
聚类任务

陈飞宇

fchen@cqu.edu.cn

办公室：软件学院529

7.1 聚类任务简介

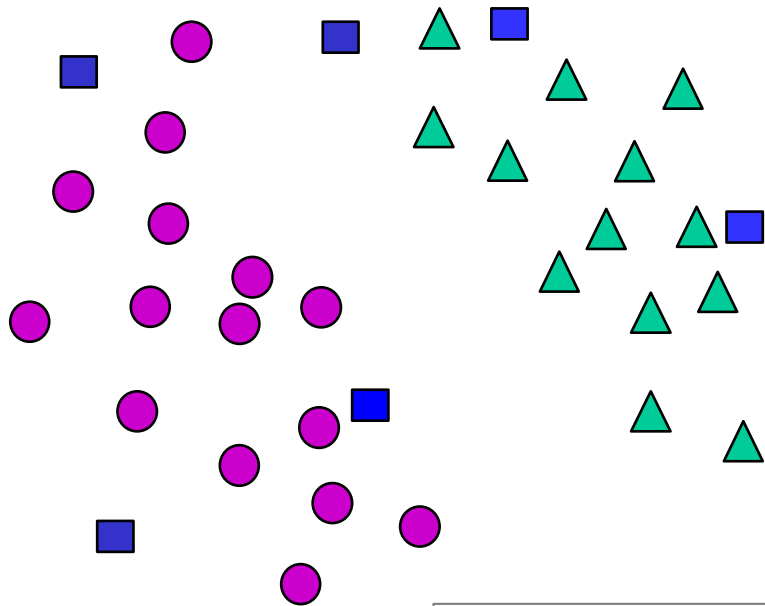
陈飞宇

fchen@cqu.edu.cn

办公室：软件学院529

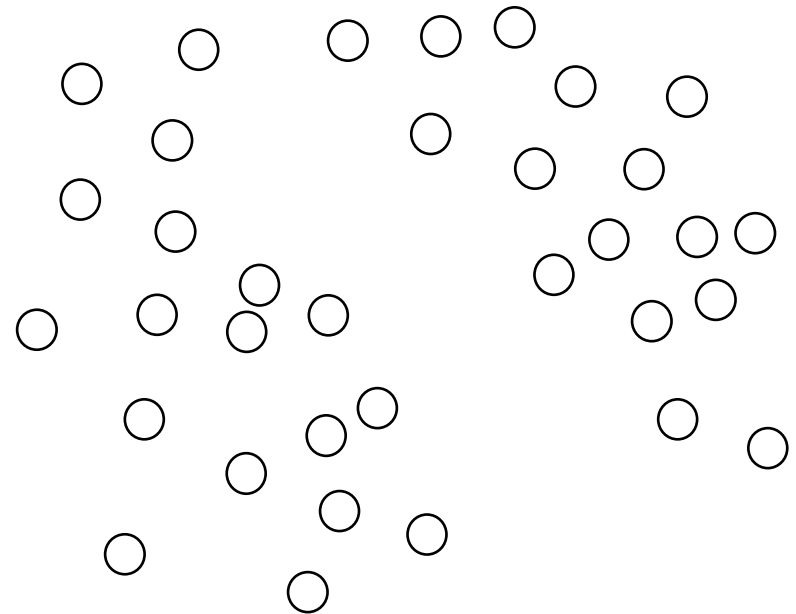
分类与聚类

分类



正例	●
反例	▲
验证集	■

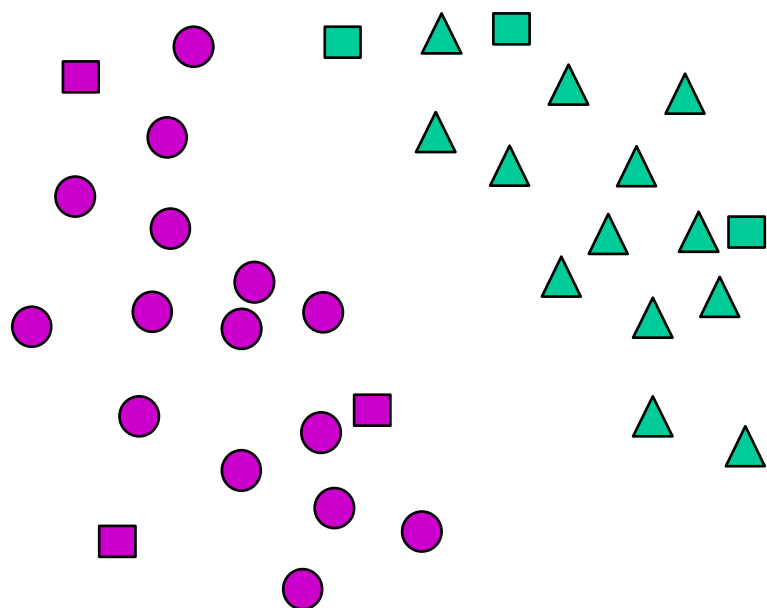
聚类



样本点	○
-----	---

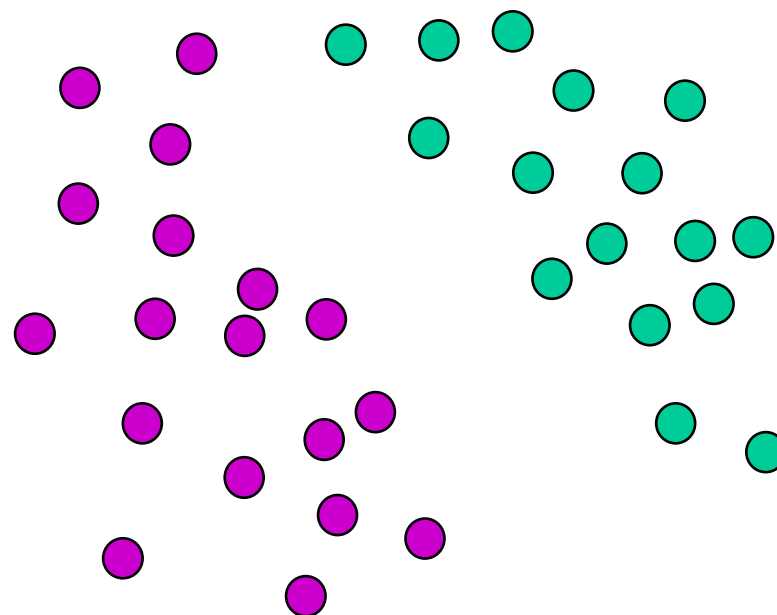
分类与聚类

分类



正例	●
反例	▲
验证集	■

聚类



簇一	●
簇二	●

聚类问题

聚类：根据某种相似性，把一组数据划分成若干个簇的过程。

难点一：相似性很难精准定义！

难点二：可能存在的划分太多！

难点三：若干个簇 = ？

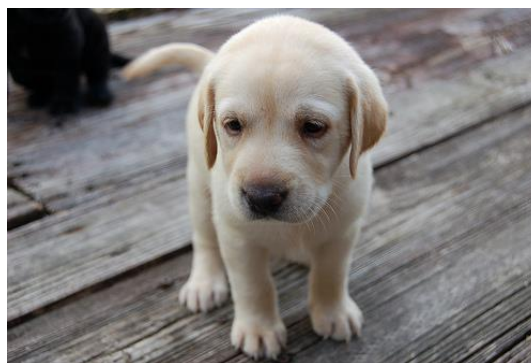
1. 相似性定义？



相似？



相似？



相似？

2. 可能的划分？

假设我们要把 $n = 5$ 个不同的数据点放入不相同的 $k = 2$ 个无差别的盒子中，有几种方案？

$$C_1 = \{1\}, C_2 = \{2, 3, 4, 5\}$$

$$C_1 = \{2\}, C_2 = \{1, 3, 4, 5\}$$

$$C_1 = \{3\}, C_2 = \{1, 2, 4, 5\}$$

$$C_1 = \{4\}, C_2 = \{1, 2, 3, 5\}$$

$$C_1 = \{5\}, C_2 = \{1, 2, 3, 4\}$$

$$C_1 = \{1, 2\}, C_2 = \{3, 4, 5\} \quad C_1 = \{1, 3\}, C_2 = \{2, 4, 5\}$$

$$C_1 = \{1, 4\}, C_2 = \{2, 3, 5\} \quad C_1 = \{1, 5\}, C_2 = \{2, 3, 4\}$$

$$C_1 = \{2, 3\}, C_2 = \{1, 4, 5\} \quad C_1 = \{2, 4\}, C_2 = \{1, 3, 5\}$$

$$C_1 = \{2, 5\}, C_2 = \{1, 3, 4\} \quad C_1 = \{3, 4\}, C_2 = \{1, 2, 5\}$$

$$C_1 = \{3, 5\}, C_2 = \{1, 2, 4\} \quad C_1 = \{4, 5\}, C_2 = \{1, 2, 3\}$$

2. 可能的划分？

第二类 Stirling 数： 集合的一个划分，表示将 n 个不同的元素拆分成 k 个集合的方案数，记为 $S(n, k)$ 。

$$S(n, k) = \frac{1}{k!} \sum_{i=0}^k (-1)^i \binom{k}{i} (k-i)^n$$

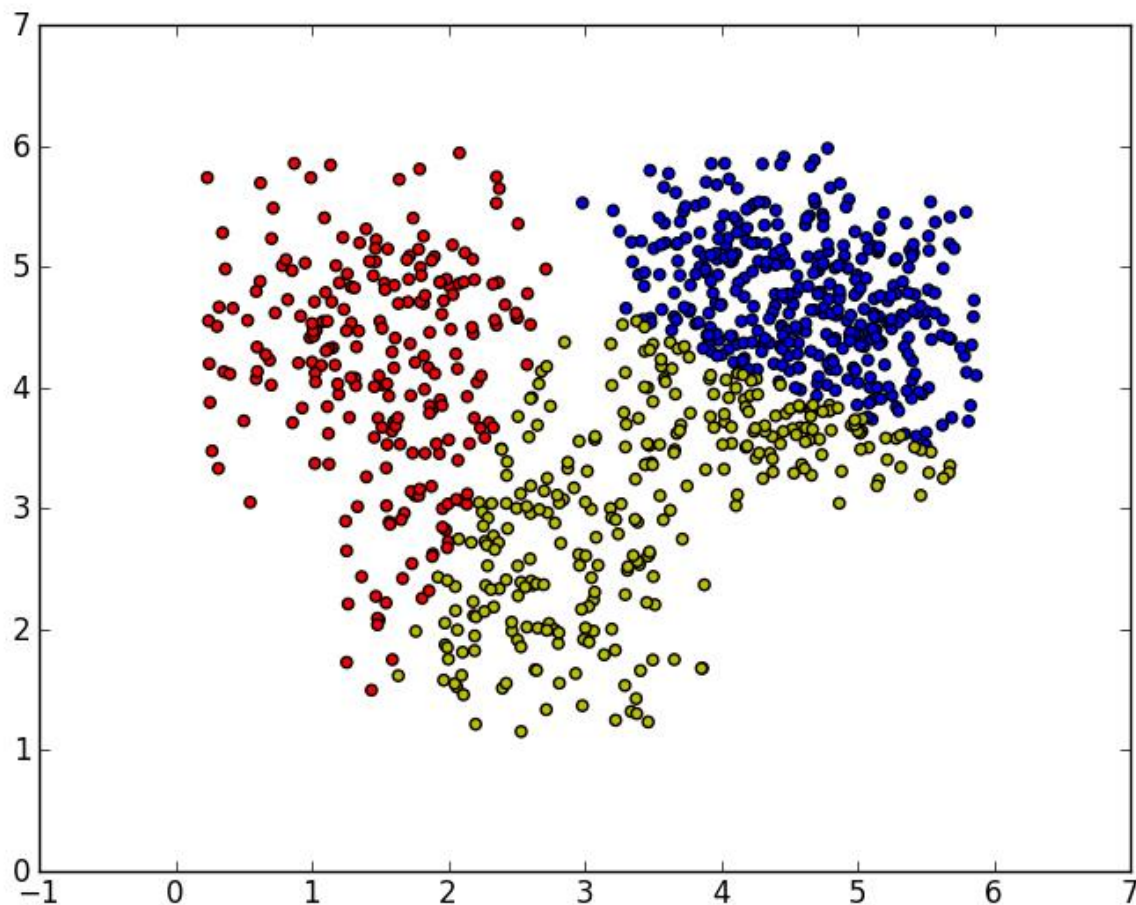
随着集合元素 n 和子集个数 k 的增加，方案数 $S(n, k)$ 呈爆炸式的增长！！！！

2. 可能的划分？

帕斯卡三角形：

n=0	1
n=1	0 1
n=2	0 1 1
n=3	0 1 3 1
n=4	0 1 7 6 1
n=5	0 1 15 25 10 1
n=6	0 1 31 90 65 15 1
n=7	0 1 63 301 350 140 21 1
n=8	0 1 127 966 1701 1050 266 28 1
n=9	0 1 255 3025 7770 6951 2646 462 36 1

3. 聚类个数？



聚类问题

聚类：根据某种相似性，把一组数据划分成若干个簇的过程。

难点一：相似性很难精准定义！

--- 各种距离，度量学习。

难点二：可能存在的划分太多！

--- 避免穷举，优化算法。

难点三：若干个簇 = ？

--- 预先给定，算法自适应。

7.2 K-means 算法

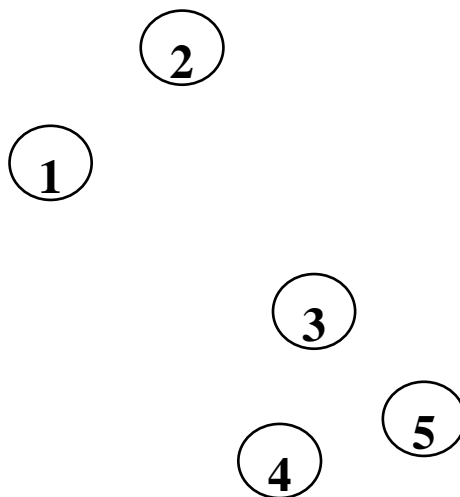
陈飞宇

fchen@cqu.edu.cn

办公室：软件学院529

人造例子

如何对下面数据进行聚类？



相似性： 欧式距离

簇个数： $k = 2$

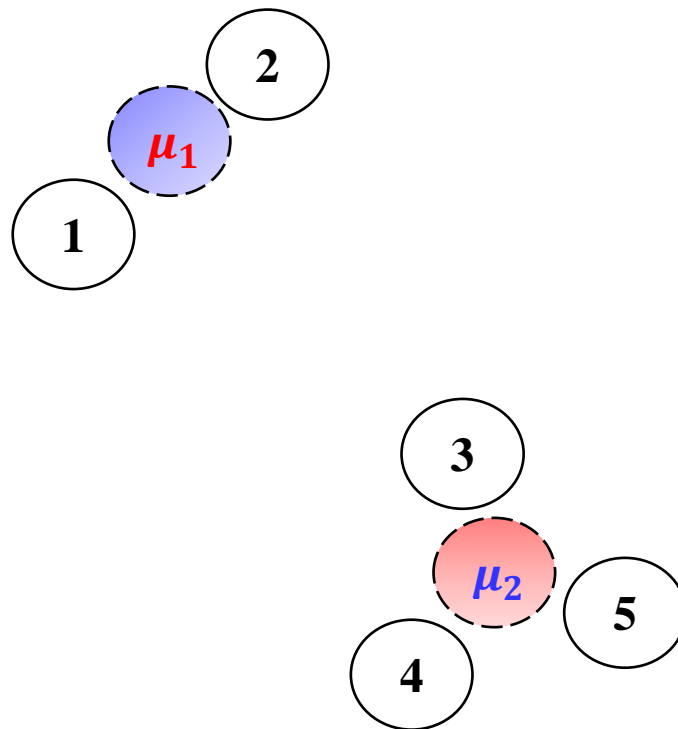
最优化分： $C_1 = \{1, 2\}$, $C_2 = \{3, 4, 5\}$

簇的中心

聚类问题可以通过为每个簇寻找合适的中心来实现。

假设每个簇的中心已经找到，可以把所有数据点分配到距离它最近的中心所在的簇。

$$j = \operatorname{argmin}_l \mathit{dist}(x_i, \mu_l)$$

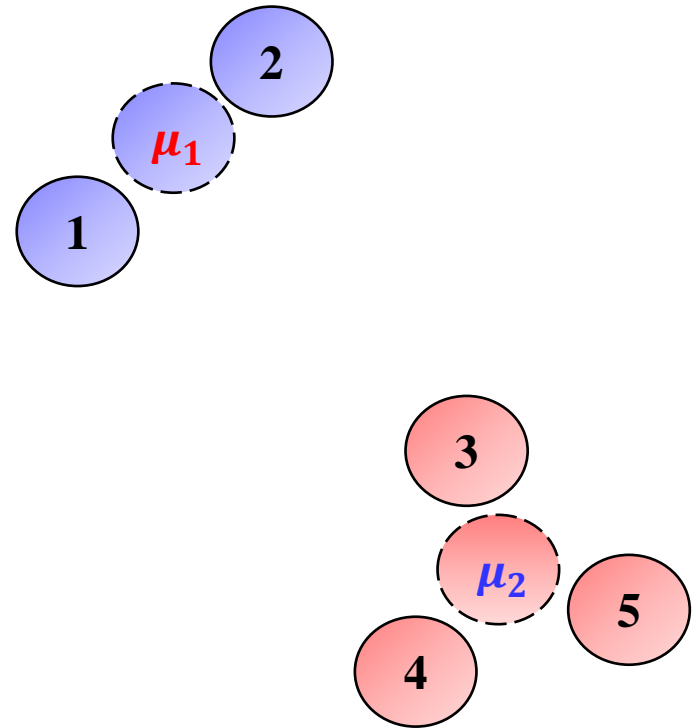


簇的中心

聚类问题可以通过为每个簇寻找合适的中心来实现。

假设每个簇的中心已经找到，可以把所有数据点分配到距离它最近的中心所在的簇。

$$j = \operatorname{argmin}_l \operatorname{dist}(x_i, \mu_l)$$



K-means

给定数据 x 以及簇的个数 k ，K-means 模型可以写成下面的优化模型：

$$\operatorname{argmin}_{\mu_i, C_i} \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$$

变量：

- C_i ：第 i 个簇。
- μ_i ：第 i 个簇的中心。

K-means

目标函数对中心 μ_i 求偏导，我们有

$$\frac{\partial (\sum_{i=1}^k \sum_{x \in \mathcal{C}_i} \|x - \mu_i\|^2)}{\partial \mu_i} = -2 \sum_{x \in \mathcal{C}_i} (x - \mu_i)$$

即

$$\mu_i = \frac{1}{|\mathcal{C}_i|} \sum_{x \in \mathcal{C}_i} x$$

因此，上述模型通常被称作 **K-means**。

K-means

K-means 模型:

$$\operatorname{argmin}_{\mathcal{C}_i} \sum_{i=1}^k \sum_{x \in \mathcal{C}_i} \left\| x - \frac{1}{|\mathcal{C}_i|} \sum_{x \in \mathcal{C}_i} x \right\|^2$$

可能的划分数:

$$S(n, k) = \frac{1}{k!} \sum_{i=0}^k (-1)^i \binom{k}{i} (k-i)^n$$

非凸组合优化问题，NP-难!!!

非凸优化问题

求解非凸组合优化问题的两种常见方法：

- **启发式方法 (Heuristic method)**：一个基于直观或经验构造的算法，在可接受的时间内下给出待解决组合优化问题每一个实例的一个可行解，该可行解与最优解的偏离程度一般不能被预计。
- **松弛方法 (Relaxation method)**：对组合优化问题进行适当的松弛，将其转化为多项式时间内可解的优化问题，松弛后问题的解不是原组合优化问题的解，需要适当的后处理。

Lloyd 算法 I

给定数据 x 以及簇的个数 k :

初始化: 随机选取 k 个簇的中心 $\{\mu_i\}_{i=1,\dots,k}$

重复下面迭代过程直到收敛:

- **划分步骤:** 对于每一个数据点 x_j , 计算其应该属于的簇

$$\operatorname{argmin}_i \|x_j - \mu_i\|_2^2$$

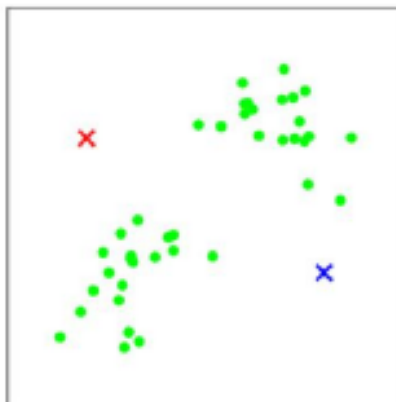
- **更新步骤:** 重新计算每个簇的中心

$$\mu_i = \frac{1}{|C_i|} \sum_{x_j \in C_i} x_j$$

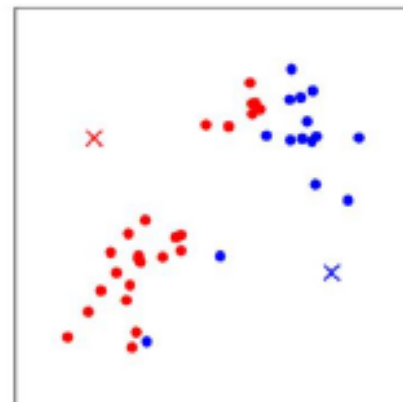
Lloyd 算法 II



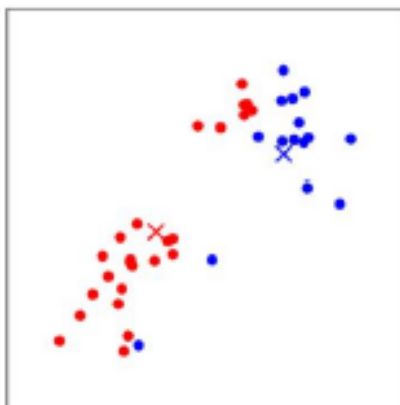
(a)



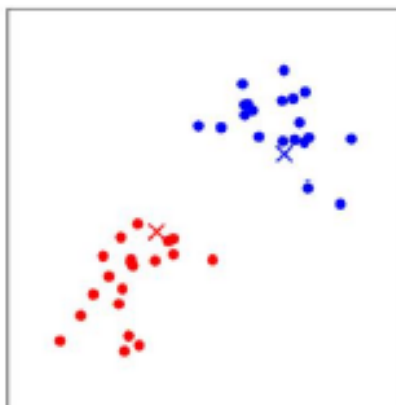
(b)



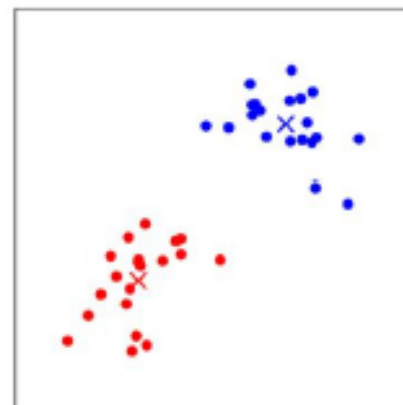
(c)



(d)



(e)



(f)

Lloyd 算法 III

优势:

- Lloyd 算法属于EM算法（期望最大化），可以保证收敛到K-means问题的局部最优解。
- Lloyd 算法的速度快，计算复杂度为 $O(nk)$ 。
- Lloyd 算法思想简单，容易实现，可拓展性强。

劣势:

- 簇的个数 k 需要预先给定。
- 聚类结果依赖于初值的选取。

7.3 K-means进阶（选读）

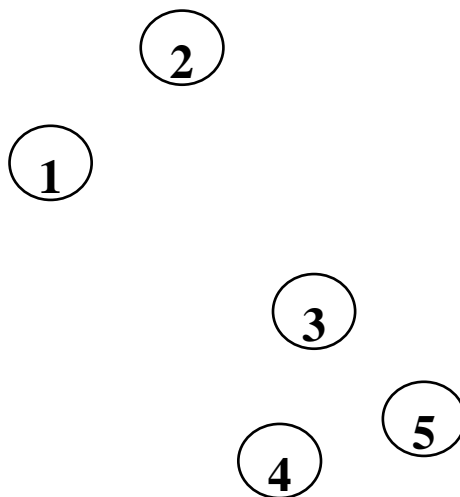
陈飞宇

fchen@cqu.edu.cn

办公室：软件学院529

人造例子

如何对下面数据进行聚类？



相似性： 欧式距离

簇个数： $k = 2$

最优化分： $C_1 = \{1, 2\}$, $C_2 = \{3, 4, 5\}$

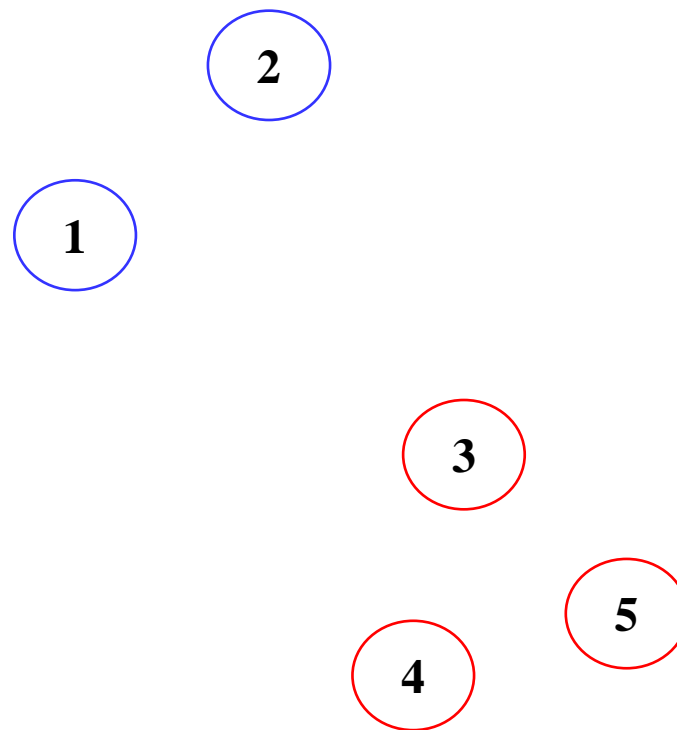
人造例子

示性矩阵 (Indicator matrix) :

$$H = \begin{pmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \\ \mathbf{0} & \mathbf{1} \\ \mathbf{0} & \mathbf{1} \end{pmatrix} \in \mathbb{R}^{n \times k}$$

其中:

$$H_{ij} = 1 \iff i \in \mathcal{C}_j$$



人造例子

标准化示性矩阵 **Normalized Indicator matrix:**

$$H = \begin{pmatrix} 1/\sqrt{2} & 0 \\ 1/\sqrt{2} & 0 \\ 0 & 1/\sqrt{3} \\ 0 & 1/\sqrt{3} \\ 0 & 1/\sqrt{3} \end{pmatrix} \in \mathbb{R}^{n \times k} \quad H_{ij} = 1/\sqrt{|\mathcal{C}_j|} \Leftrightarrow i \in \mathcal{C}_j$$

所有**标准化示性矩阵**的集合:

$$\mathcal{H} = \{H \in \mathbb{R}^{n \times k} \mid H^T H = I, H \geq 0, HH^T \mathbf{1}_n = \mathbf{1}_n\}$$

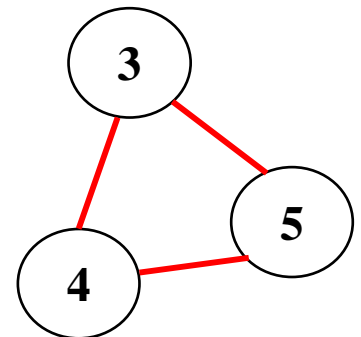
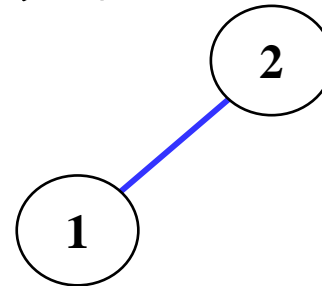
人造例子

邻接矩阵 (Adjacency matrix) :

$$W = \begin{pmatrix} \mathbf{1} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{1} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{1} & \mathbf{1} \\ \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{1} & \mathbf{1} \\ \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{1} & \mathbf{1} \end{pmatrix} \in \mathbb{R}^{n \times n}$$

其中:

$$W_{ij} = 1 \iff i, j \in \mathcal{C}_l \text{ for some } l.$$



人造例子

均一化邻接矩阵 (Adjacency matrix) :

$$W = \begin{pmatrix} \mathbf{1/2} & \mathbf{1/2} & 0 & 0 & 0 \\ \mathbf{1/2} & \mathbf{1/2} & 0 & 0 & 0 \\ 0 & 0 & \mathbf{1/3} & \mathbf{1/3} & \mathbf{1/3} \\ 0 & 0 & \mathbf{1/3} & \mathbf{1/3} & \mathbf{1/3} \\ 0 & 0 & \mathbf{1/3} & \mathbf{1/3} & \mathbf{1/3} \end{pmatrix} \in \mathbb{R}^{n \times n}$$

$$W_{ij} = 1/|C_l| \Leftrightarrow i, j \in C_l \text{ for some } l.$$

所有均一化邻接矩阵的集合:

$$\mathcal{W} = \{W \in \mathbb{R}^{n \times n} | W \mathbf{1}_n = \mathbf{1}_n, \mathbf{Tr}(W) = k,$$

$$W \geq 0, W^T = W, W^2 = W\}$$

人造例子

关系:

$$HH^T = \begin{pmatrix} 1/\sqrt{2} & 0 \\ 1/\sqrt{2} & 0 \\ 0 & 1/\sqrt{3} \\ 0 & 1/\sqrt{3} \\ 0 & 1/\sqrt{3} \end{pmatrix} \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 & 0 & 0 \\ 0 & 0 & 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \end{pmatrix}$$

$$= \begin{pmatrix} 1/2 & 1/2 & 0 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 & 0 \\ 0 & 0 & 1/3 & 1/3 & 1/3 \\ 0 & 0 & 1/3 & 1/3 & 1/3 \\ 0 & 0 & 1/3 & 1/3 & 1/3 \end{pmatrix} = W$$

$$H \in \mathcal{H} \quad \Longleftrightarrow \quad HH^T \in \mathcal{W}$$

人造例子

划分: $c_1 = \{1, 2\}$, $c_2 = \{3, 4, 5\}$

标准化示性矩阵: $H = \begin{pmatrix} 1/\sqrt{2} & 0 \\ 1/\sqrt{2} & 0 \\ 0 & 1/\sqrt{3} \\ 0 & 1/\sqrt{3} \\ 0 & 1/\sqrt{3} \end{pmatrix}$

均一化邻接矩阵: $W = \begin{pmatrix} 1/2 & 1/2 & 0 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 & 0 \\ 0 & 0 & 1/3 & 1/3 & 1/3 \\ 0 & 0 & 1/3 & 1/2 & 1/3 \\ 0 & 0 & 1/3 & 1/3 & 1/3 \end{pmatrix}$

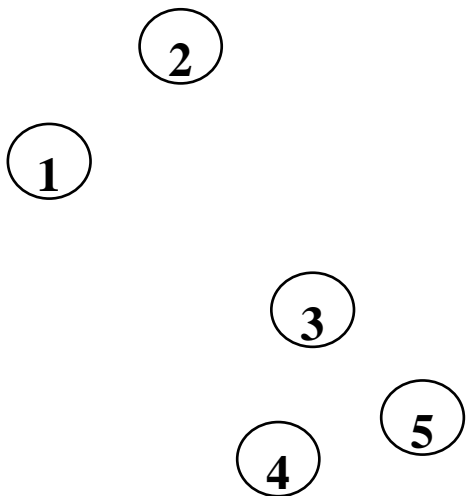
小结

对于任意一种划分都对应一个标准化示性矩阵 H （或均一化邻接矩阵 W ）。反之，任意一个标准化示性矩阵 H （或均一化邻接矩阵 W ）都对应着一种划分。

本质上说，标准化示性矩阵 H （或均一化邻接矩阵 W ）的引入并没有简化聚类问题的难度，但是为求解方法提供了更多的参考和选择。

人造例子

如何对下面数据进行聚类？



划分: $c_1 = \{1, 2\}$, $c_2 = \{3, 4, 5\}$

标准化示性矩阵:

$$H = \begin{pmatrix} 1/\sqrt{2} & 0 \\ 1/\sqrt{2} & 0 \\ 0 & 1/\sqrt{3} \\ 0 & 1/\sqrt{3} \\ 0 & 1/\sqrt{3} \end{pmatrix}$$

K-means

对于给定的数据矩阵 X ，标准化示性矩阵 H ，我们有

$$\begin{aligned} XHH^T &= (x_1 \ x_2 \ x_3 \ x_4 \ x_5) \begin{pmatrix} 1/2 & 1/2 & 0 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 & 0 \\ 0 & 0 & 1/3 & 1/3 & 1/3 \\ 0 & 0 & 1/3 & 1/2 & 1/3 \\ 0 & 0 & 1/3 & 1/3 & 1/3 \end{pmatrix} \\ &= \left(\frac{x_1 + x_2}{2} \quad \frac{x_1 + x_2}{2} \quad \frac{x_3 + x_4 + x_5}{3} \quad \frac{x_3 + x_4 + x_5}{3} \quad \frac{x_3 + x_4 + x_5}{3} \right) \\ &= (\mu_1 \quad \mu_1 \quad \mu_2 \quad \mu_2 \quad \mu_2) \end{aligned}$$

因此， $\|X - XHH^T\|_F^2 = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$

K-means

K-means 模型:

$$\operatorname{argmin}_{\mu_i, \mathcal{C}_i} \sum_{i=1}^k \sum_{x \in \mathcal{C}_i} \|x - \mu_i\|^2$$

等价于

$$\operatorname{argmin}_{H \in \mathcal{H}} \|X - XHH^T\|_F^2$$

或者

$$\operatorname{argmin}_{W \in \mathcal{W}} \|X - XW\|_F^2$$

把原来的组合优化问题改写成矩阵优化问题！

K-means

如果正交约束 $H^T H = I$ 成立，我们有

$$\begin{aligned} & \|X - XHH^T\|_F^2 \\ &= \text{Tr}\left((X - XHH^T)(X - XHH^T)^T\right) \\ &= \text{Tr}\left((X - XHH^T)(X^T - HH^T X^T)\right) \\ &= \text{Tr}(XX^T) + \text{Tr}(XH[H^T H]H^T X^T) \\ &\quad - \text{Tr}(XHH^T X^T) - \text{Tr}(X^T XHH^T) \\ &= \text{Tr}(XX^T) - \text{Tr}(XHH^T X^T) \end{aligned}$$

K-means

$$\operatorname{argmin}_{H \in \mathcal{H}} \|X - XHH^T\|_F^2 \Leftrightarrow \operatorname{argmax}_{H \in \mathcal{H}} \operatorname{Tr}(X^T X H H^T)$$

$$\mathcal{H} = \{H \in \mathbb{R}^{n \times k} \mid H^T H = I, H \geq 0, H H^T \mathbf{1}_n = \mathbf{1}_n\}$$

$$\operatorname{argmin}_{W \in \mathcal{W}} \|X - XW\|_F^2 \Leftrightarrow \operatorname{argmax}_{W \in \mathcal{W}} \operatorname{Tr}(X^T X W)$$

$$\mathcal{W} = \{W \in \mathbb{R}^{n \times n} \mid W \mathbf{1}_n = \mathbf{1}_n, \operatorname{Tr}(W) = k, \\ W \geq 0, W^T = W, W^2 = W\}$$

K-means

松弛模型 I: PCA!!!

$$\operatorname{argmax}_{H^T H = I} \operatorname{Tr}(H^T X^T X H)$$

松弛模型 II: NMF!!!

$$\operatorname{argmin}_{H \geq 0} \|X - X H H^T\|_F^2$$

**PCA和NMF都可以得到K-means模型的松弛解，
但是松弛解并不满足聚类需要!!!**

聚类任务小结

我们描述了聚类任务：与分类任务的区别，问题的难点。

我们学习了聚类算法：K-means 模型 + Lloyd 算法 + **K-means++初始化**。

我们利用**示性矩阵和邻接矩阵**推导了 K-means 模型的两种等价形式。

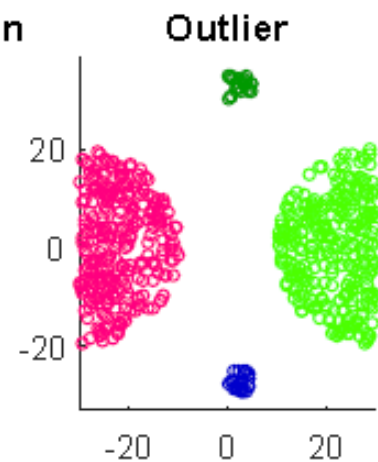
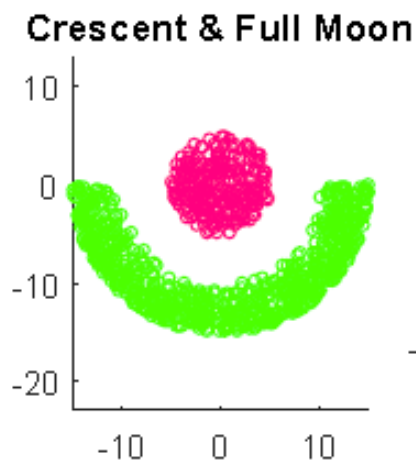
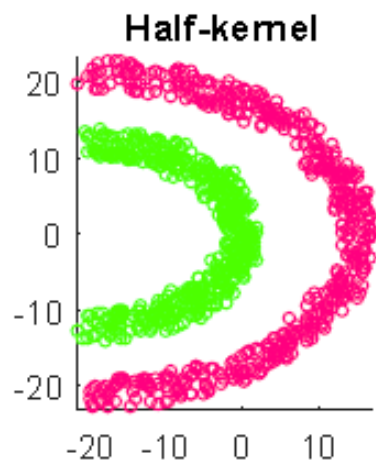
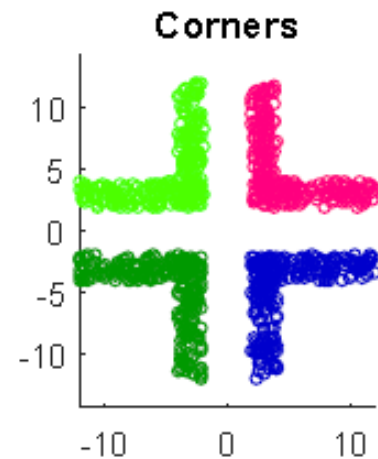
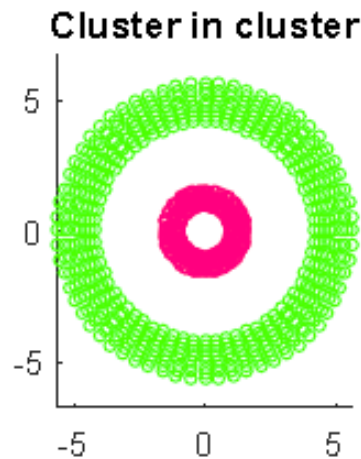
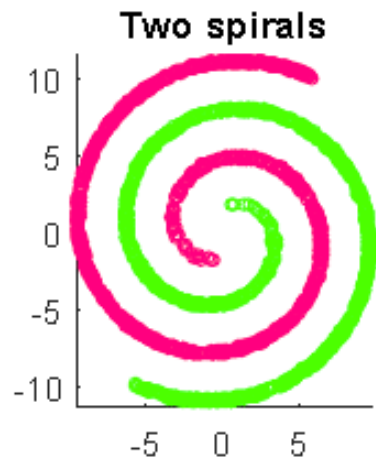
7.4 谱聚类

陈飞宇

fchen@cqu.edu.cn

办公室：软件学院529

Non-globular Clustering



Non-globular Clustering

解决方案：核化 K-means ?

$$\operatorname{argmin}_{\mu_i, \mathcal{C}_i} \sum_{i=1}^k \sum_{x \in \mathcal{C}_i} \|x - \mu_i\|^2$$

- K-means 模型仅和数据点 x 有关！
- Lloyd 算法对核矩阵选取中心？

结论：K-means 模型 + Lloyd 算法不能直接推广！

Non-globular Clustering

矩阵形式 K-means 模型:

$$\operatorname{argmax} \operatorname{Tr}(H^T X^T X H)$$

$$\text{s.t. } H^T H = I, H \geq 0, H H^T \mathbf{1}_n = \mathbf{1}_n$$

- K-means 模型和数据矩阵 $X^T X$ 有关!!!
- 求解 K-means 和 Kernel K-means 是一样的。
- 近似解法: **PALM**

谱聚类 (Spectral Clustering)

谱 (Spectral) : 矩阵的特征值 !

谱聚类: 根据图论, 把聚类问题转化为求解拉普拉斯矩阵的特征值问题 !

优势: 理论高深, 实现简单, 可适用于各种形状的数据 !

图论基础知识

通常用 $G(V, E)$ 表示一个图， V 中的元素称为节点， E 中的元素称为边。

若节点 v_i 和 v_j 之间有边相连，可以给其对应的边 e_{ij} 赋予一个非负的权重 $w_{ij} > 0$ 。反之，若 $w_{ij} = 0$ ，则说明节点 v_i 和 v_j 之间没有边相连。

节点 v_i 的度 d_i ： $d_i = \sum_{j=1}^n w_{ij}$ 。

谱聚类 (Spectral Clustering)

输入：数据矩阵 \mathbf{X} ，簇个数 k 。

- 构建邻接矩阵 \mathbf{W} 。
- 计算拉普拉斯矩阵 \mathbf{L} 。
- 计算 \mathbf{L} 最小 k 个特征值对应的特征向量 \mathbf{U} 。
- 对 \mathbf{U} 进行 **K-means** 聚类。

输出： k 个簇。

邻接矩阵 (Adjacency Matrix)

构造邻接矩阵 W 的方法:

- ε - 邻域: 连接所有距离小于 ε 的点。
- k - 近邻: 把每个点与它最近的 k 个点相连。
(可能非对称, 对称化或相互最近邻)
- 全连接图: $w_{ij} = \exp(-\|x_i - x_j\|^2 / (2\sigma^2))$
(其中 σ 控制邻域的宽度)

拉普拉斯矩阵

拉普拉斯矩阵:

$$\mathbf{L} = \mathbf{D} - \mathbf{W}$$

其中 $\mathbf{D} = \text{diag}(d_1, d_2, \dots, d_n)$ 是所有数据点的度矩阵。

Normalized 拉普拉斯矩阵:

$$\mathbf{L}_{rw} = \mathbf{D}^{-1}\mathbf{L} = \mathbf{I} - \mathbf{D}^{-1}\mathbf{W}$$

$$\mathbf{L}_{sym} = \mathbf{D}^{-1/2}\mathbf{L}\mathbf{D}^{-1/2} = \mathbf{I} - \mathbf{D}^{-1/2}\mathbf{W}\mathbf{D}^{-1/2}$$

特征值问题

特征值问题 I :

$$Lu = \lambda u$$

特征值问题 II :

$$L_{rw}u = D^{-1}Lu = \lambda u$$

$$Lu = \lambda Du$$

特征值问题 III :

$$L_{sym}u = \lambda u$$

谱聚类 (Spectral Clustering)

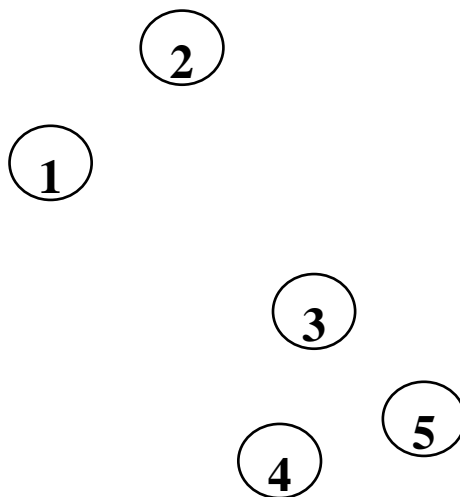
输入：数据矩阵 \mathbf{X} ，簇个数 k 。

- 构建邻接矩阵 \mathbf{W} 。
- 计算拉普拉斯矩阵 \mathbf{L} 。
- 计算 \mathbf{L} 最小 k 个特征值对应的特征向量 \mathbf{U} 。
- 对 \mathbf{U} 进行 **K-means** 聚类。

输出： k 个簇。

人造例子

如何对下面数据进行聚类？



相似性： 欧式距离

簇个数： $k = 2$

最优化分： $c_1 = \{1, 2\}$, $c_2 = \{3, 4, 5\}$

人造例子

邻接矩阵与度矩阵：

$$W = \begin{pmatrix} \mathbf{1} & \mathbf{1} & 0 & 0 & 0 \\ \mathbf{1} & \mathbf{1} & 0 & 0 & 0 \\ 0 & 0 & \mathbf{1} & \mathbf{1} & \mathbf{1} \\ 0 & 0 & \mathbf{1} & \mathbf{1} & \mathbf{1} \\ 0 & 0 & \mathbf{1} & \mathbf{1} & \mathbf{1} \end{pmatrix}$$

$$D = \begin{pmatrix} 2 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 3 \end{pmatrix}$$

人造例子

拉普拉斯矩阵:

$$L = \begin{pmatrix} 2 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 3 \end{pmatrix} - \begin{pmatrix} 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 & 1 \end{pmatrix}$$
$$= \begin{pmatrix} 1 & -1 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & -1 & -1 \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & -1 & 2 \end{pmatrix}$$

人造例子

特征向量:

$$Lu_1 = \begin{pmatrix} 1 & -1 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & -1 & -1 \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & -1 & 2 \end{pmatrix} \begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

且

$$Lu_2 = \begin{pmatrix} 1 & -1 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & -1 & -1 \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & -1 & 2 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 1/\sqrt{3} \\ 1/\sqrt{3} \\ 1/\sqrt{3} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

人造例子

聚类:

$$U = (u_1, u_2)^T = \begin{pmatrix} \mathbf{1/\sqrt{2}} & \mathbf{1/\sqrt{2}} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{1/\sqrt{3}} & \mathbf{1/\sqrt{3}} & \mathbf{1/\sqrt{3}} \end{pmatrix}$$

前2个、后3个分别重合！

问题？

1. 相比直接进行 K-means 聚类，谱聚类有什么不同？

相当于先做了一次特征提取，再聚类！

2. 谱聚类背后的机理是什么？

图论中的切割图问题！

拉普拉斯矩阵

拉普拉斯矩阵 L :

- 对于任意一个向量 f , 我们有

$$f' L f = \frac{1}{2} \sum_{i,j=1}^n w_{ij} (f_i - f_j)^2$$

- L 是一个对称半正定矩阵。
- L 的最小特征值是0, 其对应的特征向量是 $\mathbf{1}_n$
- L 有 n 个非负的实特征值。

拉普拉斯矩阵

因为 $L = D - W$; $d_i = \sum_{j=1}^n w_{ij}$, 我们有

$$\begin{aligned} f' L f &= f' D f - f' W f = \sum_{i=1}^n d_i f_i^2 - \sum_{i,j=1}^n w_{ij} f_i f_j \\ &= \frac{1}{2} \left(\sum_{i=1}^n d_i f_i^2 - 2 \sum_{i,j=1}^n w_{ij} f_i f_j + \sum_{j=1}^n d_j f_j^2 \right) \\ &= \frac{1}{2} \sum_{i,j=1}^n w_{ij} (f_i - f_j)^2 \end{aligned}$$

拉普拉斯矩阵的谱 vs 图的连同分支

定理：假设 $G(V, E)$ 是一个图， L 是它的拉普拉斯矩阵。那么 $G(V, E)$ 的连通分支 A_1, A_2, \dots, A_k 的个数等于 L 的零特征值的重数。并且，零特征值的特征空间由示性向量 $\mathbf{1}_{A_1}, \mathbf{1}_{A_2}, \dots, \mathbf{1}_{A_k}$ 张成。

切割图问题

