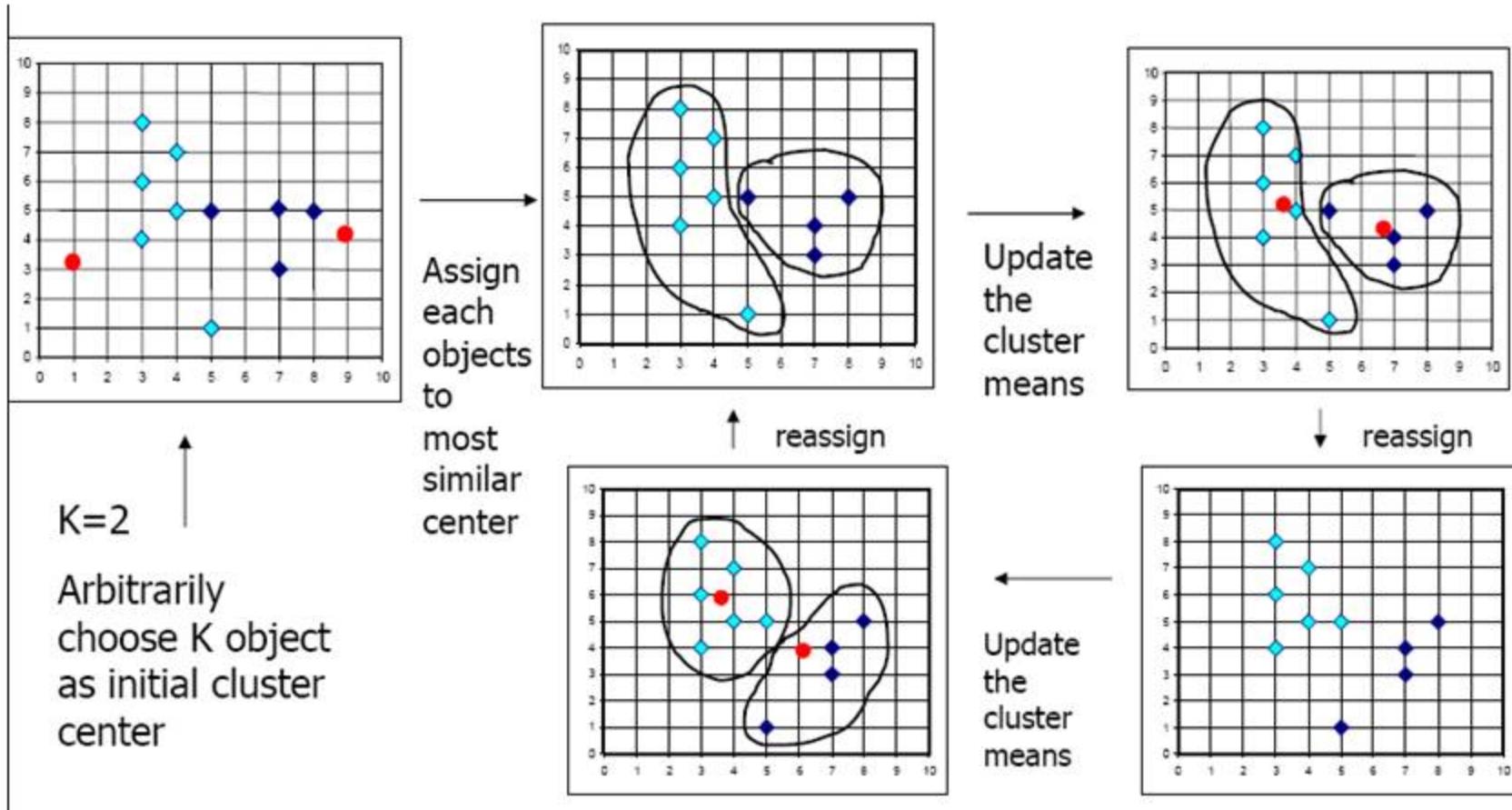
2.1划分方法

- ◆ 启发式方法 $E = \sum_{i=1}^k \sum_{p \in C_i} (d(p, c_i))^2$
 - k-均值(k-means)
 - 每个簇用该簇中对象的均值来表示
 - 基于质心的技术
 - k-中心点(k-medoids)
 - 每个簇用接近簇中心的一个对象来表示
 - 基于代表对象的技术
- ◆ 适用性
 - 这些启发式算法适合发现中小规模数据库中的球状聚类
 - 对于大规模数据库和处理任意形状的聚类，这些算法需要进一步扩展

2.1.1 KMeans

1. 从数据集中随机取K个样本作为初始的聚类中心
 $C = \{c_1, c_2, \dots, c_k\}$
2. 针对数据集中每个样本 x_i ，计算它到K个聚类中心的距离
并将其分到距离最小的聚类中心对应的类中
3. 针对每个类别 c_i ，重新计算其聚类中心 $c_i = \frac{1}{|c_i|} \sum_{x \in c_i} x$
4. 重复第2和第3步骤，直到 $\sum_{i=0}^n \min_{c_j \in C} (\|x_j - c_i\|^2)$ 小于特定的阈值

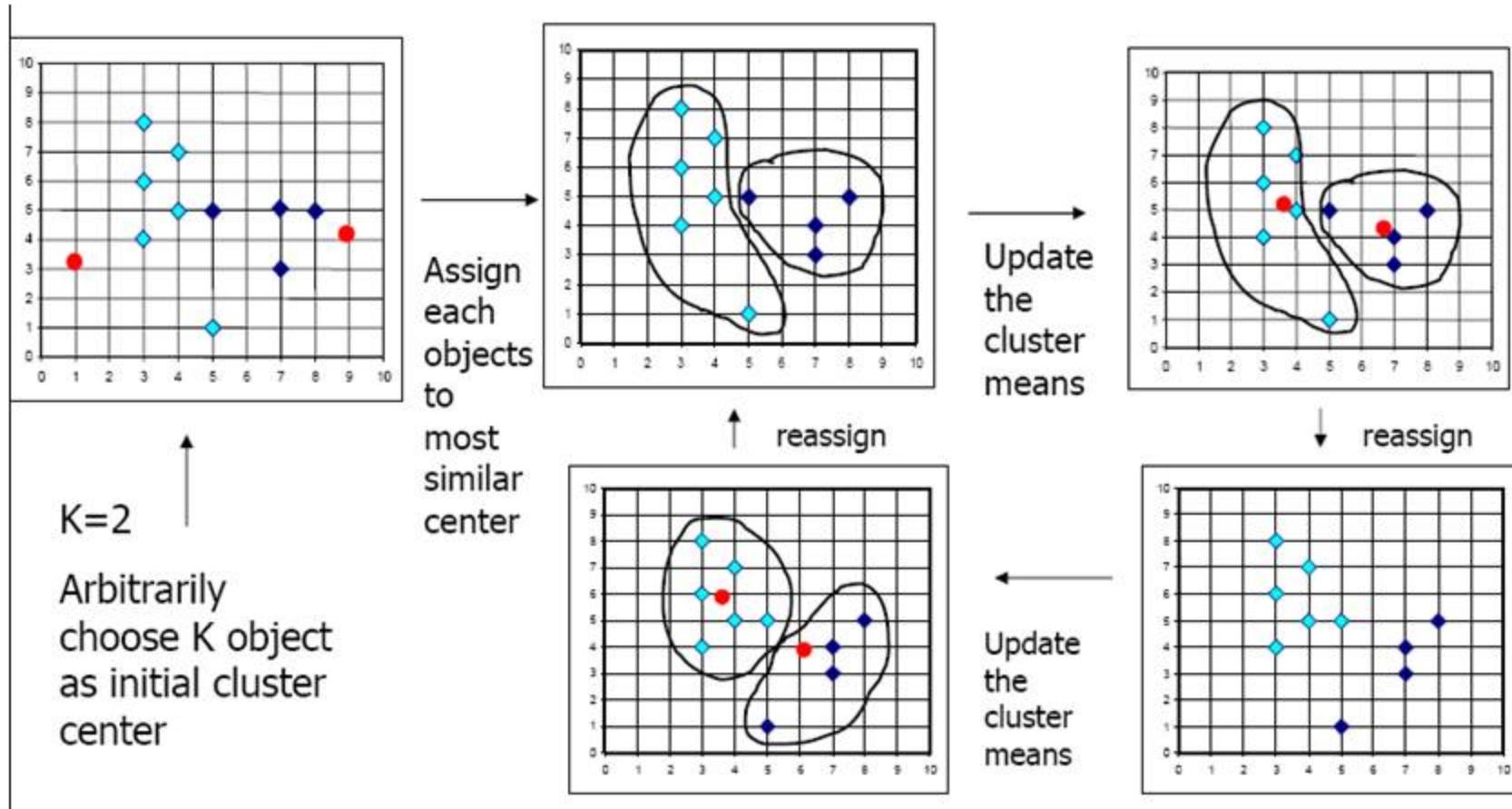
2.1.1 KMeans



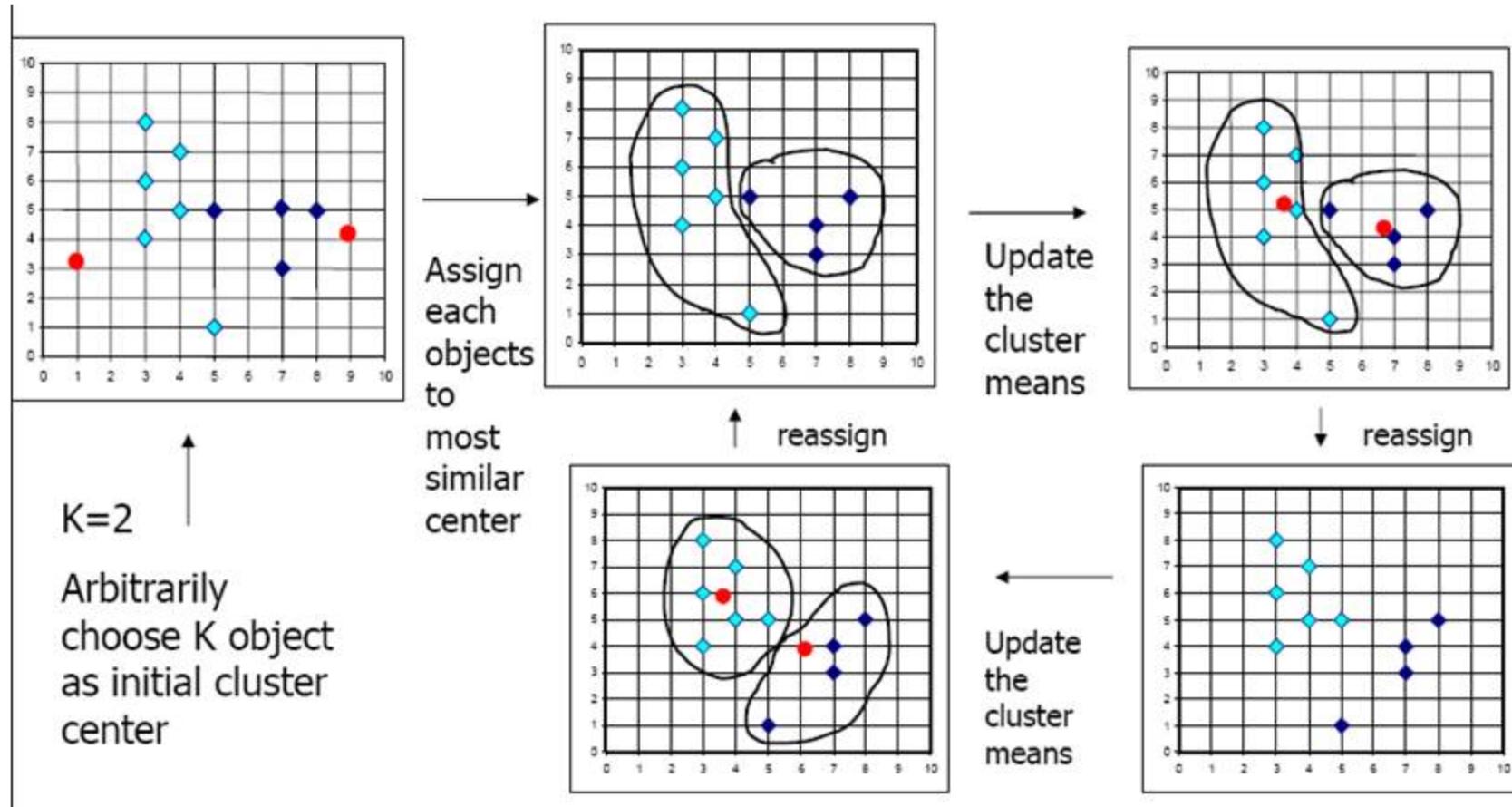
2.1.1 KMeans

1. 从数据集中随机取K个样本作为初始的聚类中心
 $C = \{c_1, c_2, \dots, c_k\}$
2. 针对数据集中每个样本 x_i ，计算它到K个聚类中心的距离
并将其分到距离最小的聚类中心对应的类中
3. 针对每个类别 c_i ，重新计算其聚类中心 $c_i = \frac{1}{|c_i|} \sum_{x \in c_i} x$
4. 重复第2和第3步骤，直到 $\sum_{i=0}^n \min_{c_j \in C} (\|x_j - c_i\|^2)$ 小于特定的阈值

2.1.1 KMeans

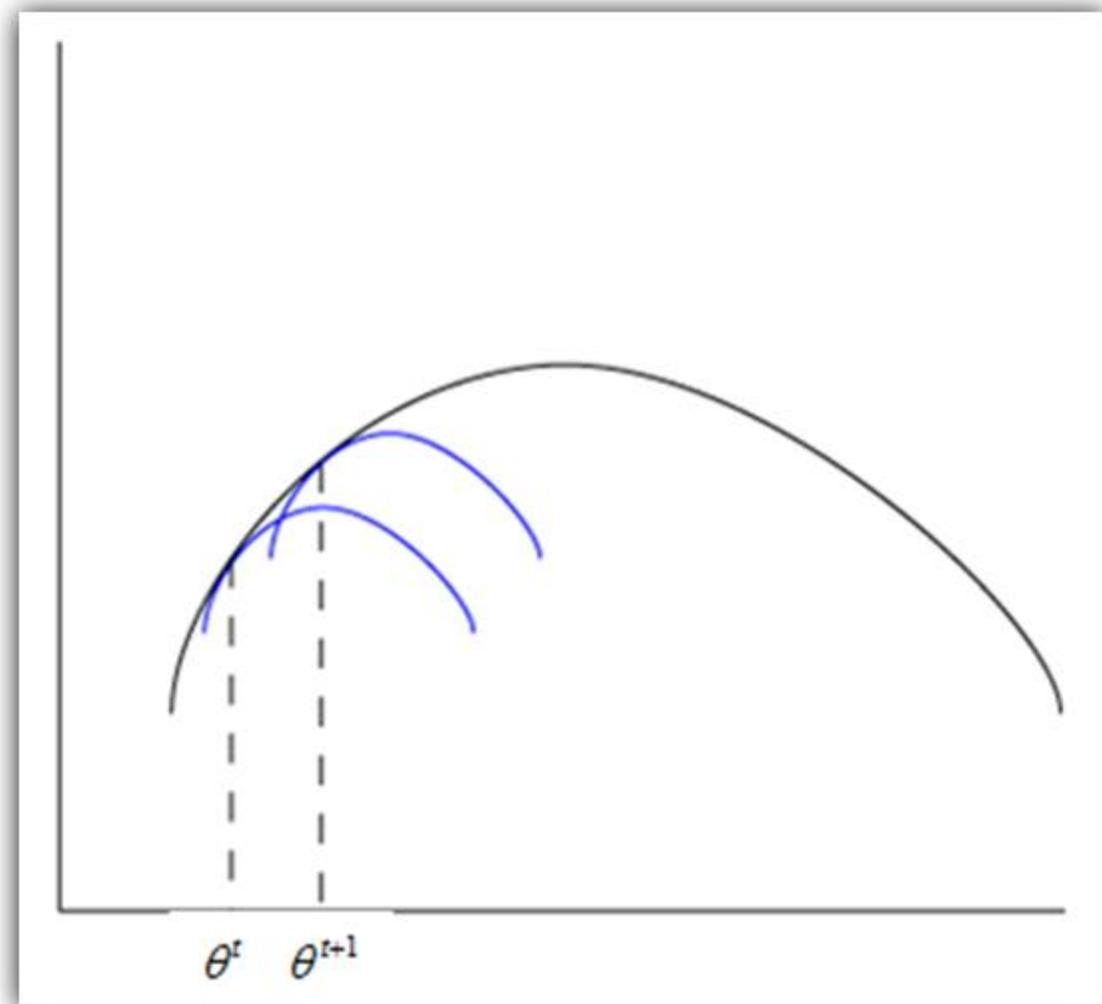
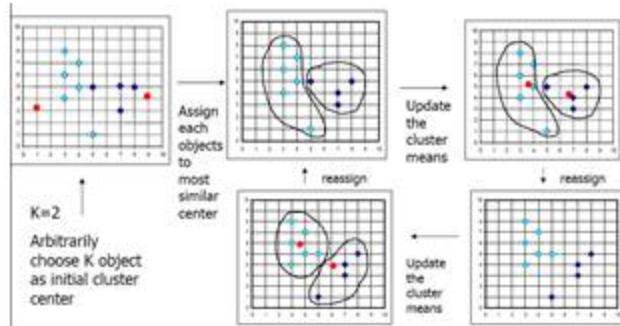


2.1.1 KMeans



Kmeans算法为启发式算法，遵循的寻优原则：每次聚类保证局部最优，随后调整聚类，利用局部最优聚类的上限来不断逼近全局最优

2.1.1 KMeans



Kmeans算法为启发式算法，遵循的寻优原则：**每次聚类保证局部最优，随后调整聚类，利用局部最优聚类的上限来不断逼近全局最优**

2.1.1 k-means计算示例

◆假设：给定如下要进行聚类的对象：

$\{2, 4, 10, 12, 3, 20, 30, 11, 25\}$, $k = 2$, 请使用k均值划分聚类

◆步骤如下：

m1	m2	K1	K2
2	4	{2,3}	{4,10,12,20,30,11,25}
2.5	16	{2,3,4}	{10,12,20,30,11,25}
3	18	{2,3,4,10}	{10,12,20,30,11,25}
4.75	19.6	{2,3,4,10,11,12}	{20,30,25}
7	25	{2,3,4,10,11,12}	{20,30,25}

2.1.1 An Example of K-Means Clustering

- 【例】假定我们对A、B、C、D四个样品分别测量两个变量和得到结果见表。

样品	变量	
	X_1	X_2
A	5	3
B	-1	1
C	1	-2
D	-3	-2

样品测量结果

试将以上的样品聚成两类。

最后形成的2类为

A (A,B)(C,D)

B (A,B,C)(D)

C (A)(B,C,D)

- ◆ 第一步：按要求取K=2，为了实施均值法聚类，我们将这些样品随意分成两类，比如 (A、B) 和 (C、D)，然后计算这两个聚类的中心坐标，见下表所示。
- ◆ 中心坐标是通过原始数据计算得来的。

样品	变量	
	X_1	X_2
A	5	3
B	-1	1
C	1	-2
D	-3	-2

聚类 中心坐标

聚类	中心坐标	
	X_1	X_2
(A、B)	2	2
(C、D)	-1	-2

提交

2.1.1 An Example of K-Means Clustering

样品	变量	
	X_1	X_2
A	5	3
B	-1	1
C	1	-2
D	-3	-2

- ◆ 第一步：按要求取K=2，为了实施均值法聚类，我们将这些样品随意分成两类，比如 (A、B) 和 (C、D)，然后计算这两个聚类的中心坐标，见下表所示。
- ◆ 中心坐标是通过原始数据计算得来的。

聚类	中心坐标	
	X_1	X_2
(A、B)	2	2
(C、D)	-1	-2

2.1.1 An Example of K-Means Clustering

第二步：计算某个样品到各类中心的欧氏平方距离，然后将该样品分配给最近的一类。对于样品有变动的类，重新计算它们的中心坐标，为下一步聚类做准备。先计算**A**到两个类的平方距离：

$$d^2(A, (AB)) = (5 - 2)^2 + (3 - 2)^2 = 10$$

$$d^2(A, (CD)) = (5 + 1)^2 + (3 + 2)^2 = 61$$

由于**A**到(**A**、**B**)的距离小于到(**C**、**D**)的距离，因此**A**不用重新分配。计算**B**到两类的平方距离：

$$d^2(B, (AB)) = (-1 - 2)^2 + (1 - 2)^2 = 10$$

$$d^2(B, (CD)) = (-1 + 1)^2 + (1 + 2)^2 = 9$$

2.1.1 An Example of K-Means Clustering

- 由于B到(A、B)的距离大于到(C、D)的距离，因此B要分配给(C、D)类，得到新的聚类是(A)和(B、C、D)。更新中心坐标如下表所示。

聚类	中心坐标	
	X_1	X_2
(A)	5	3
(B、C、D)	-1	-1

更新后的中心坐标

2.1.1 An Example of K-Means Clustering

第三步：再次检查每个样品，以决定是否需要重新分类。计算各样品到各中心的距离平方，结果见下表。

聚类	样品到中心的距离平方			
	A	B	C	D
(A)	0	40	41	89
(B、C、D)	52	4	5	5

- 到现在为止，每个样品都已经分配给距离中心最近的类，因此聚类过程就此结束。最终得到K=2的聚类结果是A独自成一类，B、C、D聚成一类。

2.1.1 k-means算法效率分析

- ◆ 算法的计算复杂度为 $O(nkt)$
- ◆ 其中
 - n 为数据集中对象的数目
 - k 为期望得到的簇的数目
 - t 为迭代的次数
- ◆ 在处理大数据库时也是相对有效的(可扩展性)

2.1.1 k-means优缺点

- ◆ 优点
 - 聚类时间快
 - 当结果簇是密集的，而簇与簇之间区别明显时，效果较好
 - 相对可扩展和有效，能对大数据集进行高效划分
- ◆ 缺点
 - 用户必须事先指定聚类簇的个数
 - 常常终止于局部最优
 - 只适用于数值属性聚类(计算均值有意义)
 - 对噪声和异常数据也很敏感
 - 不同的初始值，结果可能不同
 - 不适合发现非凸面形状的簇

2.1.1 k-means编程实践

- ◆ <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html#sklearn.cluster.KMeans>
 - 同学们可以尝试利用python读入本地iris数据集，来完成k-means聚类，分析其聚类效果

第6次课后作业

- ◆ 第六次课后作业-在educoder平台上完成作业
- ◆ <https://www.educoder.net/shixuns/85h6qo42/challenges>
- ◆ <https://www.educoder.net/shixuns/n38eh9ji/challenges>
- ◆ <https://www.educoder.net/shixuns/p87sflg2/challenges>
- ◆ <https://www.educoder.net/shixuns/8jrb659h/challenges>
- ◆ <https://www.educoder.net/shixuns/k6fp4saq/challenges>

提交作业截至时间：2020年3月08日

2.1.1 KMeans的问题

◆ 1. KMeans中初始簇规模估计正确

```
import numpy as np
import matplotlib.pyplot as plt

from sklearn.cluster import KMeans
from sklearn.datasets import make_blobs

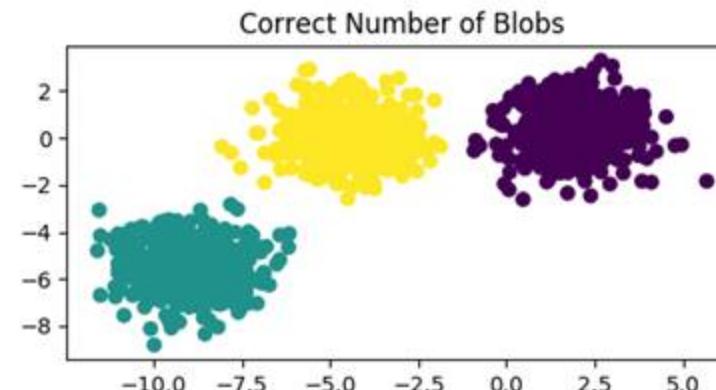
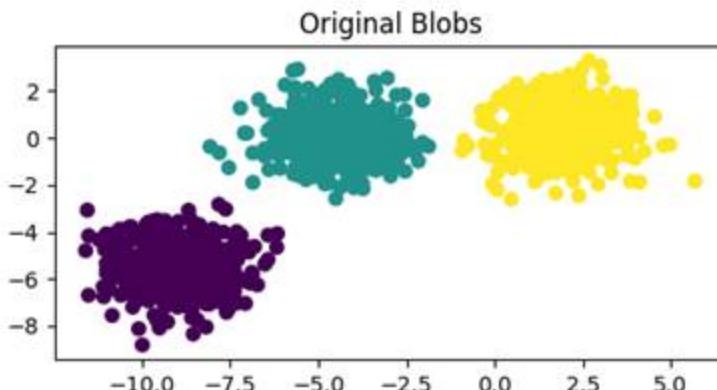
plt.figure(figsize=(12, 12))

n_samples = 1500
random_state = 170
X, y = make_blobs(n_samples=n_samples, random_state=random_state)

plt.subplot(321)
plt.scatter(X[:, 0], X[:, 1], c=y)
plt.title("Original Blobs")

y_pred = KMeans(n_clusters=3, random_state=random_state).fit_predict(X)

plt.subplot(322)
plt.scatter(X[:, 0], X[:, 1], c=y_pred)
plt.title("Correct Number of Blobs")
```



2.1.1 KMeans的问题

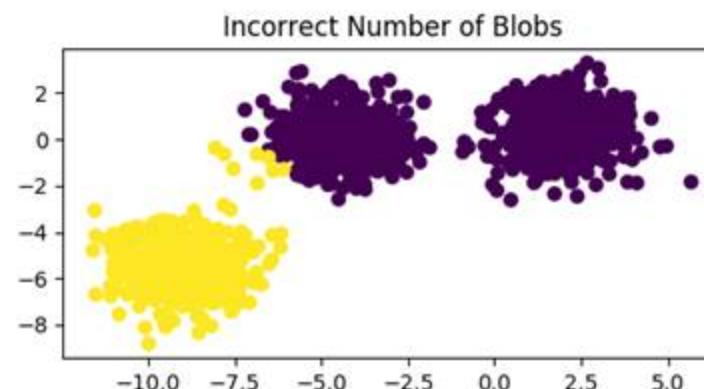
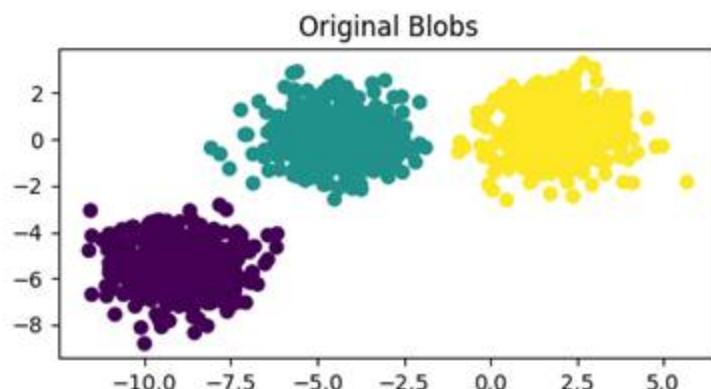
◆ 2. KMeans中初始簇规模估计错误

```
n_samples = 1500
random_state = 170
X, y = make_blobs(n_samples=n_samples, random_state=random_state)

plt.subplot(321)
plt.scatter(X[:, 0], X[:, 1], c=y)
plt.title("Original Blobs")

# Incorrect number of clusters
y_pred = KMeans(n_clusters=2, random_state=random_state).fit_predict(X)

plt.subplot(322)
plt.scatter(X[:, 0], X[:, 1], c=y_pred)
plt.title("Incorrect Number of Blobs")
```



2.1.1 KMeans 的问题

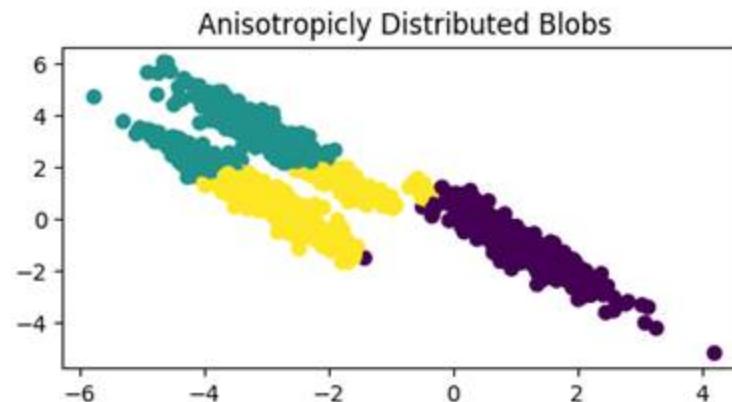
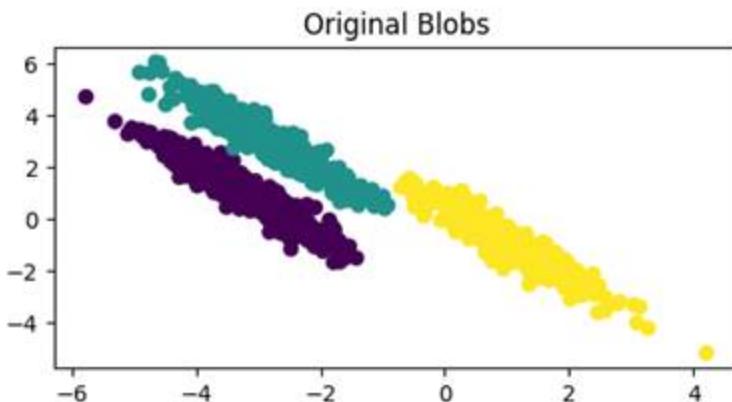
◆ 3. 数据分布形状对分簇结果的影响

```
X, y = make_blobs(n_samples=n_samples, random_state=random_state)
# Anisotropically distributed data
transformation = [[0.60834549, -0.63667341], [-0.40887718, 0.85253229]]
X_aniso = np.dot(X, transformation)

plt.subplot(321)
plt.scatter(X_aniso[:, 0], X_aniso[:, 1], c=y)
plt.title("Original Blobs")

y_pred = KMeans(n_clusters=3, random_state=random_state).fit_predict(X_aniso)

plt.subplot(322)
plt.scatter(X_aniso[:, 0], X_aniso[:, 1], c=y_pred)
plt.title("Anisotropically Distributed Blobs")
```



2.1.1 KMeans 的问题

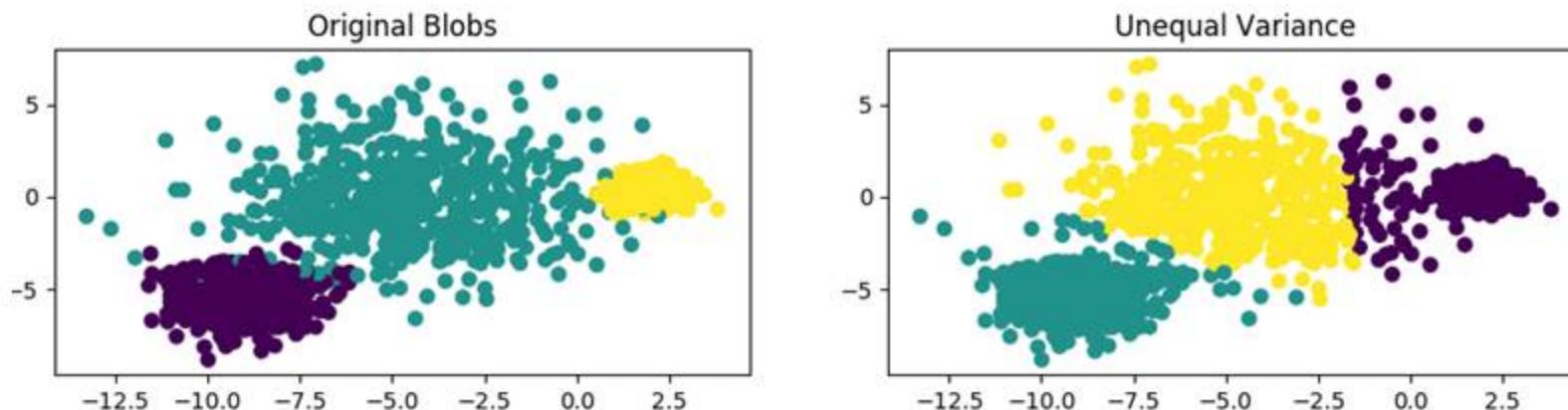
◆ 4. 数据分散程度对分簇结果的影响

```
# Different variance
X_varied, y_varied = make_blobs(n_samples=n_samples,
                                  cluster_std=[1.0, 2.5, 0.5],
                                  random_state=random_state)

plt.subplot(321)
plt.scatter(X_varied[:, 0], X_varied[:, 1], c=y_varied)
plt.title("Original Blobs")

y_pred = KMeans(n_clusters=3, random_state=random_state).fit_predict(X_varied)

plt.subplot(322)
plt.scatter(X_varied[:, 0], X_varied[:, 1], c=y_pred)
plt.title("Unequal Variance")
```



2.1.1 KMeans 的问题

◆ 5. 随机初始种子的影响

```
n_samples = 30
random_state = 50
X, y = make_blobs(n_samples=n_samples, random_state=random_state)

# Different variance
X_varied, y_varied = make_blobs(n_samples=n_samples,
                                  cluster_std=[1, 2, 1],
                                  random_state=random_state)
plt.subplot(321)
plt.scatter(X_varied[:, 0], X_varied[:, 1], c=y)
plt.title("Original Blobs")

y_pred = KMeans(init='random', n_clusters=3, n_init=1).fit_predict(X_varied)
plt.subplot(322)
plt.scatter(X_varied[:, 0], X_varied[:, 1], c=y_pred)
plt.title("Unequal Variance")

y_pred = KMeans(init='random', n_clusters=3, n_init=1).fit_predict(X_varied)
plt.subplot(323)
plt.scatter(X_varied[:, 0], X_varied[:, 1], c=y_pred)
plt.title("Unequal Variance")

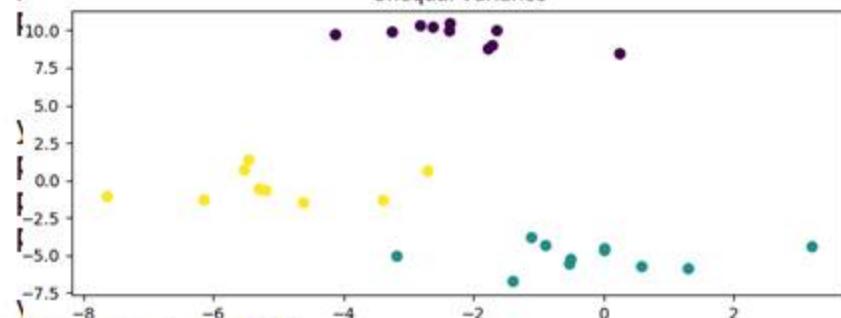
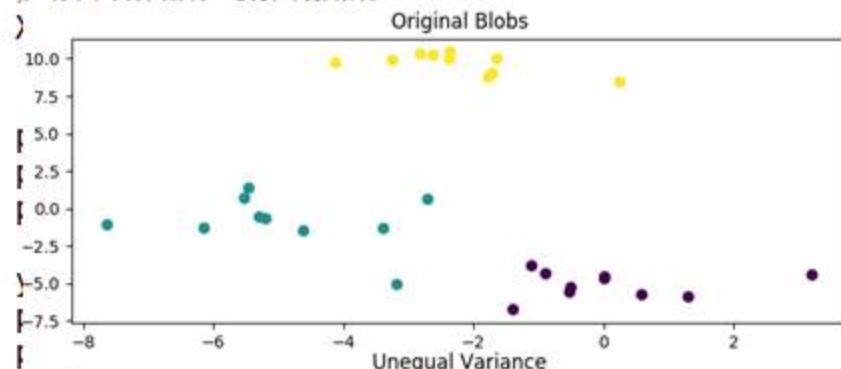
y_pred = KMeans(init='random', n_clusters=3, n_init=1).fit_predict(X_varied)
plt.subplot(324)
plt.scatter(X_varied[:, 0], X_varied[:, 1], c=y_pred)
plt.title("Unequal Variance")
```

2.1.1 KMeans 的问题

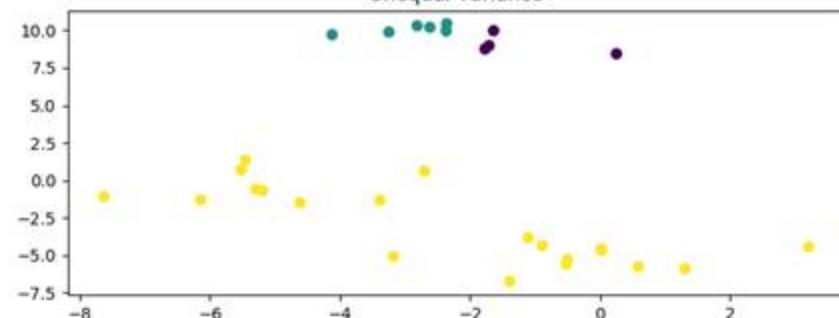
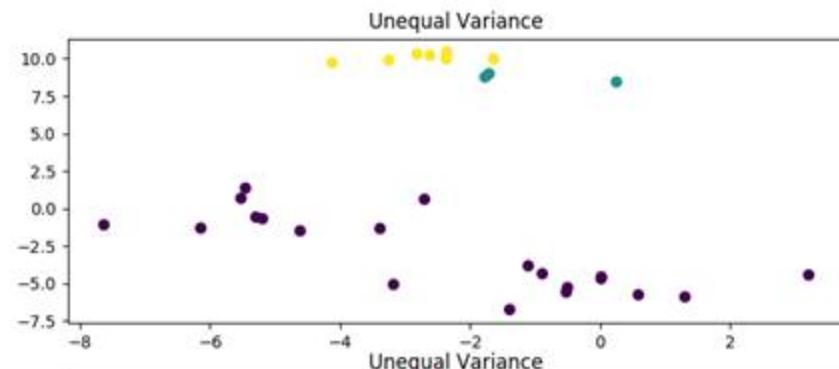
◆ 5. 随机初始种子的影响

```
n_samples = 30  
random_state = 50  
X, y = make_blobs(n_samples=n_samples, random_state=random_state)
```

Different variance



```
plt.subplot(324)  
plt.scatter(X_varied[:, 0], X_varied[:, 1], c=y_pred)  
plt.title("Unequal Variance")
```

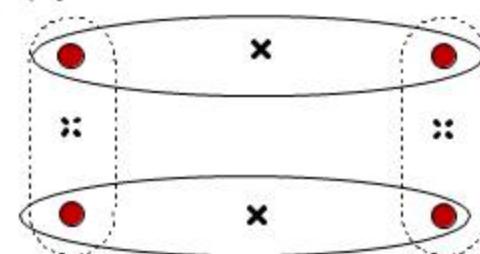


2.1.1 k-means优缺点

- ◆ 优点
 - 聚类时间快
 - 当结果簇是密集的，而簇与簇之间区别明显时，效果较好
 - 相对可扩展和有效，能对大数据集进行高效划分
- ◆ 缺点
 - 用户必须事先指定聚类簇的个数
 - 常常终止于局部最优
 - 只适用于数值属性聚类(计算均值有意义)
 - 对噪声和异常数据也很敏感
 - 不同的初始值，结果可能不同
 - 不适合发现非凸面形状的簇

2.1.2 k-modes算法-解决数据敏感的问题

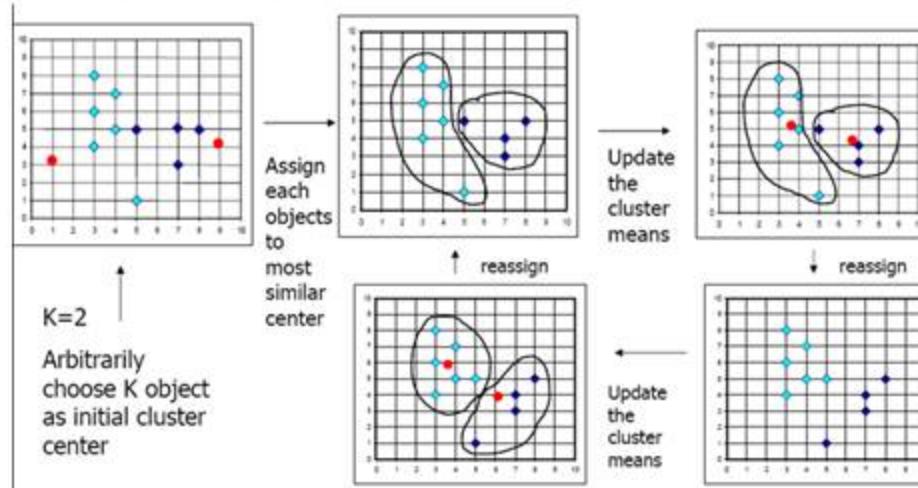
- Most of the variants of the *k-means* which differ in
 - Selection of the initial *k* means (初始k均值选择)
 - Dissimilarity calculations (相异度计算)
 - Strategies to calculate cluster means (计算均值方法)
- Handling categorical data: *k-modes*
 - Replacing means of clusters with modes
 - Using new dissimilarity measures to deal with categorical objects



2.1.3KMeans++算法-解决初始点选择问题

◆ 基本原理

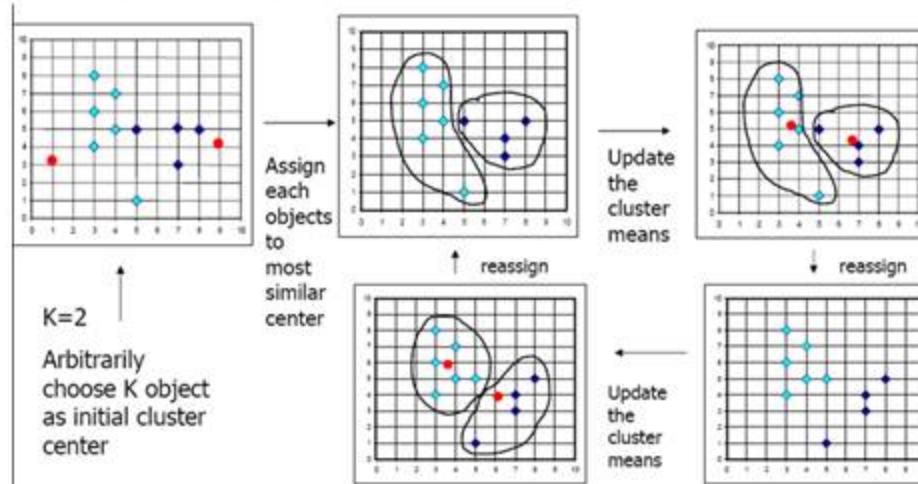
- 1.从输入的数据点集合中随机选择一个点作为第一个聚类中心；
- 2.对于数据集中的每一个点X，计算其与聚类中心的距离D(X)；
- 3.选择一个D(X)最大的点作为新的聚类中心；
- 4.重复2和3步直到K个聚类中心被选出；
- 5.利用K个初始聚类中心运行K-Means



2.1.3KMeans++算法-解决初始点选择问题

◆ 基本原理

- 1.从输入的数据点集合中随机选择一个点作为第一个聚类中心；
- 2.对于数据集中的每一个点X，计算其与聚类中心的距离D(X)；
- 3.选择一个D(X)最大的点作为新的聚类中心；
- 4.重复2和3步直到K个聚类中心被选出；
- 5.利用K个初始聚类中心运行K-Means



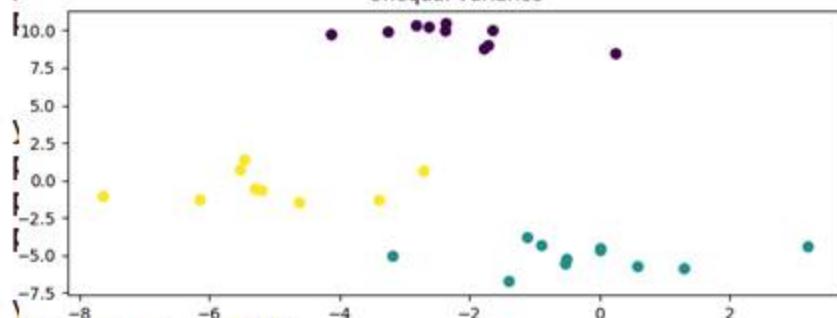
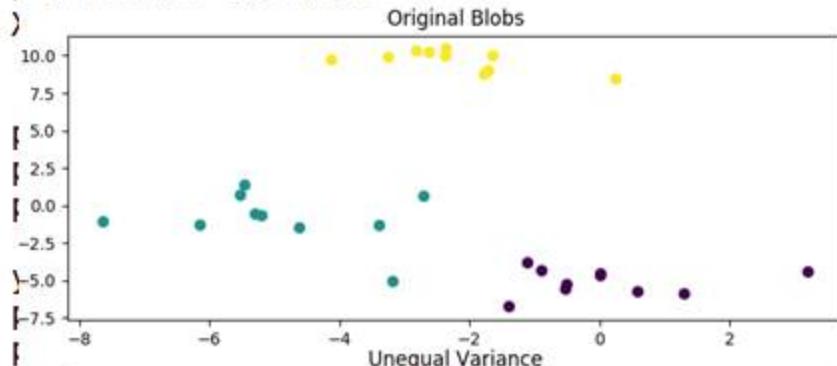
尽可能选择距离远的点来作为初始种子节点

2.1.3KMeans++算法-解决初始点选择问题

◆ 5. 随机初始种子的影响

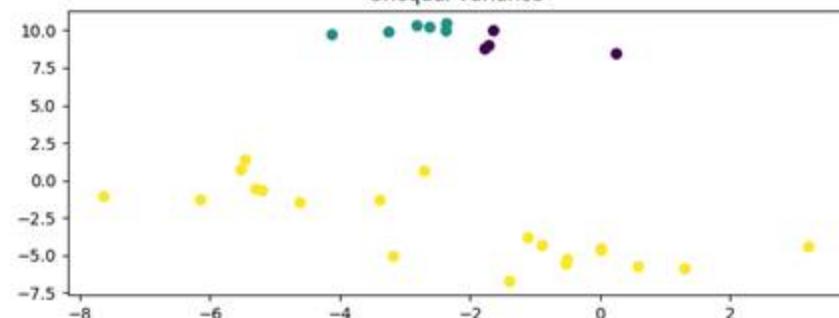
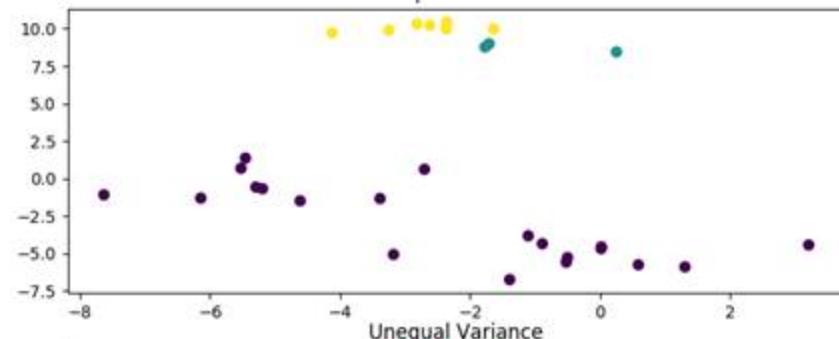
```
n_samples = 30  
random_state = 50  
X, y = make_blobs(n_samples=n_samples, random_state=random_state)
```

Different variance



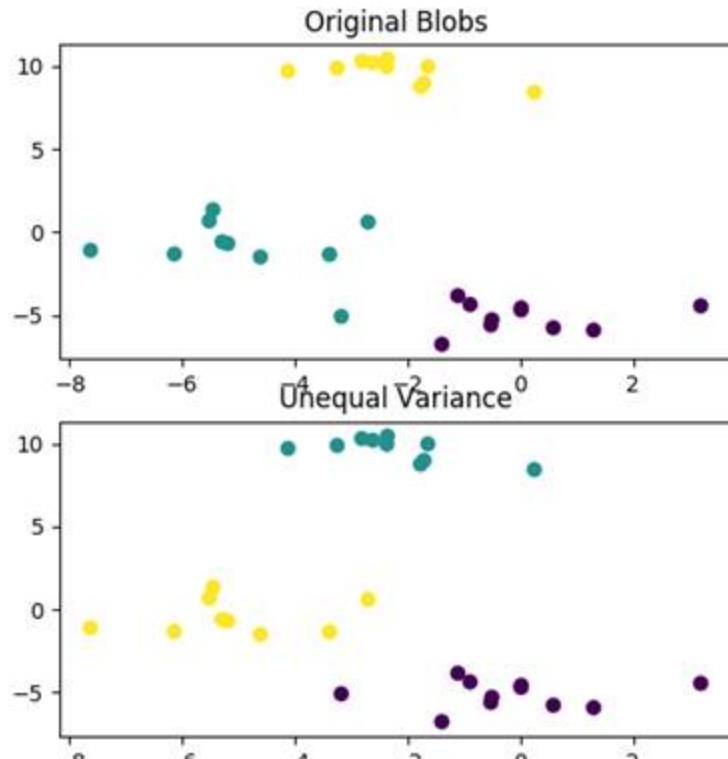
```
plt.subplot(324)  
plt.scatter(X_varied[:, 0], X_varied[:, 1], c=y_pred)  
plt.title("Unequal Variance")
```

Unequal Variance

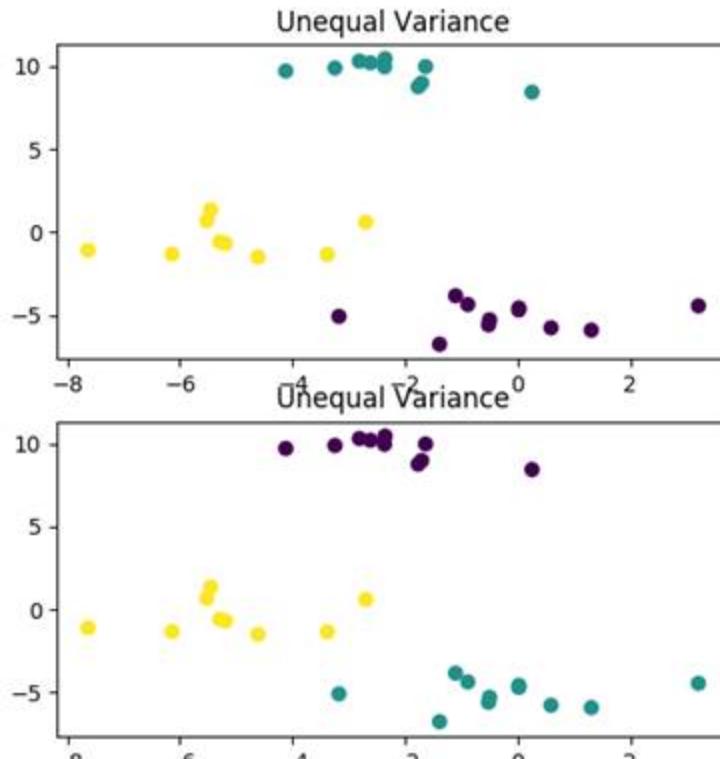


2.1.3KMeans++的效果

```
n_samples = 30  
random_state = 50  
X, y = make_blobs(n_samples=n_samples, random_state=random_state)
```



```
plt.title("Unequal Variance")
```



2.1.4 k-中心点方法解决对离群点敏感问题

一维空间的7个点 1 2 3 8 9 10 25 离群点

2.1.4 k-中心点方法解决对离群点敏感问题

一维空间的7个点 1 2 3 8 9 10 25 离群点

设置k=2，直觉应该划分为(1, 2, 3) (8, 9, 10, 25) 两个类别

2.1.4 k-中心点方法解决对离群点敏感问题

一维空间的7个点 1 2 3 8 9 10 25 离群点

设置k=2，直觉应该划分为(1, 2, 3) (8, 9, 10, 25) 两个类别

划分方法聚类质量评价
准则：最小化E值

$$E = \sum_{i=1}^k \sum_{p \in C_i} (d(p, c_i))^2$$

一维空间的7个点1 2 3 8 9 10 25，根据划分方法聚类质量评价准则，设置k为2时，(1,2,3)(8,9,10,25)聚类的E值为：[填空1]，(1,2,3,8)(9,10,25)的E值为：[填空2]

划分方法聚类质量评价
准则：最小化E值

$$E = \sum_{i=1}^k \sum_{p \in C_i} (d(p, c_i))^2$$

正常使用填空题需3.0以上版本雨课堂

作答

一维空间的7个点1 2 3 8 9 10 25，根据划分方法聚类质量评价准则，设置k为2时，将划分为如下哪2个类别

- A (1,2,3)(8,9,10,25)
- B (1,2,3,8)(9,10,25)

划分方法聚类质量评价
准则：最小化E值

$$E = \sum_{i=1}^k \sum_{p \in C_i} (d(p, c_i))^2$$

提交

2.1.4 k-中心点方法解决对离群点敏感问题

一维空间的7个点 1 2 3 8 9 10 25 离群点

设置k=2，直觉应该划分为(1, 2, 3) (8, 9, 10, 25) 两个类别

划分方法聚类质量评价
准则：最小化E值

$$E = \sum_{i=1}^k \sum_{p \in C_i} (d(p, c_i))^2$$

(1, 2, 3) (8, 9, 10, 25) 聚类的E值为：

$$(1 - 2)^2 + (2 - 2)^2 + (3 - 2)^2 + (8 - 13)^2 + (9 - 13)^2 + (10 - 13)^2 + (25 - 13)^2 = 196$$

(1, 2, 3, 8) (9, 10, 25) 的E值为

$$\begin{aligned}(1 - 3.5)^2 + (2 - 3.5)^2 + (3 - 3.5)^2 + (8 - 3.5)^2 + (9 - 14.67)^2 + (10 - 14.67)^2 \\ + (25 - 14.67)^2 = 189.67\end{aligned}$$

2.1.4 k-中心点方法解决对离群点敏感问题

一维空间的7个点 1 2 3 8 9 10 25 离群点

设置k=2，直觉应该划分为(1, 2, 3) (8, 9, 10, 25) 两个类别

划分方法聚类质量评价
准则：最小化E值

$$E = \sum_{i=1}^k \sum_{p \in C_i} (d(p, c_i))^2$$

(1, 2, 3) (8, 9, 10, 25) 聚类的E值为：

$$(1 - 2)^2 + (2 - 2)^2 + (3 - 2)^2 + (8 - 13)^2 + (9 - 13)^2 + (10 - 13)^2 + (25 - 13)^2 = 196$$

(1, 2, 3, 8) (9, 10, 25) 的E值为

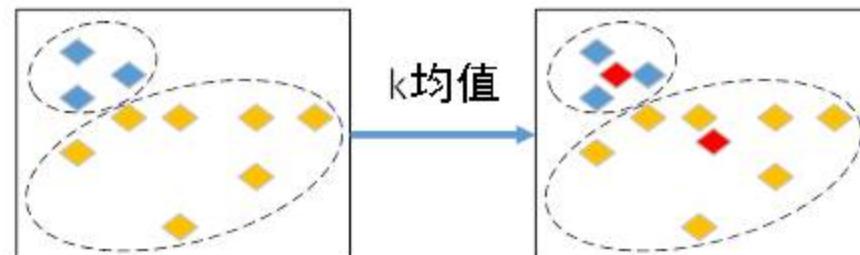
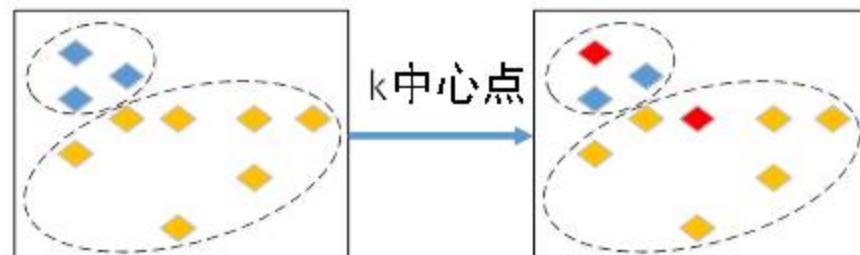
$$\begin{aligned}(1 - 3.5)^2 + (2 - 3.5)^2 + (3 - 3.5)^2 + (8 - 3.5)^2 + (9 - 14.67)^2 + (10 - 14.67)^2 \\ + (25 - 14.67)^2 = 189.67\end{aligned}$$

2.1.4 k-中心点

◆ k-中心点

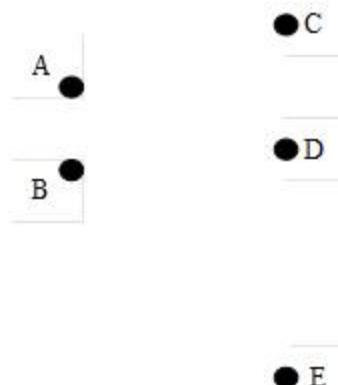
- 选用簇中位置**最中心的实际对象**即中心点作为参照点
- 基于最小化所有对象与其参照点之间的相异度之和的原则来划分(使用绝对误差标准)

$$E = \sum_{i=1}^k \sum_{p \in C_i} |p - o_i|$$



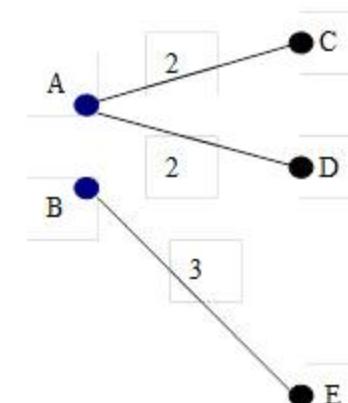
2.1.4k-中心点PAM算法示例

- 假如空间中的五个点 {A、B、C、D、E} 如图所示，各点之间的距离关系如表所示，根据所给的数据对其运行 PAM 算法实现划分聚类（设 $k=2$ ）。样本点间距离如下表所示：



样本点	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0

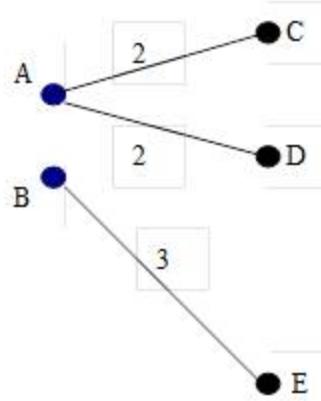
样本点



起始中心点

2.1.4 k-中心点PAM算法示例

- 第一步 建立阶段：假如从5个对象中随机抽取的2个中心点为{A, B},则样本被划分为{A、 C、 D}和{B、 E}
- 第二步 交换阶段：假定中心点A、 B分别被非中心点C、 D、 E替换，根据PAM算法需要计算下列代价 TC_{AC} 、 TC_{AD} 、 TC_{AE} 、 TC_{BC} 、 TC_{BD} 、 TC_{BE} ，即验证中心点变换后，代价如何变化
 - 以 TC_{AC} 为例说明计算过程



- a) 当A被C替换以后，A不再是一个中心点，因为A离B比A离C近，A被分配到B中心点代表的簇， $C_{AAC}=d(A,B)-d(A,A)=1$
- b) B是一个中心点，当A被C替换以后，B不受影响， $C_{BAC}=0$

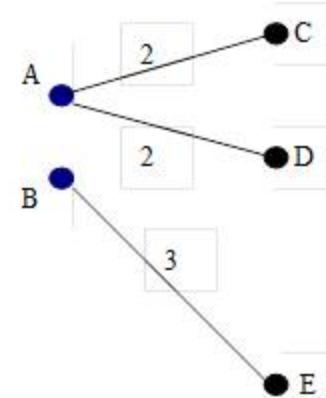
样本点	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0

2.1.4 k-中心点PAM算法示例

- a) $C_{AAC} = d(A, B) - d(A, A) = 1$
- b) $C_{BAC} = 0$
- c) C原先属于A中心点所在的簇，当A被C替换以后，C是新中心点，
 $C_{cAC} = d(C, C) - d(C, A) = 0 - 2 = -2$
- d) D原先属于A中心点所在的簇，当A被C替换以后，离D最近的中心点是C，
 $C_{DAC} = d(D, C) - d(D, A) = 1 - 2 = -1$
- e) E原先属于B中心点所在的簇，当A被C替换以后，离E最近的中心仍然是B，根据PAM算法代价函数的第三种情况

因此， $T C_{AAC} = C_{AAC} + C_{BAC} + C_{cAC} + C_{DAC} + C_{EAC} = 1 + 0 - 2 - 1 + 0 = -2$

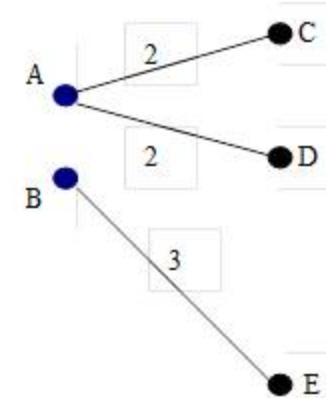
代价函数变小，说明C点替换以前的中心点A，聚类更加紧凑，那么C作为中心点更加合理



样本点	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0

2.1.4 k-中心点PAM算法示例

- a) $C_{AAC} = d(A, B) - d(A, A) = 1$
- b) $C_{BAC} = 0$
- c) C原先属于A中心点所在的簇，当A被C替换以后，C是新中心点，
 $C_{cAC} = d(C, C) - d(C, A) = 0 - 2 = -2$
- d) D原先属于A中心点所在的簇，当A被C替换以后，离D最近的中心点是C，
 $C_{DAC} = d(D, C) - d(D, A) = 1 - 2 = -1$
- e) E原先属于B中心点所在的簇，当A被C替换以后，离E最近的中心仍然是B，根据PAM算法代价函数的第三种情况



样本点	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0

□ 因此， $TC_{AAC} = C_{AAC} + C_{BAC} + C_{cAC} + C_{DAC} + C_{EAC} = 1 + 0 - 2 - 1 + 0 = -2$

代价函数变小，说明C点替换以前的中心点A，聚类更加紧凑，那么C作为中心点更加合理

同理计算 TC_{AC} 、 TC_{AD} 、 TC_{AE} 、 TC_{BC} 、 TC_{BD} 、 TC_{BE} ，代价函数最小值作为新中

2.1.4 k-中心点方法

◆ 基本思想

- 首先为每个簇随意选择一个代表对象，剩余的对象根据其与代表对象的**距离**分配给最近的一个簇
- 然后迭代地用非代表对象来替代代表对象，以改进聚类的质量(**找更好的代表对象**)
- 聚类结果的质量用一个**代价函数**来估算，该函数评估了对象与其参照对象之间的平均相异度

2.1.4 PAM算法

- ◆ PAM算法(Partitioning Around Medoids)
- ◆ 最早提出的k-中心点算法之一
 - (1) 随机选择 k 个对象作为初始的代表对象
repeat
 - (2) 指派每个剩余的对象给离它最近的代表对象所代表的簇
 - (3) 随机地选择一个非代表对象 O_{random}
 - (4) 计算用 O_{random} 代替 O_j 的总代价 S
 - (5) 如果 $S < 0$, 则用 O_{random} 替换 O_j 形成新的 k 个代表对象的集合
until 1不发生变化

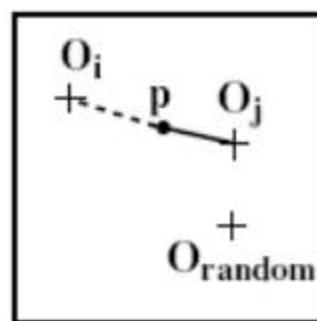
2.1.4 代价函数

◆ 代价函数 $TC_{jh} = \sum_{p=1}^n C_{pjh}$

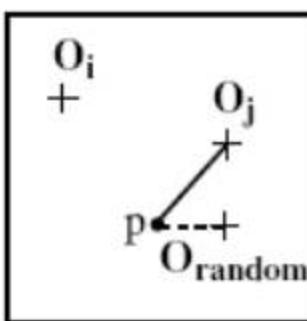
- 其中： n 是数据集中样本的个数； C_{pjh} 表示中心点 O_j 被非中心点 O_h 替代后，样本点 p 的代价。
- 问题：如何计算每个样本点 p 产生的代价 C_{pjh} ？

2.1.4 代价函数

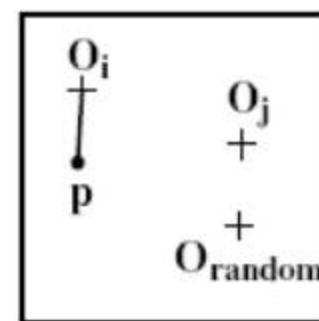
- 当非代表样本 O_{random} 替代代表样本 O_j 后，对于数据集中的每一个样本 p ，它所属的簇的类别将有以下四种可能的变化：



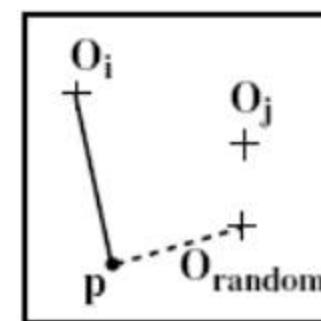
1. Reassigned to O_i
中心点变换后，
 p 另寻他人
• data object
+ cluster center
— before swapping
--- after swapping



2. Reassigned to O_{random}
中心点变换
后， p 加入新
中心点



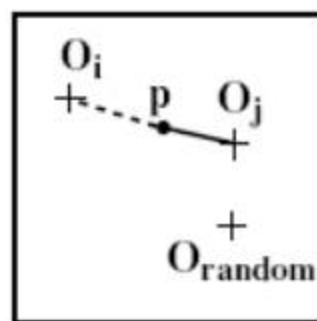
3. No change
中心点变换后，
 p 保持原位



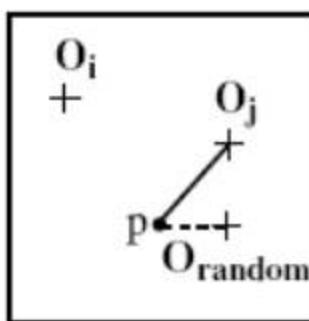
4. Reassigned to O_{random}
中心点变换
后， p 加入新
中心点

2.1.4 代价函数

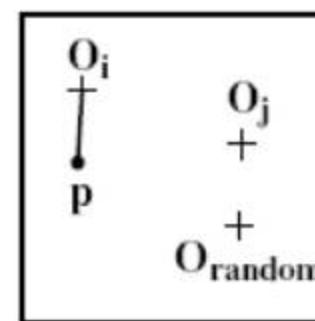
- 当非代表样本 O_{random} 替代代表样本 O_j 后，对于数据集中的每一个样本 p ，它所属的簇的类别将有以下四种可能的变化：



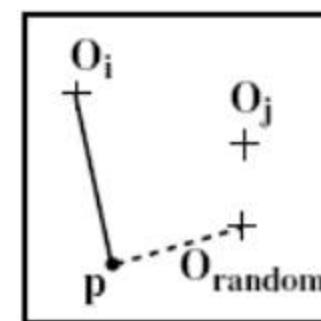
1. Reassigned to O_i
中心点变换后，
p另寻他人
• data object
+ cluster center
— before swapping
--- after swapping



2. Reassigned to O_{random}
中心点变换
后，p加入新
中心点
p最初属于被
替换的点



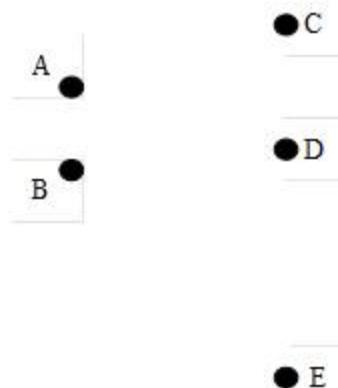
3. No change
中心点变换后，
p保持原位



4. Reassigned to O_{random}
中心点变换
后，p加入新
中心点
p最初不属于
被替换的点

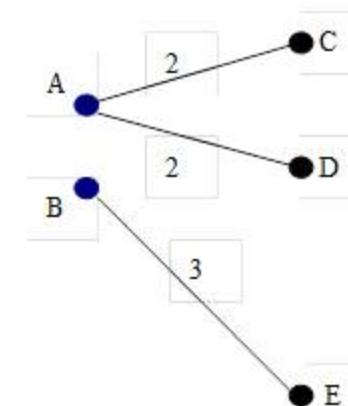
2.1.4k-中心点PAM算法示例

- 假如空间中的五个点 {A、B、C、D、E} 如图所示，各点之间的距离关系如表所示，根据所给的数据对其运行 PAM 算法实现划分聚类（设 $k=2$ ）。样本点间距离如下表所示：



样本点	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0

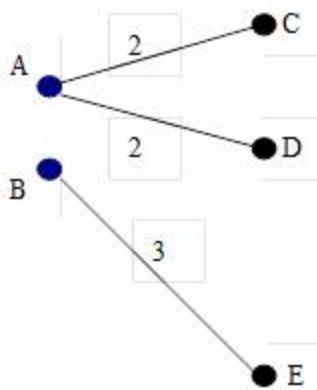
样本点



起始中心点

2.1.4 k-中心点PAM算法示例

- 第一步 建立阶段：假如从5个对象中随机抽取的2个中心点为{A, B},则样本被划分为{A、 C、 D}和{B、 E}
- 第二步 交换阶段：假定中心点A、 B分别被非中心点C、 D、 E替换，根据PAM算法需要计算下列代价 TC_{AC} 、 TC_{AD} 、 TC_{AE} 、 TC_{BC} 、 TC_{BD} 、 TC_{BE} ，即验证中心点变换后，
 - 以 TC_{AC} 为例说明计算过程

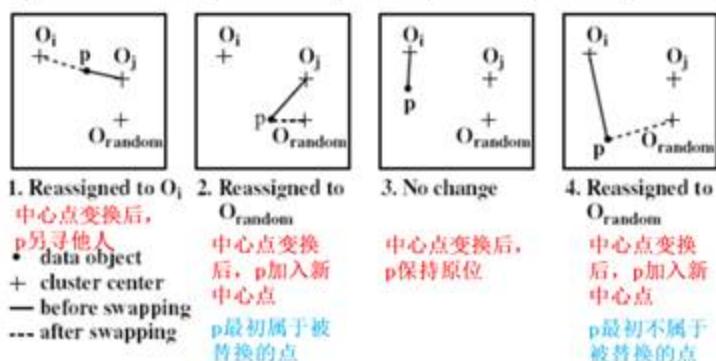


起始中心点

a) 当A被C替换以后，A不再是一个中心点，因为A离B比A离C近，A被分配到B中心点代表的簇， $C_{AAC}=d(A,B)-d(A,A)=1$

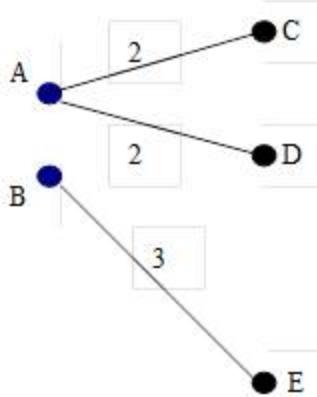
b) B是一个中心点，当A被C替换以后，B不受影响， $C_{BAC}=0$

	代价如何变化				
	B	C	D	E	
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0



2.1.4 k-中心点PAM算法示例

- 第一步 建立阶段：假如从5个对象中随机抽取的2个中心点为{A, B},则样本被划分为{A、 C、 D}和{B、 E}
- 第二步 交换阶段：假定中心点A、 B分别被非中心点C、 D、 E替换，根据PAM算法需要计算下列代价 TC_{AC} 、 TC_{AD} 、 TC_{AE} 、 TC_{BC} 、 TC_{BD} 、 TC_{BE} ，即验证中心点变换后，
 - 以 TC_{AC} 为例说明计算过程

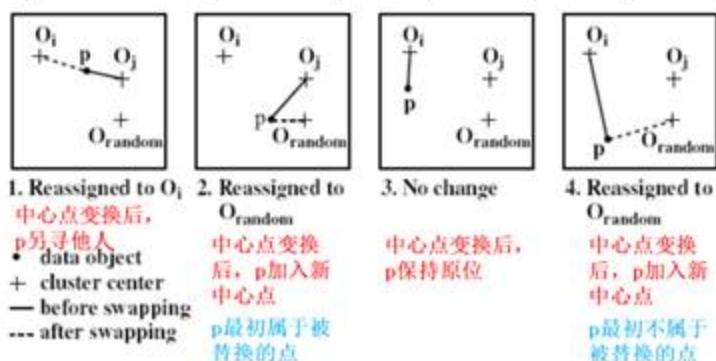


起始中心点

a) 当A被C替换以后，A不再是一个中心点，因为A离B比A离C近，A被分配到B中心点代表的簇， $C_{AAC}=d(A,B)-d(A,A)=1$

b) B是一个中心点，当A被C替换以后，B不受影响， $C_{BAC}=0$

	代价如何变化				B	C	D	E
A	0	1	2	2	3			
B	1	0	2	4	3			
C	2	2	0	1	5			
D	2	4	1	0	3			
E	3	3	5	3	0			



2.1.4 k-中心点PAM算法示例

a) $C_{AAC} = d(A, B) - d(A, A) = 1$

b) $C_{BAC} = 0$

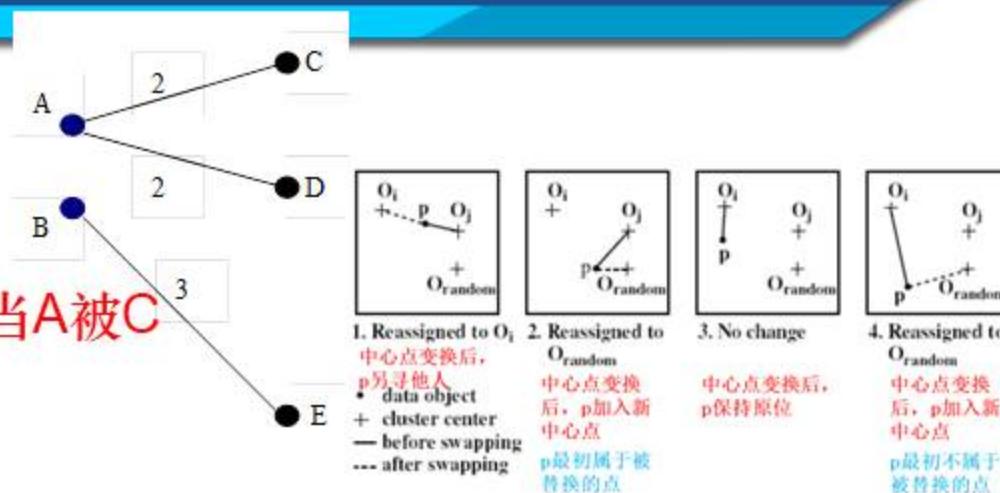
c) C原先属于A中心点所在的簇，当A被C替换以后，C是新中心点，
 $C_{cAC} = d(C, C) - d(C, A) = 0 - 2 = -2$

d) D原先属于A中心点所在的簇，当A被C替换以后，离D最近的中心点是C，
 $C_{DAC} = d(D, C) - d(D, A) = 1 - 2 = -1$

e) E原先属于B中心点所在的簇，当A被C替换以后，离E最近的中心仍然是B，根据PAM算法代价函数的第三种情况

$C_{EAC} = 0$

因此， $TC_{AC} = C_{AAC} + C_{BAC} + C_{cAC} + C_{DAC} + C_{EAC} = 1 + 0 - 2 - 1 + 0 = -2$ 。



样本点	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0

2.1.4 k-中心点PAM算法示例

a) $C_{AAC} = d(A, B) - d(A, A) = 1$

b) $C_{BAC} = 0$

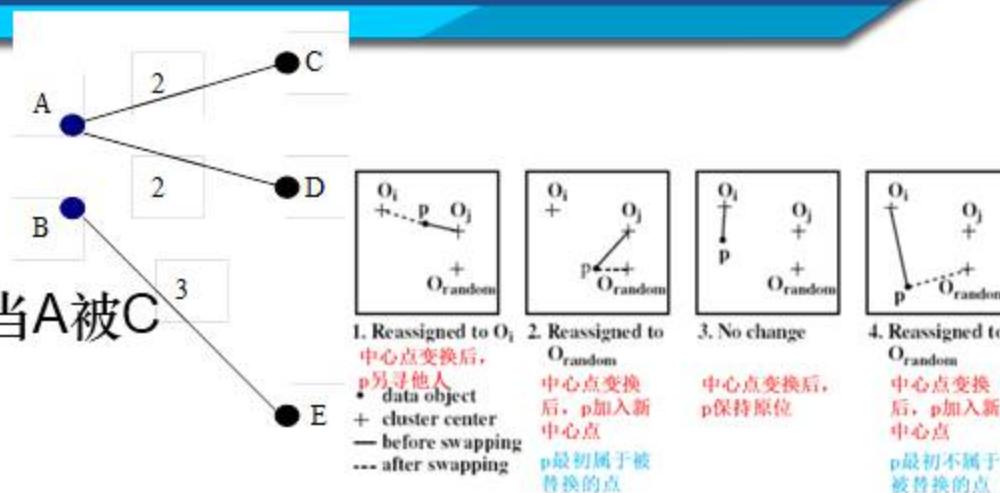
c) C原先属于A中心点所在的簇，当A被C替换以后，C是新中心点，
 $C_{cAC} = d(C, C) - d(C, A) = 0 - 2 = -2$

d) D原先属于A中心点所在的簇，当A被C替换以后，离D最近的中心点是C，
 $C_{DAC} = d(D, C) - d(D, A) = 1 - 2 = -1$

e) E原先属于B中心点所在的簇，当A被C替换以后，离E最近的中心仍然是B，根据PAM算法代价函数的第三种情况

$C_{EAC} = 0$

因此， $TC_{AC} = C_{AAC} + C_{BAC} + C_{cAC} + C_{DAC} + C_{EAC} = 1 + 0 - 2 - 1 + 0 = -2$ 。



样本点	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0

2.1.4 k-中心点PAM算法示例

a) $C_{AAC} = d(A, B) - d(A, A) = 1$

b) $C_{BAC} = 0$

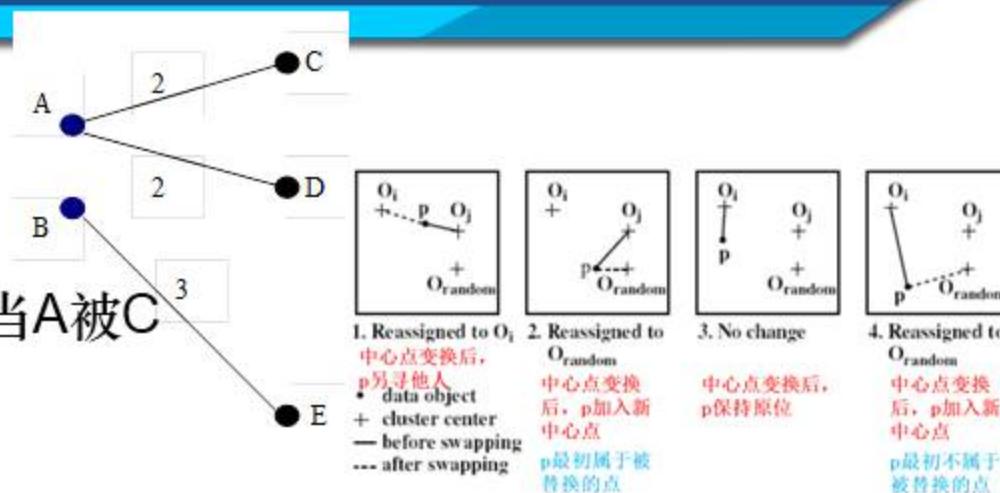
c) C原先属于A中心点所在的簇，当A被C替换以后，C是新中心点，
 $C_{cAC} = d(C, C) - d(C, A) = 0 - 2 = -2$

d) D原先属于A中心点所在的簇，当A被C替换以后，离D最近的中心点是C，
 $C_{DAC} = d(D, C) - d(D, A) = 1 - 2 = -1$

e) E原先属于B中心点所在的簇，当A被C替换以后，离E最近的中心仍然是B，根据PAM算法代价函数的第三种情况

$C_{EAC} = 0$

因此， $TC_{AC} = C_{AAC} + C_{BAC} + C_{cAC} + C_{DAC} + C_{EAC} = 1 + 0 - 2 - 1 + 0 = -2$ 。



样本点	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0

2.2.1k-中心点PAM算法示例

a) $C_{AAC} = d(A, B) - d(A, A) = 1$

b) $C_{BAC} = 0$

c) C原先属于A中心点所在的簇，当A被C替换以后，C是新中心点，
 $C_{cAC} = d(C, C) - d(C, A) = 0 - 2 = -2$

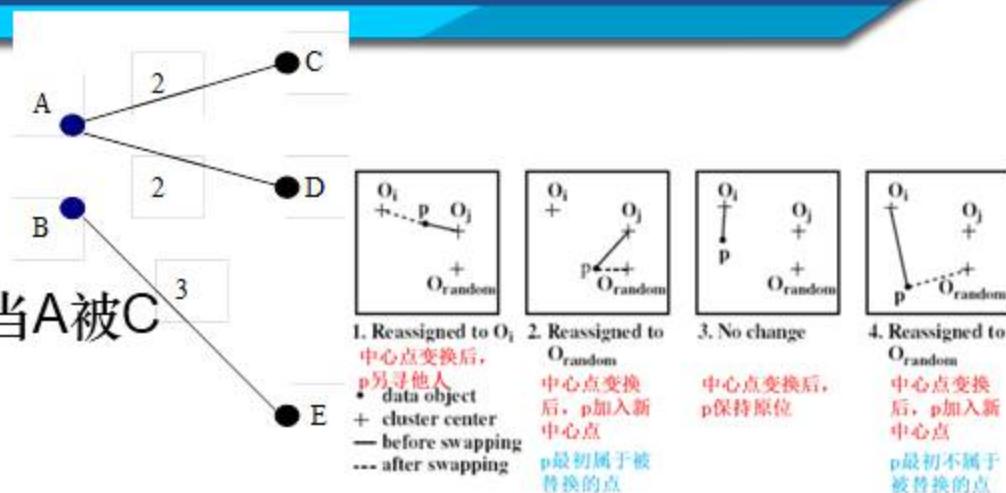
d) D原先属于A中心点所在的簇，当A被C替换以后，离D最近的中心点是C，
 $C_{DAC} = d(D, C) - d(D, A) = 1 - 2 = -1$

e) E原先属于B中心点所在的簇，当A被C替换以后，离E最近的中心仍然是B，根据PAM算法代价函数的第三种情况

$C_{EAC} = 0$

因此， $TC_{AC} = C_{AAC} + C_{BAC} + C_{cAC} + C_{DAC} + C_{EAC} = 1 + 0 - 2 -$

代价函数变小，说明C点替换以前的中心点A，聚类更加紧凑，那么C作为中心点更加合理



样本点	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0

2.2.1k-中心点PAM算法示例

a) $C_{AAC} = d(A, B) - d(A, A) = 1$

b) $C_{BAC} = 0$

c) C原先属于A中心点所在的簇，当A被C替换以后，C是新中心点，
 $C_{cAC} = d(C, C) - d(C, A) = 0 - 2 = -2$

d) D原先属于A中心点所在的簇，当A被C替换以后，离D最近的中心点是C，
 $C_{DAC} = d(D, C) - d(D, A) = 1 - 2 = -1$

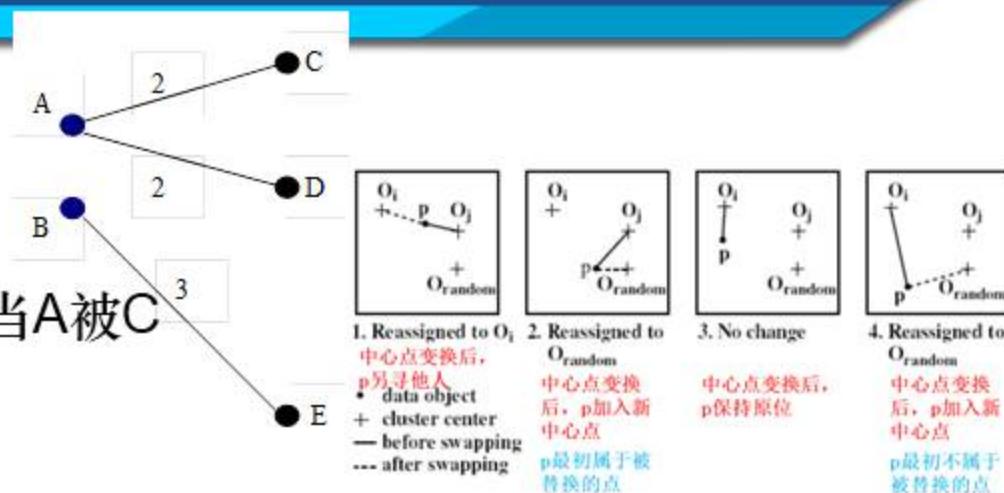
e) E原先属于B中心点所在的簇，当A被C替换以后，离E最近的中心仍然是B，根据PAM算法代价函数的第三种情况

$C_{EAC} = 0$

因此， $TC_{AC} = C_{AAC} + C_{BAC} + C_{cAC} + C_{DAC} + C_{EAC} = 1 + 0 - 2 -$

代价函数变小，说明C点替换以前的中心点A，聚类更加紧凑，那么C作为中心点更加合理

同理计算 TC_{AC} 、 TC_{AD} 、 TC_{AE} 、 TC_{BC} 、 TC_{BD} 、 TC_{BE} ，代价函数最小值作为新中

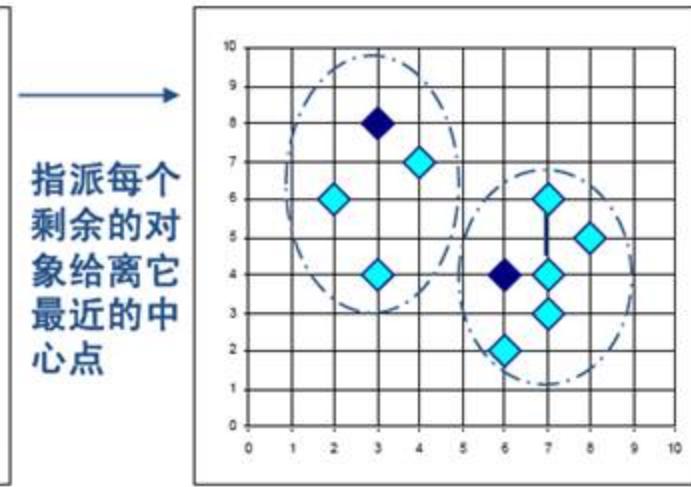
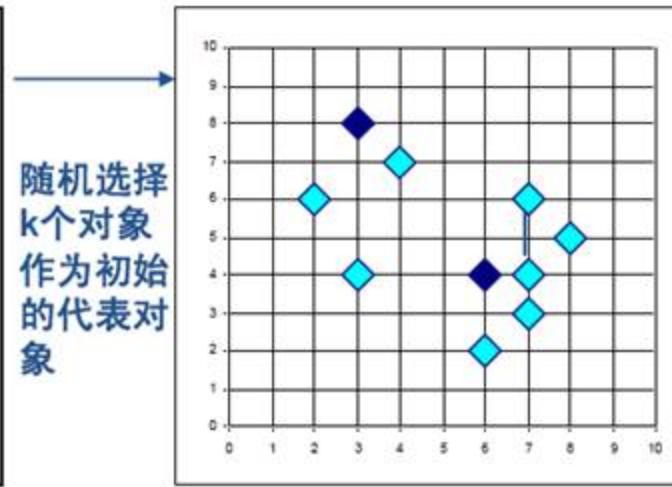
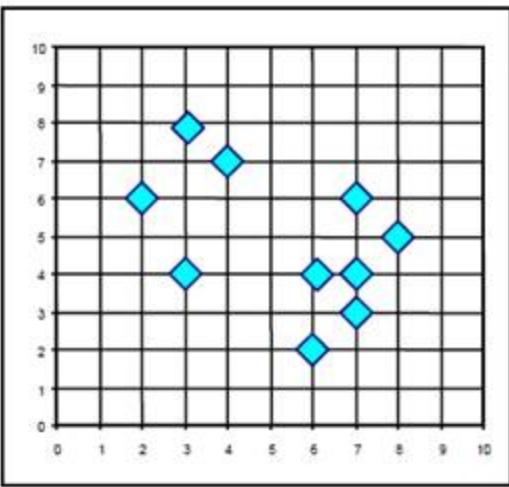


样本点	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0

2.1.4 PAM算法

- ◆ PAM算法(Partitioning Around Medoids)
- ◆ 最早提出的k-中心点算法之一
 - (1) 随机选择 k 个对象作为初始的代表对象
repeat
 - (2) 指派每个剩余的对象给离它最近的代表对象所代表的簇
 - (3) 随机地选择一个非代表对象 O_{random}
 - (4) 计算用 O_{random} 代替 O_j 的总代价 S
 - (5) 如果 $S < 0$, 则用 O_{random} 替换 O_j 形成新的 k 个代表对象的集合
until 1不发生变化

2.1.4 PAM算法示意

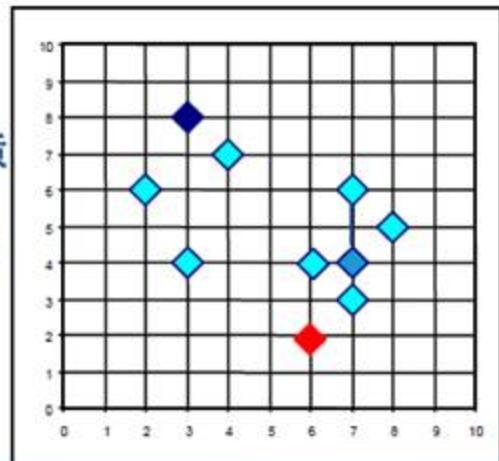


K=2

循环

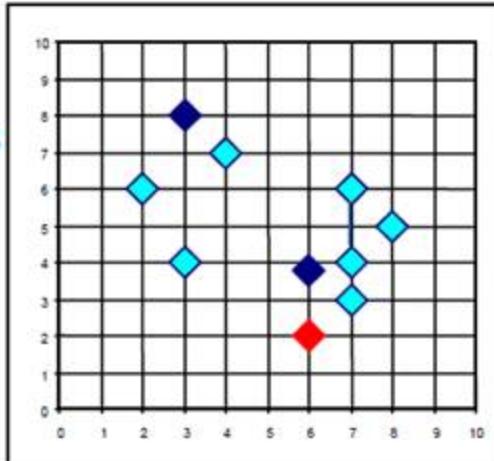
直到不发生变化

如果能提高质量，则交换
 O_{random} 和 O_j



随机地选择一个非代表对象 O_{random}

计算交换的总代价 S



2.1.4 PAM算法

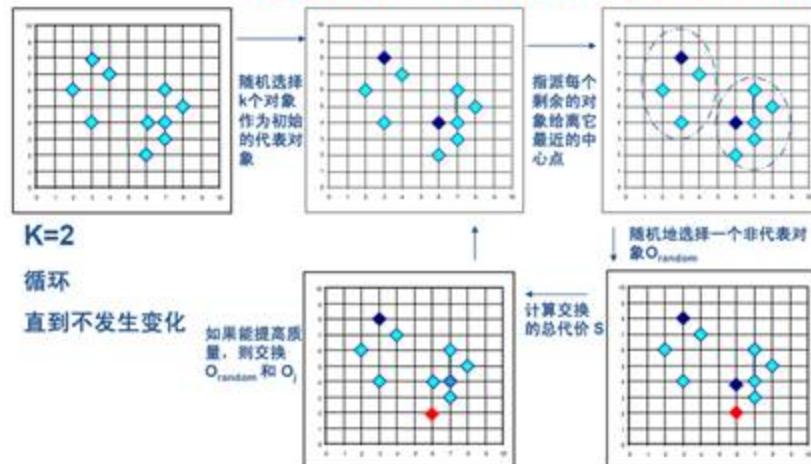
- ◆ 优点

- 当存在噪音和孤立点时，PAM 比 k-均值方法更健壮

- ◆ 缺点

- PAM 对于较小的数据集非常有效，但不能很好地扩展到大型数据集，原因是计算复杂度较高

- 每次中心点交换的时候，就要计算数据中每个点的代价



2.1.4 PAM算法

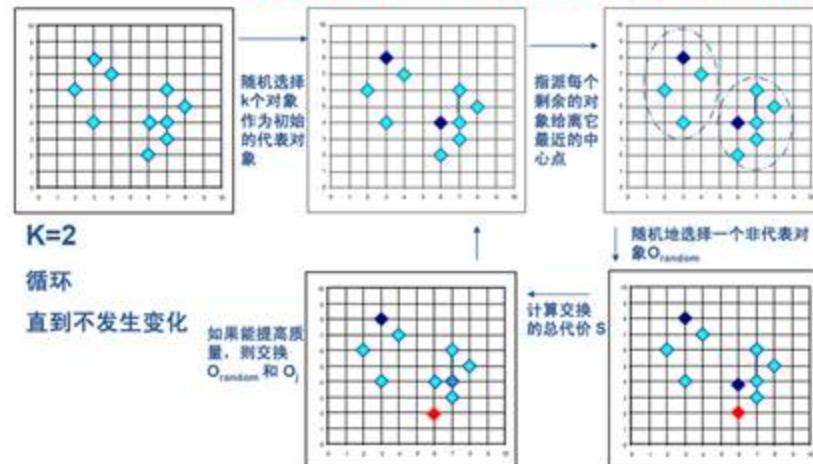
◆ 优点

- 当存在噪音和孤立点时，PAM 比 k-均值方法更健壮

◆ 缺点

- PAM 对于较小的数据集非常有效，但不能很好地扩展到大型数据集，原因是计算复杂度较高

- 每次中心点交换的时候，就要计算数据中每个点的代价



- PAM 算法中，每一次迭代都需要确定出 $k(n-k)$ 个交换对
- 对于交换对 i, h 需要计算 Tc^{ih}
- 在计算 TC^{ih} 过程中，需要计算 $n-k$ 个非中心点的代价
- 所以每次迭代的复杂度是 $k(n-k)^2$

下列哪些属于基于划分的聚类算法

- A K-means
- B K-modes
- C K-means++
- D K中心点

提交

下列说法正确的是

- A K-means算法能够解决有离群点的聚类问题
- B K-modes能够解决离散数据的聚类问题
- C K-means++能够解决初始点影响聚类效果的问题
- D K中心点能够解决有离群点的聚类问题

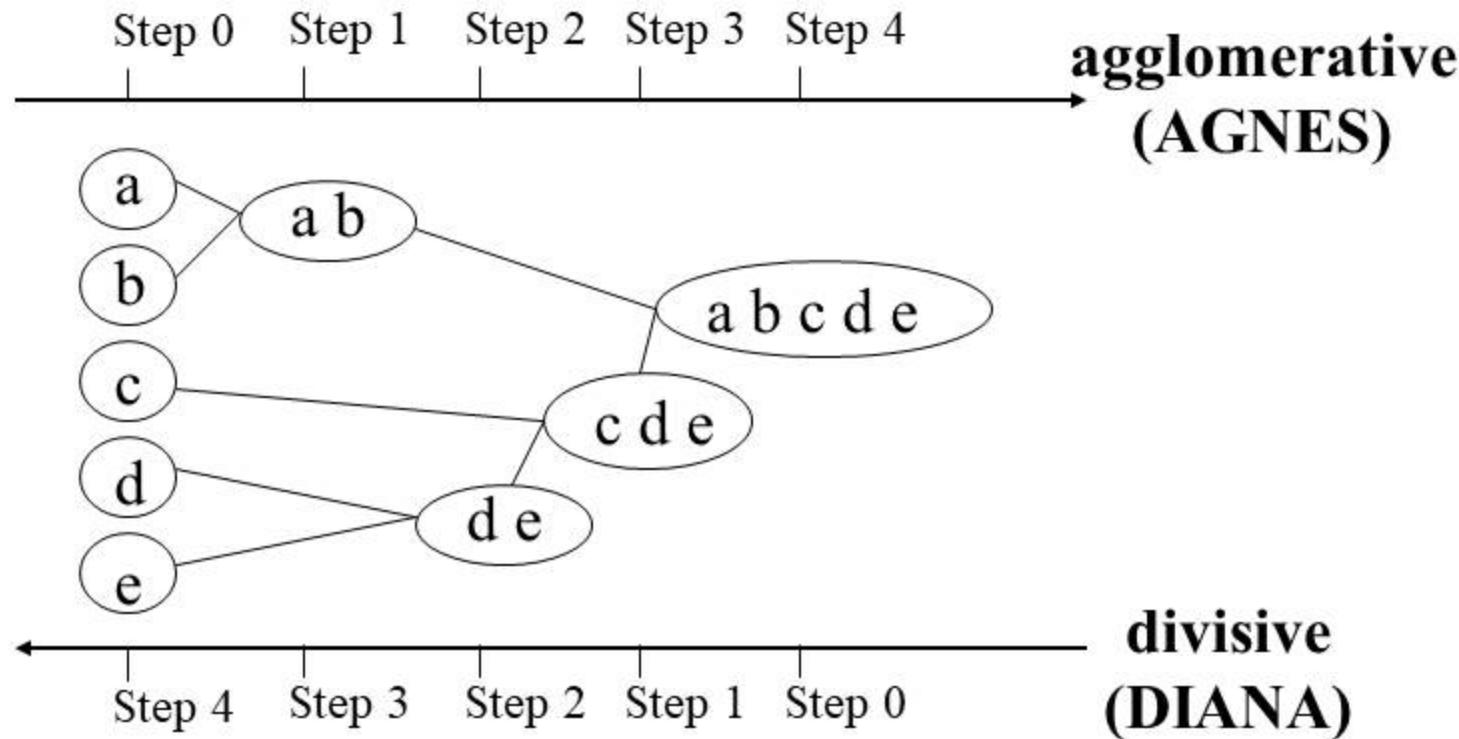
提交

2.2 层次方法

- ◆ 层次方法
 - 对给定数据对象集进行层次的分解
 - 使用距离矩阵作为聚类标准
 - 不需要输入聚类数目k，但需要终止条件
- ◆ 两种层次方法
 - 自底向上方法(凝聚)
 - 初始将每个对象作为单独的一个簇，然后相继的合并相近的对象或簇，直到所有的簇合并为一个，或者达到一个终止条件
 - 代表算法：AGNES算法
 - 自顶向下方法(分裂)
 - 初始将所有的对象置于一个簇中，在迭代的每一步，一个簇被分裂为多个更小的簇，直到最终每个对象在一个单独的簇中，或达到一个终止条件
 - 代表算法：DIANA算法

2.2 层次聚类过程图示

- ◆ 凝聚的(agglomerative)和分裂的(divisive)层次聚类图示



2.2.1AGNES算法

- ◆ AGNES (Agglomerative Nesting)算法
 - 首先，将数据集中的每个样本作为一个簇；
 - 然后，根据某些准则将这些簇逐步合并；
 - 合并的过程反复进行，直至不能再合并或者达到结束条件为止。
- ◆ 合并准则
 - 每次找到距离最近的两个簇进行合并。
 - 两个簇之间的距离由这两个簇中距离最近的样本点之间的距离来表示。

2.2.1AGNES算法

AGNES算法（自底向上合并算法）

输入：包含n个样本的数据集，终止条件簇的数目k。

输出：k个簇，达到终止条件规定的簇的数目。

(1) 初始时，将每个样本当成一个簇；

(2) REPEAT

 根据不同簇中最近样本间的距离找到最近的两个簇；

 合并这两个簇，生成新的簇的集合；

(3) UNTIL 达到定义的簇的数目。

算法终止条件是什么？

2.2.1AGNES算法

AGNES算法（自底向上合并算法）

输入：包含n个样本的数据集，终止条件簇的数目k。

输出：k个簇，达到终止条件规定的簇的数目。

(1) 初始时，将每个样本当成一个簇；

(2) REPEAT

根据不同簇中最近样本间的距离找到最近的两个簇；

合并这两个簇，生成新的簇的集合；

(3) UNTIL 达到定义的簇的数目。

算法终止条件是什么？

(1) 指定簇的数目k

2.2.1AGNES算法

AGNES算法（自底向上合并算法）

输入：包含n个样本的数据集，终止条件簇的数目k。

输出：k个簇，达到终止条件规定的簇的数目。

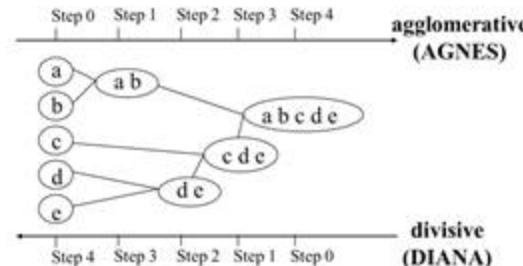
(1) 初始时，将每个样本当成一个簇；

(2) REPEAT

根据不同簇中最近样本间的距离找到最近的两个簇；

合并这两个簇，生成新的簇的集合；

(3) UNTIL 达到定义的簇的数目。



算法终止条件是什么？

(1) 指定簇的数目k

(2) 簇之间的距离超过一定阈值

2.2.1AGNES算法

AGNES算法（自底向上合并算法）

输入：包含n个样本的数据集，终止条件簇的数目k。

输出：k个簇，达到终止条件规定的簇的数目。

(1) 初始时，将每个样本当成一个簇；

(2) REPEAT

根据不同簇中最近样本间的距离找到最近的两个簇；

合并这两个簇，生成新的簇的集合；

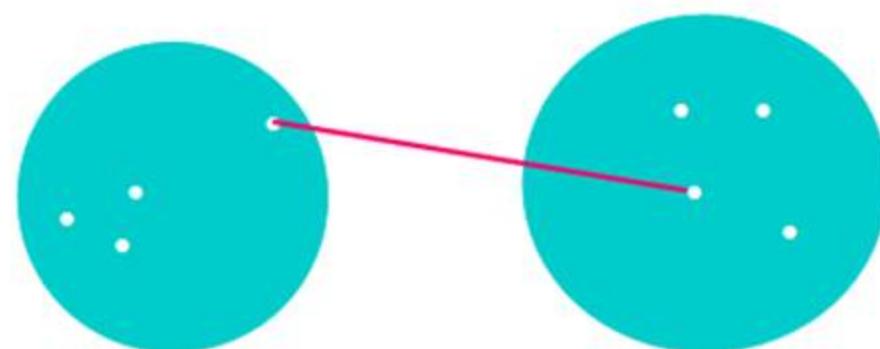
(3) UNTIL 达到定义的簇的数目。 簇之间的距离如何计算？

2.2.1AGNES算法-最小距离

- ◆ **单链接 (single-link) 方法**, 其每个簇可以用簇中所有对象代表, 簇间的相似度用属于不同簇中**最近的**数据点对之间的相似度来度量
- ◆ 也称为**最短距离法**, 定义簇的邻近度为取自不同簇的所有点对的俩个最近的点之间的邻近度

设 d_{ij} 表示样本 $X_{(i)}$ 和 $X_{(j)}$ 之间的距离, D_{ij} 表示类 G_i 和 G_j 之间距离

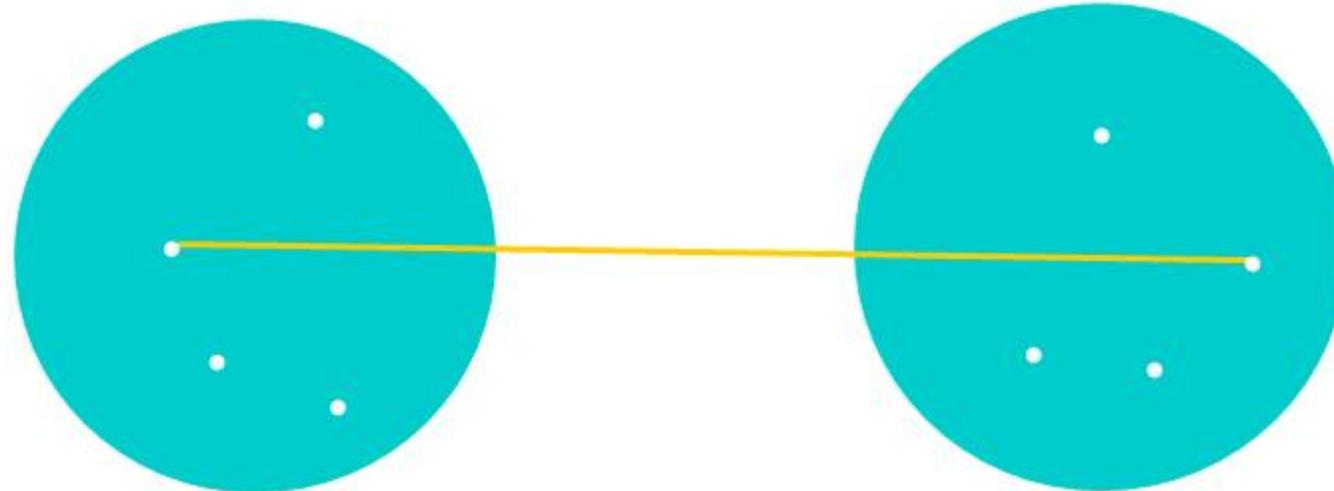
$$D_{ij} = \min_{X_{(i)} \in D_i, X_{(j)} \in D_j} d_{ij}$$



2.2.1AGNES算法-最大距离

- ◆ 全链

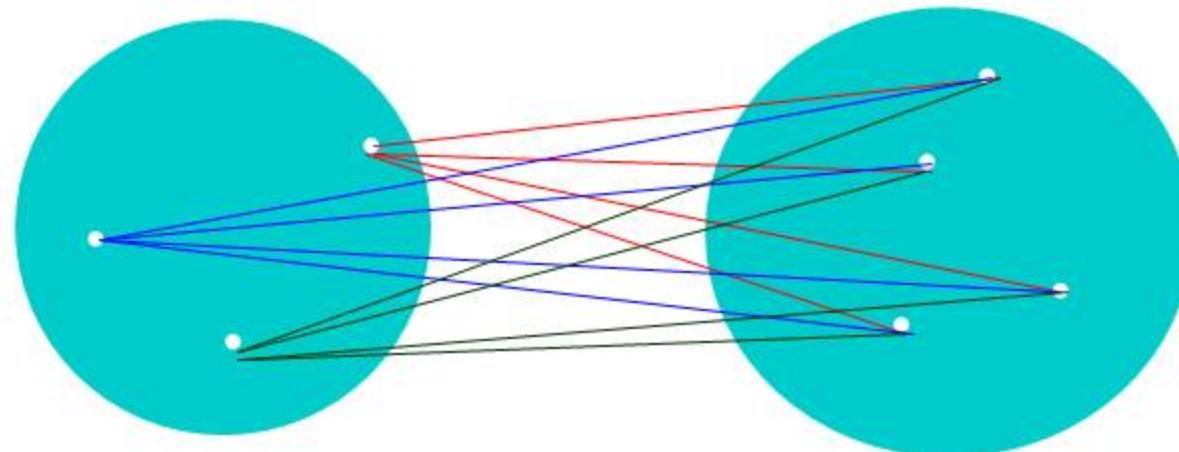
- 取自不同簇中的两个**最远的**点之间邻近度作为簇的邻近度，或者使用图的术语，不同的结点子集中两个结点之间的最长边 $D_{ij} = \max_{X_{(i)} \in G_i, X_{(j)} \in G_j} d_{ij}$



2.2.1AGNES算法-平均距离

- ◆ 组平均 (average linkage method) 类间所有样本点的平均距离
 - 该法利用了所有样本的信息，被认为是较好的系统聚类法

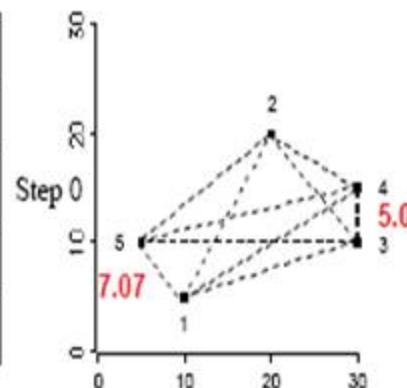
$$d_{avg}(C_i, C_j) = \frac{1}{n_i n_j} \sum_{p \in C_i} \sum_{p' \in C_j} \|p - p'\|$$



2.2.1单连接算法

	x1	x2
1	10	5
2	20	20
3	30	10
4	30	15
5	5	10

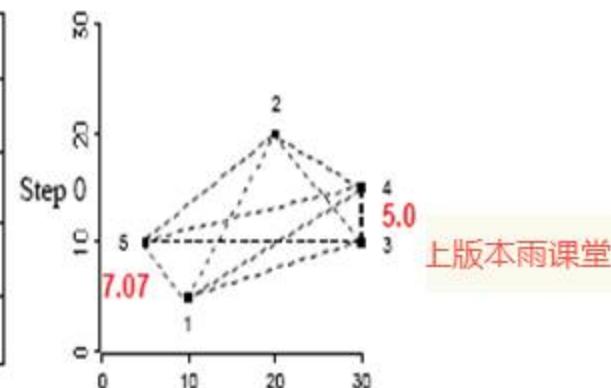
	1	2	(3)	4	5
1	0.00				
2	18.0	0.00			
3	20.6	14.1	0.00		
4	22.4	11.2	5.00	0.00	
5	7.07	18.0	25.0	25.5	0.00



- 先将五个样本都分别看成是一个簇，最靠近的两个簇是3和4，因为他们具有最小的簇间距离 $D(3, 4) = [填空1]$

	x1	x2
1	10	5
2	20	20
3	30	10
4	30	15
5	5	10

	1	2	(3)	4	5
1	0.00				
2	18.0	0.00			
3	20.6	14.1	0.00		
4	22.4	11.2	5.00	0.00	
5	7.07	18.0	25.0	25.5	0.00



作答

2.2.1 单连接算法

- 先将五个样本都分别看成是一个簇，最靠近的两个簇是3和4，因为他们具有最小的簇间距离 $D(3, 4) = 5.0$

1

x1

10

5

2

x2

20

20

3

30

10

4

30

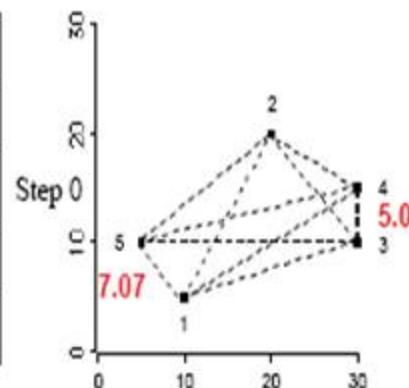
15

5

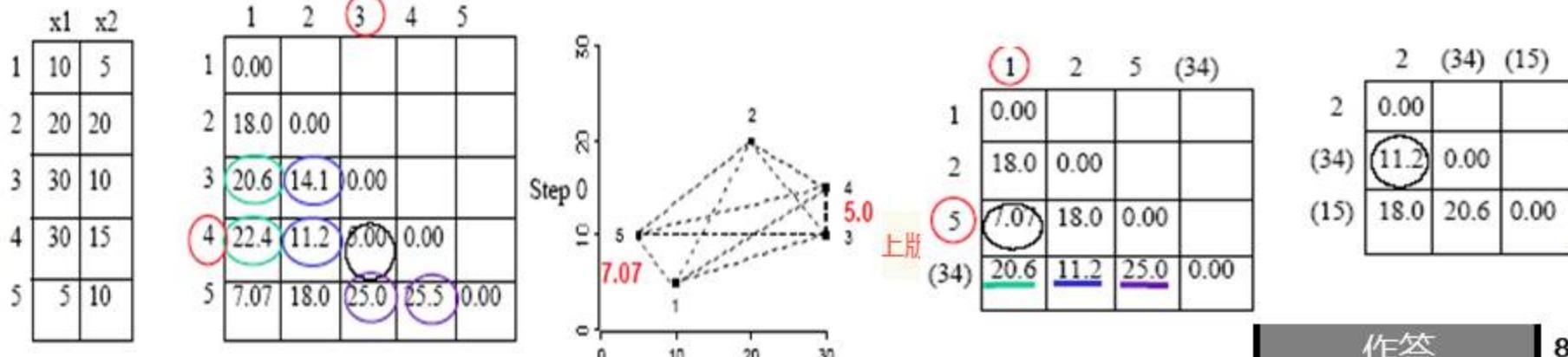
5

10

	1	2	(3)	4	5
1	0.00				
2	18.0	0.00			
3	20.6	14.1	0.00		
4	22.4	11.2	5.00	0.00	
5	7.07	18.0	25.0	25.5	0.00



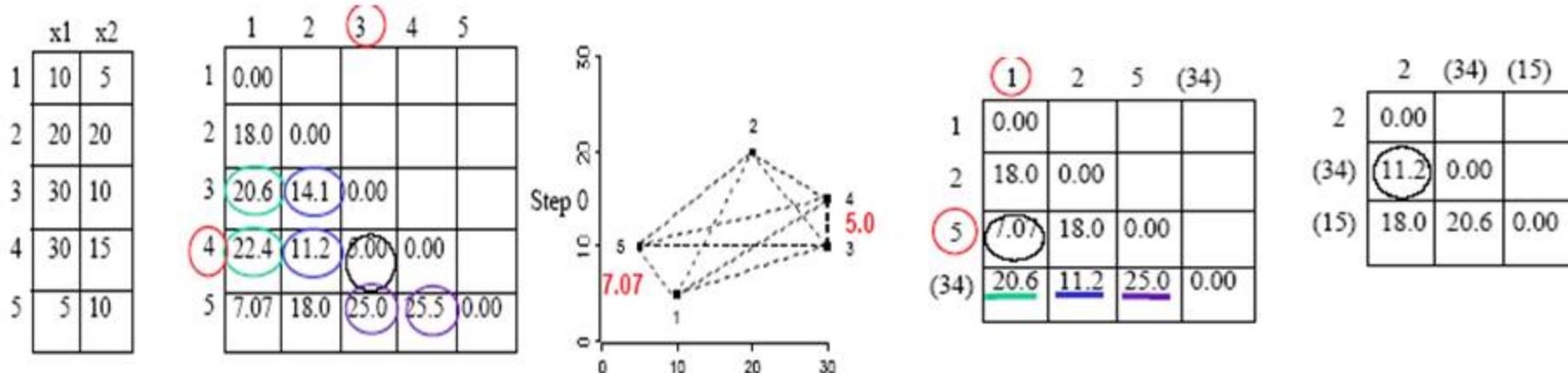
- 先将五个样本都分别看成是一个簇，最靠近的两个簇是3和4，因为他们具有最小的簇间距离 $D(3, 4) = 5.0$
- 合并簇3和4，得到新簇集合1, 2, (34), 5更新距离矩阵
 - $D(1, (34)) = \min(D(1, 3), D(1, 4)) = \min(20.6, 22.4) = [\text{填空1}]$
 - $D(2, (34)) = \min(D(2, 3), D(2, 4)) = \min(14.1, 11.2) = [\text{填空2}]$
 - $D(5, (34)) = \min(D(3, 5), D(4, 5)) = \min(25.0, 25.5) = [\text{填空3}]$



作答

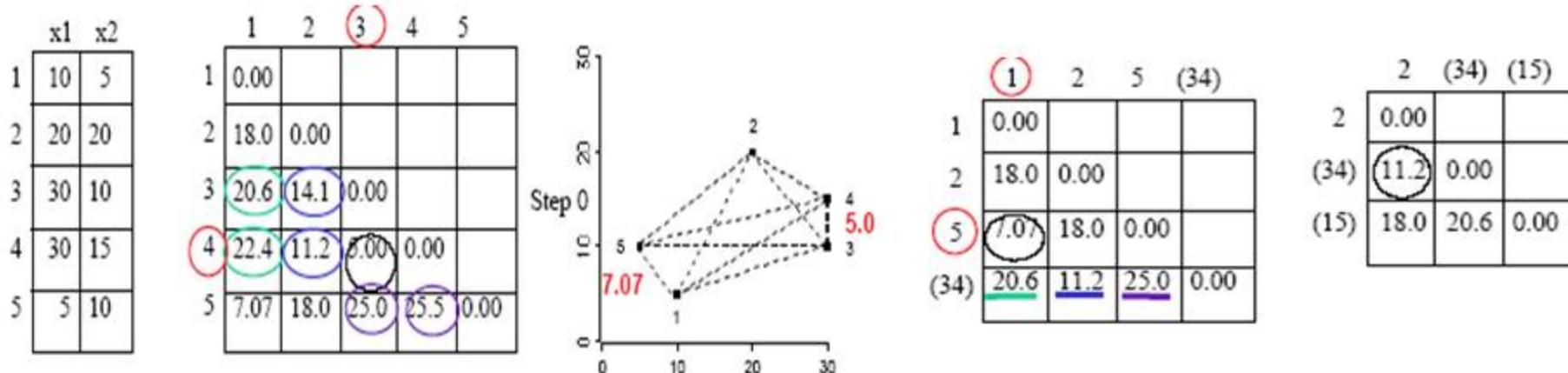
2.2.1 单连接算法

- 先将五个样本都分别看成是一个簇，最靠近的两个簇是3和4，因为他们具有最小的簇间距离 $D(3, 4) = 5.0$
- 合并簇3和4，得到新簇集合1, 2, (34), 5更新距离矩阵
 - $D(1, (34)) = \min(D(1, 3), D(1, 4)) = \min(20.6, 22.4) = 20.6$
 - $D(2, (34)) = \min(D(2, 3), D(2, 4)) = \min(14.1, 11.2) = 11.2$
 - $D(5, (34)) = \min(D(5, 3), D(5, 4)) = \min(25.0, 25.5) = 25.0$



2.2.1 单连接算法

- 先将五个样本都分别看成是一个簇，最靠近的两个簇是3和4，因为他们具有最小的簇间距离 $D(3, 4) = 5.0$
- 合并簇3和4，得到新簇集合1, 2, (34), 5更新距离矩阵
 - $D(1, (34)) = \min(D(1, 3), D(1, 4)) = \min(20.6, 22.4) = 20.6$
 - $D(2, (34)) = \min(D(2, 3), D(2, 4)) = \min(14.1, 11.2) = 11.2$
 - $D(5, (34)) = \min(D(3, 5), D(4, 5)) = \min(25.0, 25.5) = 25.0$
 - 原有簇1, 2, 5间的距离不变，故在四个簇1, 2, (34), 5中，最靠近的两个簇是1和5，它们具有最小簇间距离 $D(1, 5) = 7.07$

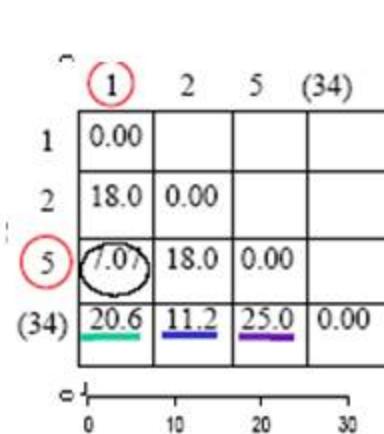


2.2.1单连接算法

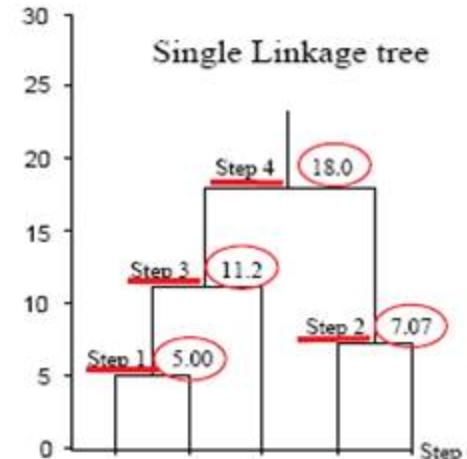
- 先将五个样本都分别看成是一个簇，最靠近的两个簇是3和4，因为他们具有最小的簇间距离 $D(3, 4) = 5.0$
- 合并簇3和4，得到新簇集合1, 2, (34), 5更新距离矩阵
 - $D(1, (34)) = \min(D(1, 3), D(1, 4)) = \min(20.6, 22.4) = 20.6$
 - $D(2, (34)) = \min(D(2, 3), D(2, 4)) = \min(14.1, 11.2) = 11.2$
 - $D(5, (34)) = \min(D(3, 5), D(4, 5)) = \min(25.0, 25.5) = 25.0$
 - 原有簇1, 2, 5间的距离不变，故在四个簇1, 2, (34), 5中，最靠近的两个簇是1和5，它们具有最小簇间距离 $D(1, 5) = 7.07$

	x1	x2
1	10	5
2	20	20
3	30	10
4	30	15
5	5	10

	1	2	(3)	4	5
1	0.00				
2	18.0	0.00			
3	20.6	14.1	0.00		
4	22.4	11.2	5.00	0.00	
5	7.07	18.0	25.0	25.5	0.00



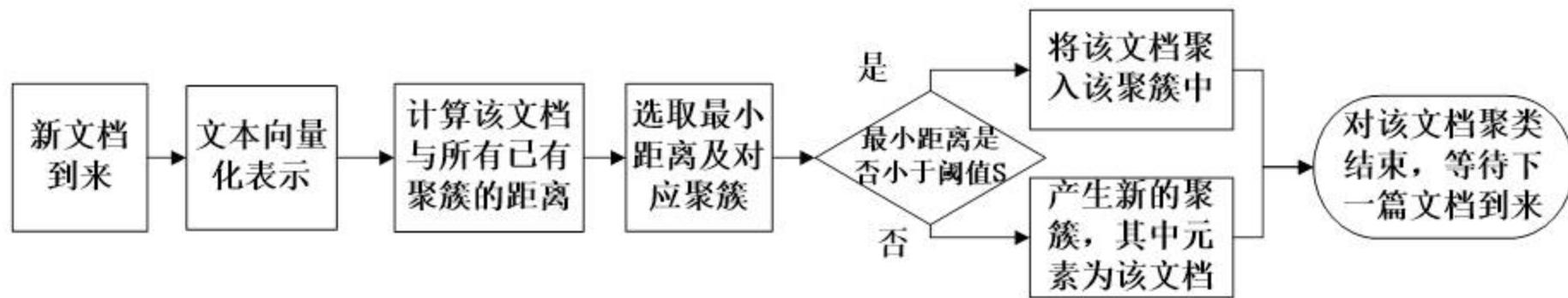
	2	(34)	(15)
2	0.00		
(34)	11.2	0.00	
(15)	18.0	20.6	0.00



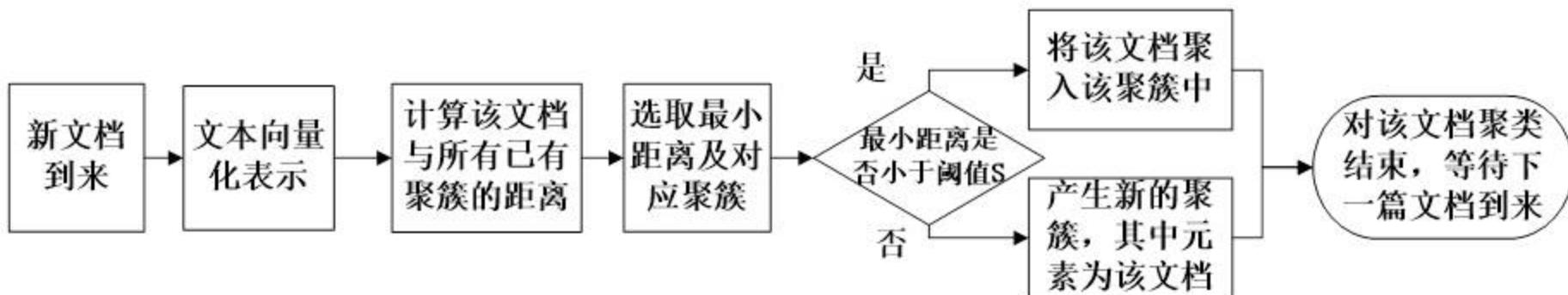
2.2.2 层次方法的问题及改进

- ◆ 层次聚类存在的主要问题
 - 合并或分裂的决定需要检查和估算大量的对象或簇
 - 一个步骤一旦完成便不能被撤消
 - 避免考虑选择不同的组合，减少计算代价
 - 不能更正错误的决定
 - 不具有很好的可伸缩性
- ◆ 改进方法：将层次聚类和其他的聚类技术进行集成，形成多阶段聚类
 - BIRCH：使用CF-tree对对象进行层次划分，然后采用其他的聚类算法对聚类结果进行求精
 - CURE：采用固定数目的代表对象来表示每个簇，然后依据一个指定的收缩因子向着聚类中心对它们进行收缩
 - CHAMELEON：使用动态模型进行层次聚类

2.2.2 层次方法变形 (single-pass单遍聚类)



2.2.2 层次方法变形 (single-pass单遍聚类)



讨论：

- (1) 如何计算文档距离
- (2) 如何计算文档与聚类的距离
- (3) 算法停止条件

层次聚类法计算2个簇之间的距离有哪些方法

- A 最小值
- B 最大值
- C 平均值

提交

2.2.3 层次聚类编程实践

- ◆ <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.AgglomerativeClustering.html#sklearn.cluster.AgglomerativeClustering>
- 同学们可以尝试利用python读入本地iris数据集，来完成层次聚类，分析其聚类效果

2.3 基于密度的方法

- ◆ 基于密度聚类

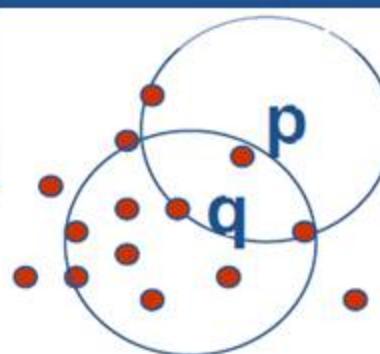
- 根据密度条件对邻近对象分组形成簇，簇的增长或者根据邻域密度，或者根据特定的密度函数(只要临近区域的密度超过某个阈值，就继续聚类)

- ◆ 主要特点

- 发现任意形状的聚类
- 处理噪音
- 一遍扫描
- 需要密度参数作为终止条件



2.3.1密度概念 1/2



MinPts = 5

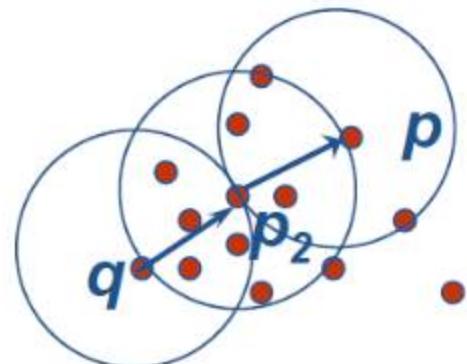
$\varepsilon = 1 \text{ cm}$

- ◆ ε -邻域
 - 给定对象半径 ε 内的邻域称为该对象的 ε -邻域
- ◆ 核心对象
 - 如果对象的 ε -邻域至少包含最小数目MinPts个对象，则称该对象为核心对象
- ◆ 直接密度可达
 - 给定对象集合D，如果p是在q 的 ε -邻域内，而q 是核心对象，则称对象p 是从对象q 关于 ε 和MinPts直接密度可达的

2.3.2密度概念 2/2

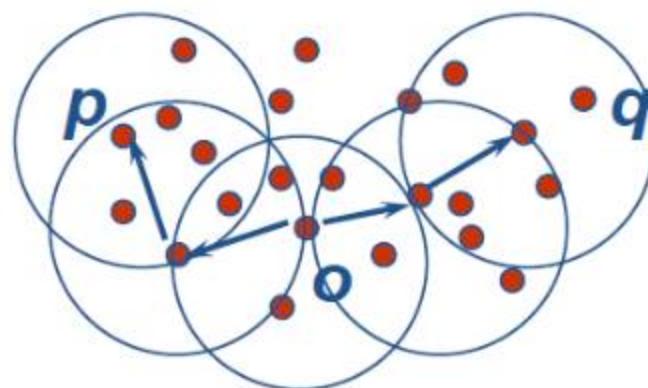
- ◆ 密度可达

- 如果存在一个对象链 p_1, \dots, p_n , $p_1 = q, p_n = p$, 使得 p_{i+1} 是从 p_i 直接密度可达的, 则称对象 p 是从对象 q 关于 ε 和 MinPts(间接)密度可达的



- ◆ 密度相连的

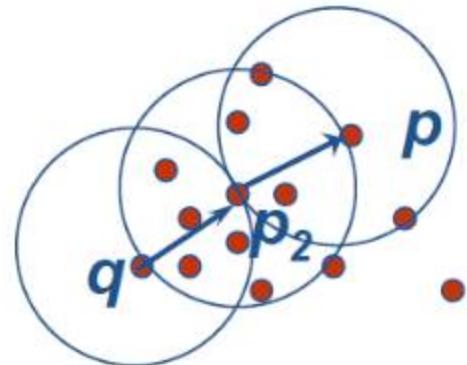
- 如果存在对象o 使得, p 和q 都是从o 关于 ε 和 MinPts 密度可达的, 则称对象p 与q 关于 ε 和 MinPts 是密度相连的



2.3.2 密度概念 2/2

◆ 密度可达

- 如果存在一个对象链 p_1, \dots, p_n , $p_1 = q$, $p_n = p$, 使得 p_{i+1} 是从 p_i 直接密度可达的, 则称对象 p 是从对象 q 关于 ε 和 MinPts(间接)密度可达的

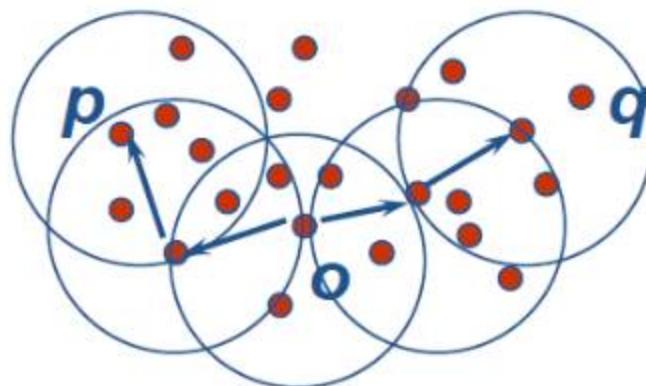


▪ 讨论：对称性？

◆ 密度相连的

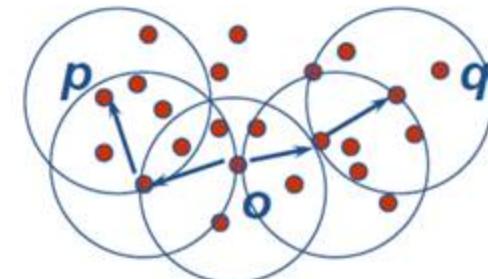
- 如果存在对象 o 使得, p 和 q 都是从 o 关于 ε 和 MinPts 密度可达的, 则称对象 p 与 q 关于 ε 和 MinPts 是密度相连的

▪ 讨论：对称性

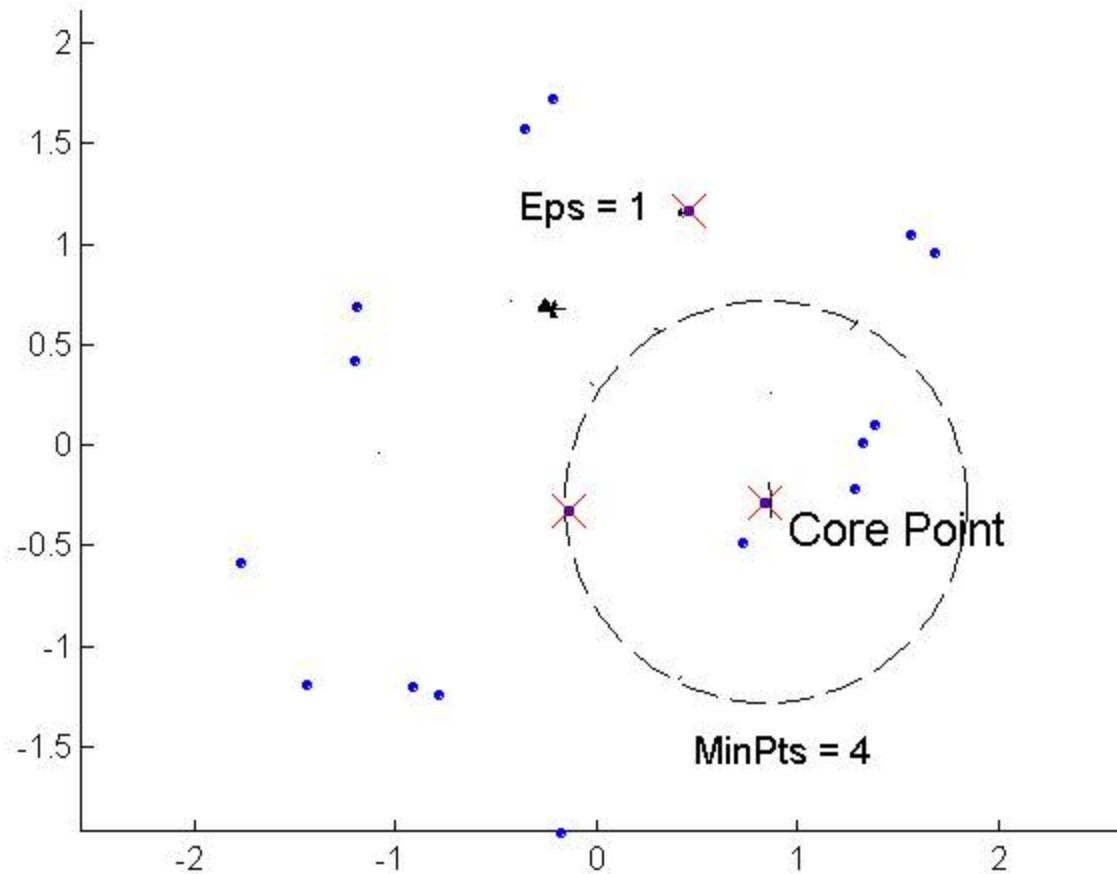


2.3.2 基于密度的簇

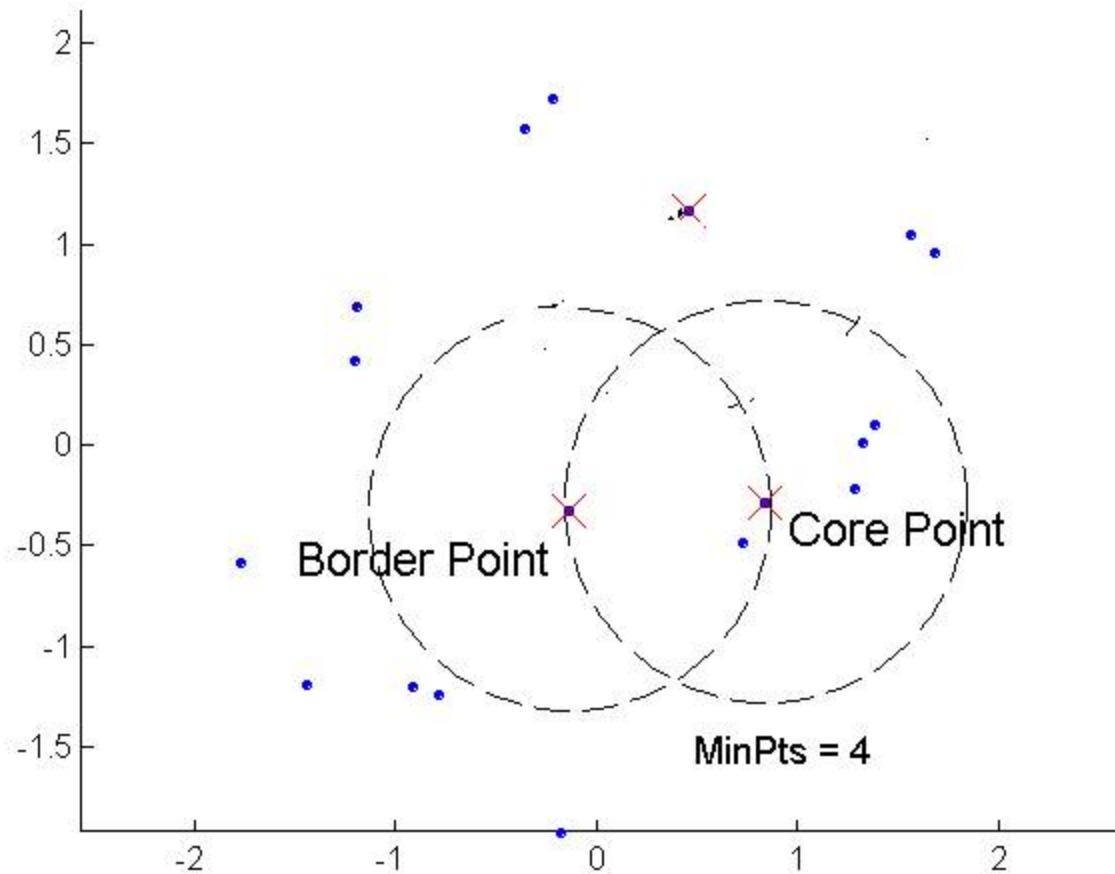
- ◆ 基于密度的簇是基于密度可达性的**最大的密度相连**对象的集合，即簇C 满足
 - 连通性：对于C 中任意的p 和q，p 与q 是关于 ε 和MinPts密度相连的
 - 极大性：对于任意的p 和q，如果p 属于C 簇，并且q 是从p 出发关于 ε 和MinPts密度可达的，则q 也属于C 簇
- ◆ 边界点（Border point）的Eps 邻域有少于MinPts 个对象，但它的邻域中有核心对象
- ◆ 噪声点（Noise point）是除核心对象和边界点之外的点
- ◆ 典型算法
 - DBSCAN、DENCLUE、OPTICS



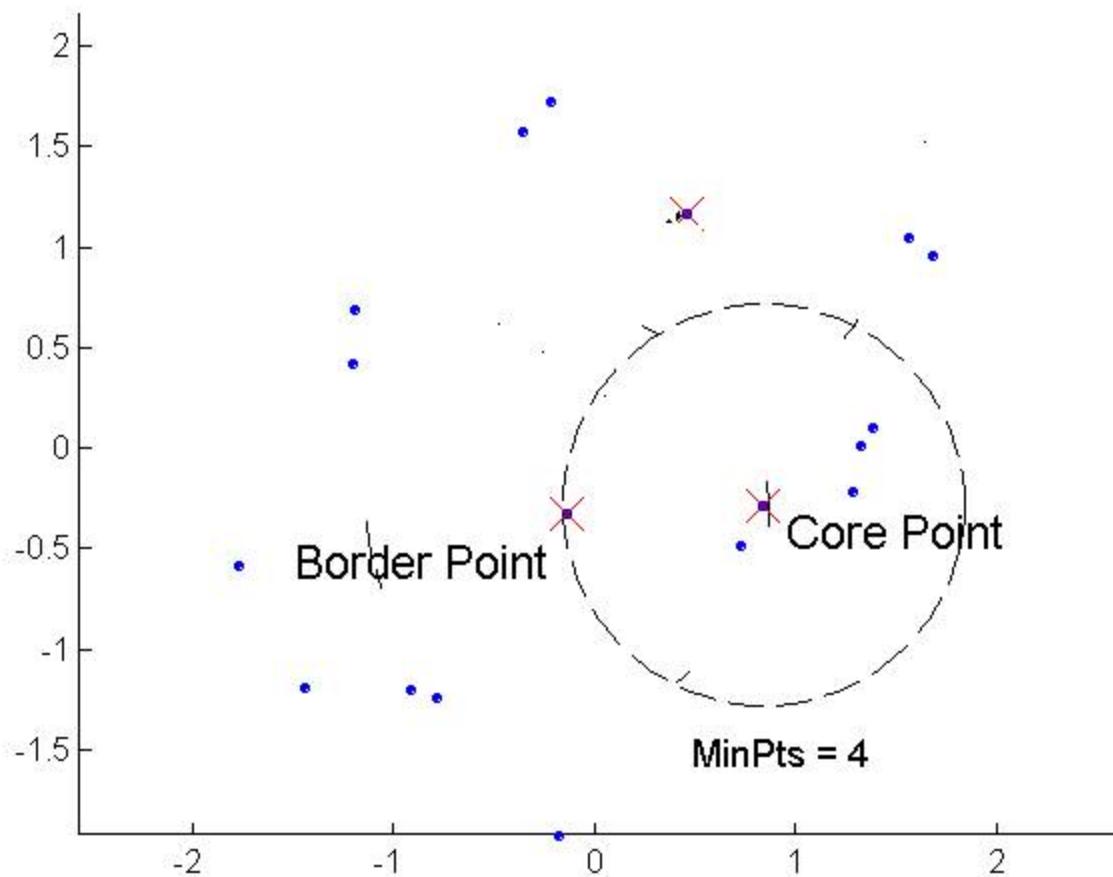
2.3.3 基本概念理解



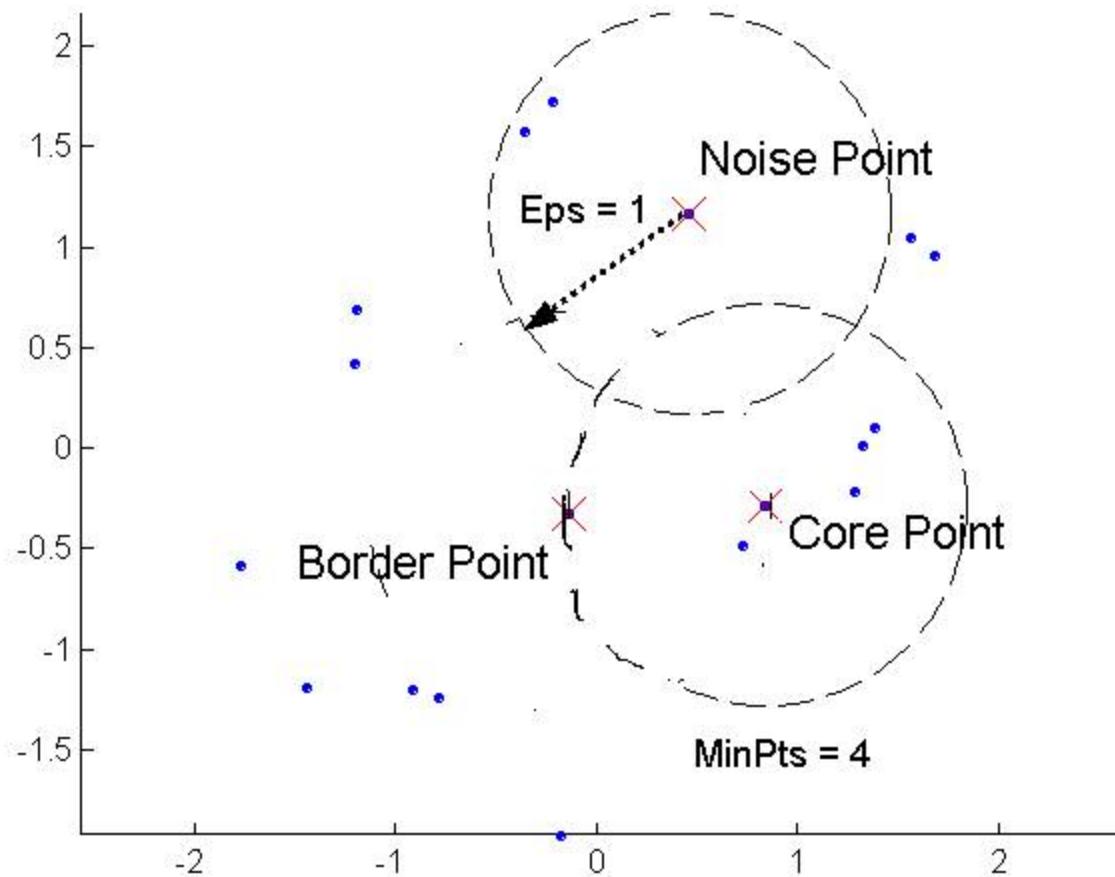
2.3.3 基本概念理解



2.3.3 基本概念理解



2.3.3 基本概念理解



2.3.4 DBSCAN

◆ 算法描述

输入：包含n个对象的数据库，半径 ϵ ，最少数目MinPts。

输出：所有生成的簇，达到密度要求。

REPEAT

 从数据库中抽取一个未处理过的点；

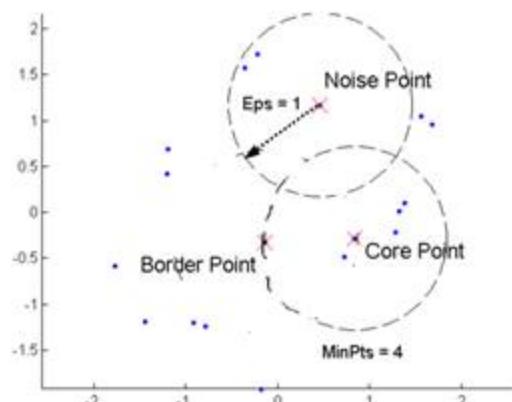
 IF 抽出的点是核心点 THEN

 找出所有从该点密度可达的对象，形成一个簇

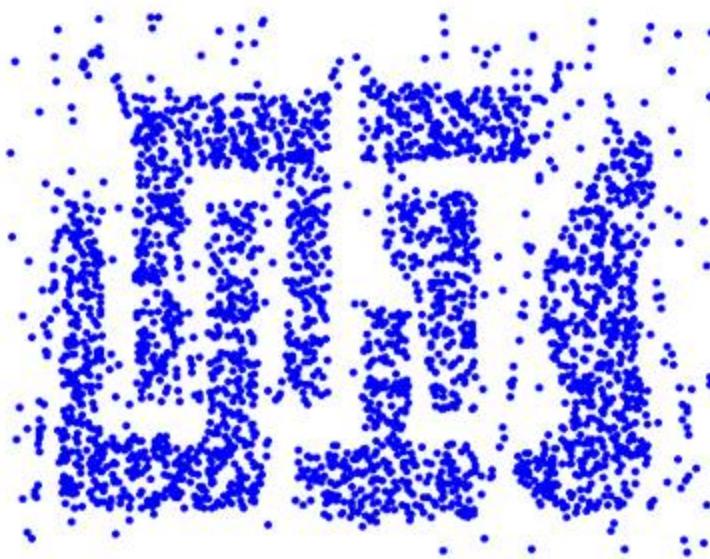
 ELSE

 抽出的点是边缘点(非核心对象)，跳出本次循环，寻找下一点；

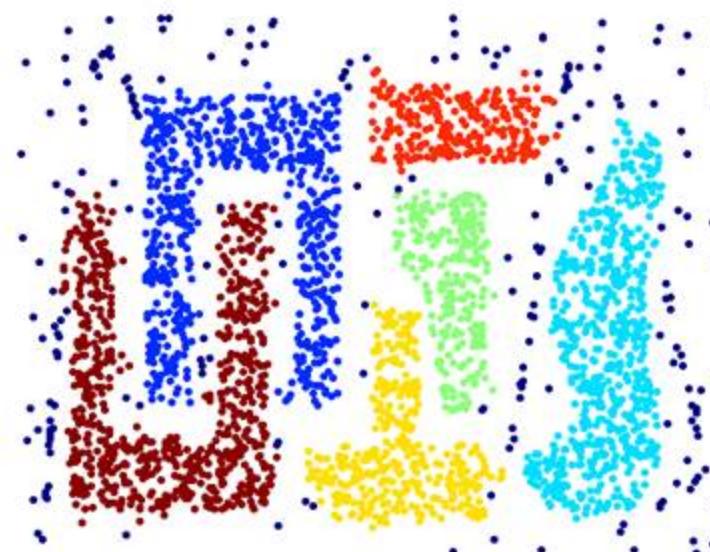
UNTIL 所有点都被处理；



2.3.5 DBSCAN优点



Original Points



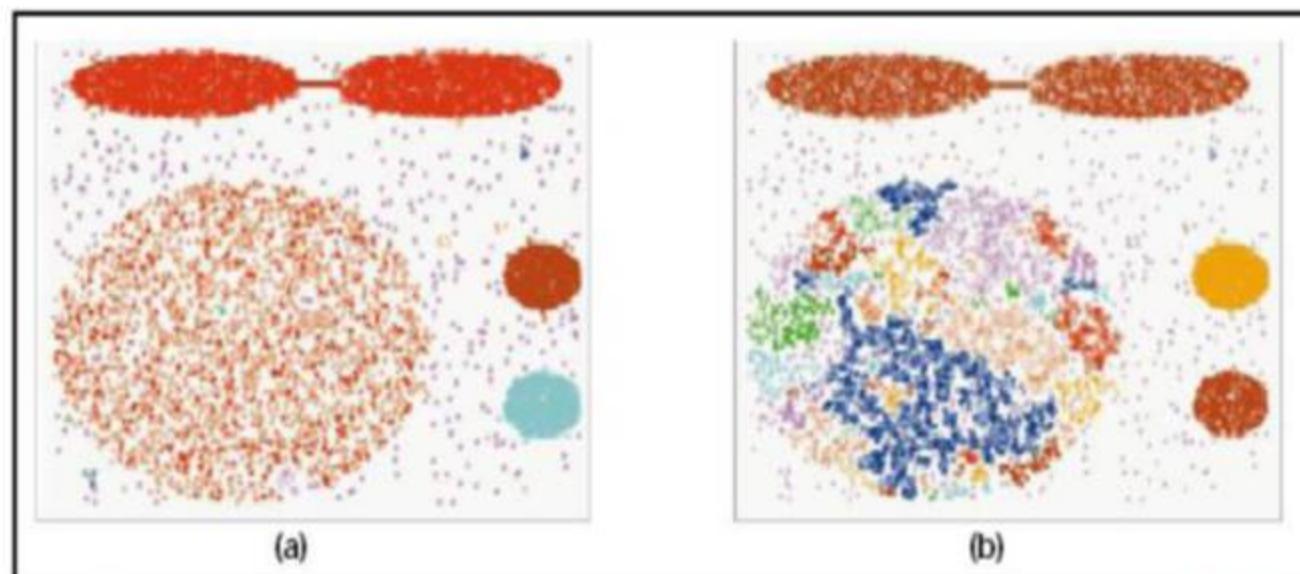
Clusters

- 抗噪声
- 能处理各种形状和大小集群

2.3.6 DBSCAN缺点

- ◆ 缺点

- 对用户定义的参数是敏感的，参数难以确定(特别是对于高维数据)，设置的细微不同可能导致差别很大的聚类。全局密度参数不能刻画内在的聚类结构



$\epsilon = 0.5$

MinPts = 4

$\epsilon = 0.4$

MinPts = 4

2.3.6 DBSCAN缺点

- ◆ 缺点

- 对用户定义的参数是敏感的，参数难以确定(特别是对于高维数据)，设置的细微不同可能导致差别很大的聚类。全局密度参数不能刻画内在的聚类结构

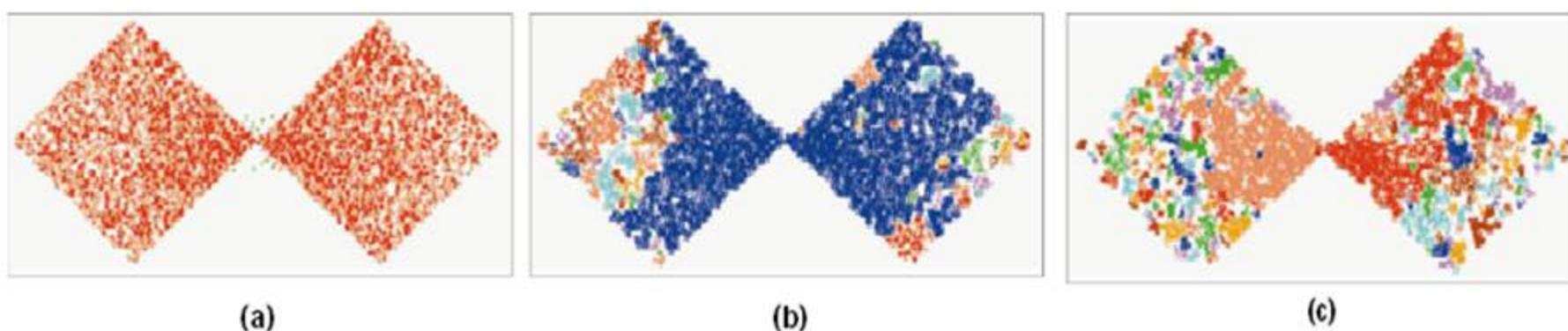


Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.

下列哪种聚类算法对聚簇的形状不敏感

- A K-means
- B K-中心点
- C AGNES
- D DBSCAN

提交

2.2.3密度聚类编程实践

- ◆ <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.DBSCAN.html#sklearn.cluster.DBSCAN>
- 同学们可以尝试利用python读入本地iris数据集，来完成密度聚类，分析其聚类效果

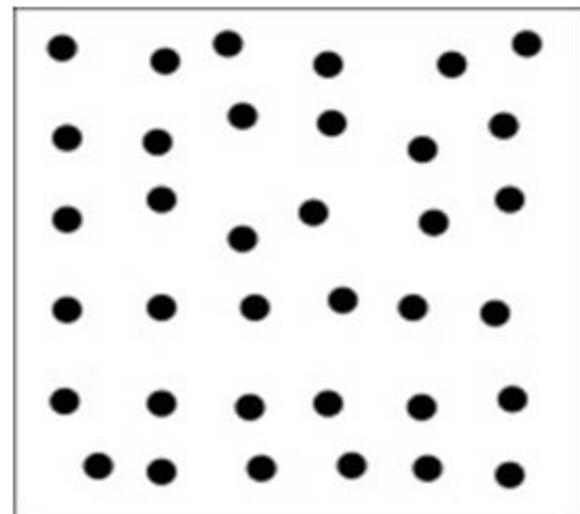
3. 聚类评估

- ◆ 估计在数据集上进行聚类的可行性和被聚类方法产生的结果的质量
- ◆ 聚类评估的任务
 - 3.1 估计聚类趋势：评估数据集是否存在非随机结构
 - 3.2 确定数据集中的簇数：在聚类之前，估计簇数
 - 3.3 测定聚类质量：聚类之后，评估结果簇的质量

3.1 估计聚类趋势

- ◆ 从统计学的角度，就是检测数据是否是随机或线性分布的

如果数据服从均匀分布，
显然对其进行的聚类操作都是没有意义的



3.2 确定数据集中的簇数

- ◆ 实验方法

- 对于n个点的数据集，簇数 $\sqrt{n/2}$ ，每个簇约有 $\sqrt{2n}$ 个点

- ◆ 肘方法 (Elbow method)

- 增加簇数可以降低簇内方差之和，但是如果形成太多的簇，降低簇内方差之和的边缘效应可能下降。

$$\alpha(\mathbf{o}) = \frac{\sum_{\mathbf{o}' \in C_i, \mathbf{o} \neq \mathbf{o}'} dist(\mathbf{o}, \mathbf{o}')}{|C_i| - 1}$$

$$b(\mathbf{o}) = \min_{C_j: 1 \leq j \leq k, j \neq i} \left\{ \frac{\sum_{\mathbf{o}' \in C_j} dist(\mathbf{o}, \mathbf{o}')}{|C_j|} \right\}$$

3.3 测定聚类质量

- ◆ 外在方法：有监督的
 - 用某种聚类质量度量对聚类结果和基准进行比较
 - 基准：一种理想的聚类，由专家构建
- ◆ 内在方法：无监督的
 - 通过考察簇的分离情况和簇的紧凑情况来评估聚类
 - 例：轮廓系数

$$\alpha(o) = \frac{\sum_{o' \in C_i, o \neq o'} dist(o, o')}{|C_i| - 1}$$

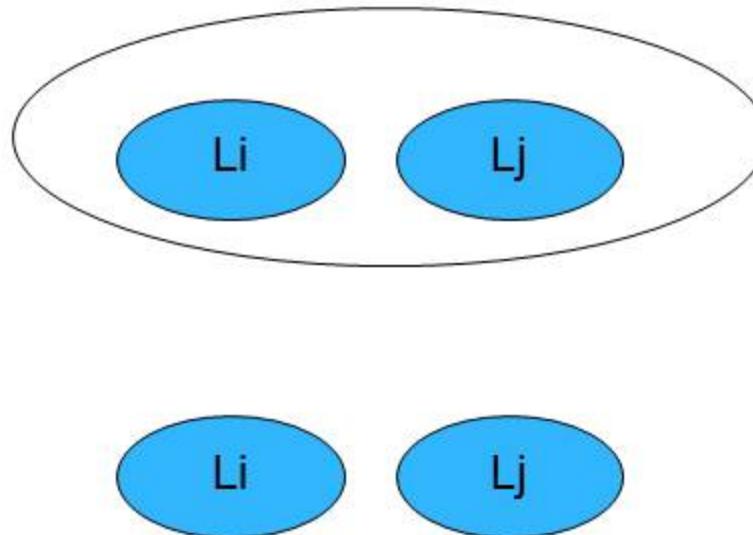
$$b(o) = \min_{C_j: 1 \leq j \leq k, j \neq i} \left\{ \frac{\sum_{o' \in C_j} dist(o, o')}{|C_j|} \right\}$$

3.4 测定聚类质量

- ◆ 簇的同质性
- ◆ 簇的完全性
- ◆ 碎布袋性
- ◆ 小簇保持性

3.4 测定聚类质量

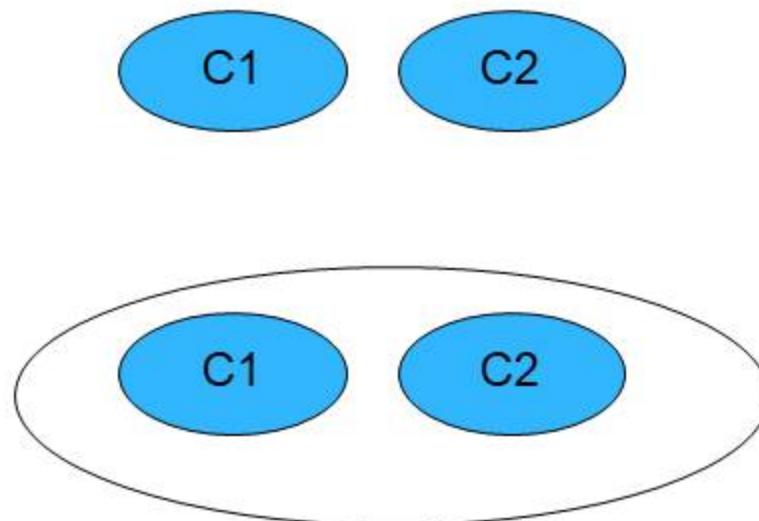
- ◆ 簇的同质性
- ◆ 簇的完全性
- ◆ 碎布袋性
- ◆ 小簇保持性



- 簇的同质性 (cluster homogeneity)。这要求，聚类中的簇越纯，聚类越好。假设基准是说数据集 D 中的对象可能属于类别 L_1, \dots, L_n 。考虑一个聚类 C_1 ，其中簇 $C \in C_1$ 包含来自两个类 L_i 和 L_j ($1 \leq i < j \leq n$) 的对象。再考虑一个聚类 C_2 ，除了把 C 划分成分别包含 L_i 和 L_j 中对象的两个簇之外，它等价于 C_1 。关于簇的同质性，聚类质量度量 Q 应该赋予 C_2 比 C_1 更高的得分，即 $Q(C_2, C_g) > Q(C_1, C_g)$ 。

3.4 测定聚类质量

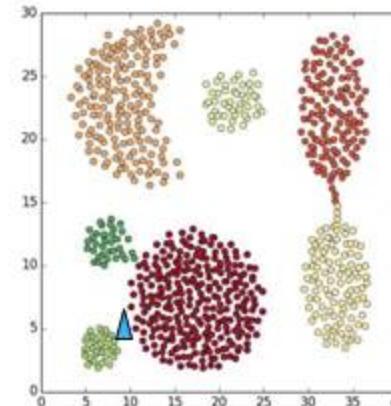
- ◆ 簇的同质性
- ◆ 簇的完全性
- ◆ 碎布袋性
- ◆ 小簇保持性



- 簇的完全性 (cluster completeness)。这与簇的同质性相辅相成。簇的完全性要求对于聚类来说，根据基准，如果两个对象属于相同的类别，则它们应该被分配到相同的簇。簇的完全性要求聚类把（根据基准）属于相同类别的对象分配到相同的簇。考虑聚类 C_1 ，它包含簇 C_1 和 C_2 ，根据基准，它们的成员属于相同的类别。假设 C_2 除 C_1 和 C_2 在 C_2 中合并到一个簇之外，它等价于聚类 C_1 。关于簇的完全性，聚类质量度量 Q 应该赋予 C_2 更高的得分，即 $Q(C_2, C_g) > Q(C_1, C_g)$ 。

3.4 测定聚类质量

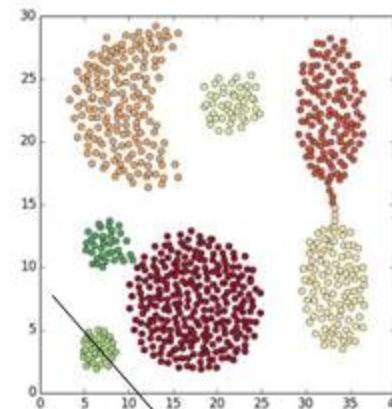
- ◆ 簇的同质性
- ◆ 簇的完全性
- ◆ 碎布袋性



- 把一个异构对象放到一个纯聚簇中要比把它放到一个碎布袋里（例如，“杂项”或“其他”类别）更糟糕
- 也就是说更希望噪音不被聚到已存在的聚簇中
- putting a heterogeneous object into a pure cluster should be penalized more than putting it into a rag bag (i.e., “miscellaneous” or “other” category)
- ◆ 小簇保持性

3.4 测定聚类质量

- ◆ 簇的同质性
- ◆ 簇的完全性
- ◆ 碎布袋性
- ◆ 小簇保持性



- 把一个小类别聚簇分成更小的聚簇比把一个大类别分成小聚簇更有害
- 也就是我们希望小类别聚簇不再被划分
- splitting a small category into pieces is more harmful than splitting a large category into pieces

2.2.3聚类评估编程实践

- ◆ <https://scikit-learn.org/stable/modules/clustering.html#clustering-performance-evaluation>
- 同学们可以尝试利用python读入本地iris数据集，根据前面实践的聚类算法评估其效果

主要内容

- ◆ 1. 聚类分析概述
- ◆ 2. 基本聚类方法
- ◆ 3. 聚类评估

第7次课后作业

- ◆ 第七次课后作业-在educoder平台上完成作业
- ◆ <https://www.educoder.net/shixuns/jfzo8xcu/challenges>
- ◆ <https://www.educoder.net/shixuns/qy9gozt8/challenges>
- ◆ <https://www.educoder.net/shixuns/qy9gozt8/challenges>
- ◆ <https://www.educoder.net/shixuns/ik6c3hpa/challenges>
- ◆ <https://www.educoder.net/shixuns/xqu69twm/challenges>

提交作业截至时间：2020年3月12日

Q&A