

# 法律声明

---

□ 本课件包括演示文稿、示例、代码、题库、视频和声音等内容，小象学院和主讲老师拥有完全知识产权的权利；只限于善意学习者在本课程使用，不得在课程范围外向任何第三方散播。任何其他人或机构不得盗版、复制、仿造其中的创意及内容，我们保留一切通过法律手段追究违反者的权利。

□ 课程详情请咨询

■ 微信公众号：小象

■ 新浪微博：ChinaHadoop



# 聚类(下)

---



小象学院  
ChinaHadoop.cn

邹博

# 本次目标

---

- 理解密度聚类并能够应用于实践
  - DBSCAN
  - DensityPeak密度最大值聚类
- 掌握谱聚类的算法
  - 考虑谱聚类和PCA的关系
- 半监督LPA算法
  - 社区发现

# 密度聚类方法

---

- 密度聚类方法的指导思想是，只要样本点的密度大于某阈值，则将该样本添加到最近的簇中。
- 这类算法能克服基于距离的算法只能发现“类圆形”（凸）的聚类的缺点，可发现任意形状的聚类，且对噪声数据不敏感。但计算密度单元的计算复杂度大，需要建立空间索引来降低计算量。
  - DBSCAN
  - 密度最大值算法

# DBSCAN算法

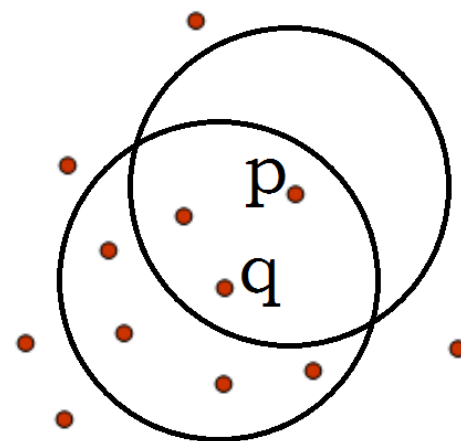
---

- DBSCAN(Density-Based Spatial Clustering of Applications with Noise)
- 一个比较有代表性的基于密度的聚类算法。与划分和层次聚类方法不同，它将簇定义为**密度相连的点的最大集合**，能够把具有足够高密度的区域划分为簇，并可在有“噪声”的数据中发现任意形状的聚类。

# DBSCAN算法的若干概念

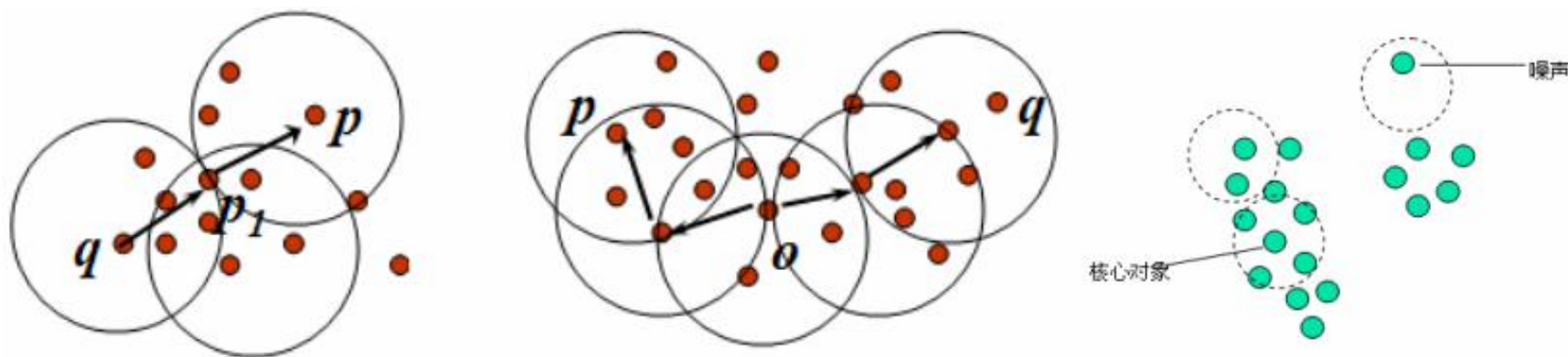
- 对象的 $\varepsilon$ -邻域：给定对象在半径 $\varepsilon$ 内的区域。
- 核心对象：对于给定的数目 $m$ ，如果一个对象的 $\varepsilon$ -邻域至少包含 $m$ 个对象，则称该对象为核心对象。
- 直接密度可达：给定一个对象集合 $D$ ，如果 $p$ 是在 $q$ 的 $\varepsilon$ -邻域内，而 $q$ 是一个核心对象，我们说对象 $p$ 从对象 $q$ 出发是直接密度可达的。

- 如图 $\varepsilon=1\text{cm}$ ， $m=5$ ， $q$ 是一个核心对象，从对象 $q$ 出发到对象 $p$ 是直接密度可达的。

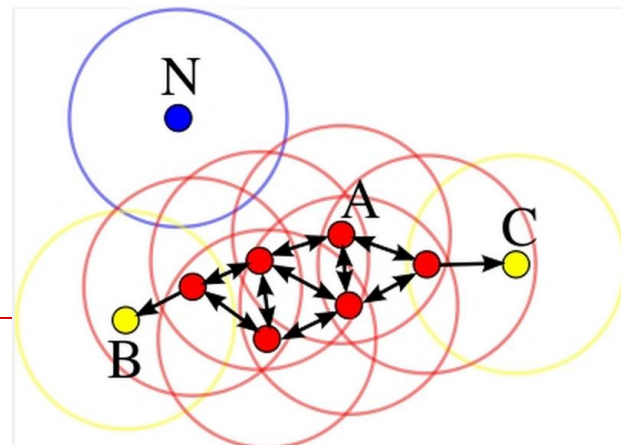


# DBSCAN算法的若干概念

- 密度可达：如果存在一个对象链 $p_1p_2\dots p_n$ ,  $p_1=q$ ,  $p_n=p$ , 对  $p_i \in D$ ,  $(1 \leq i \leq n)$ ,  $p_{i+1}$  是从  $p_i$  关于  $\epsilon$  和  $m$  直接密度可达的, 则对象  $p$  是从对象  $q$  关于  $\epsilon$  和  $m$  密度可达的。
- 密度相连：如果对象集合  $D$  中存在一个对象  $o$ , 使得对象  $p$  和  $q$  是从  $o$  关于  $\epsilon$  和  $m$  密度可达的, 那么对象  $p$  和  $q$  是关于  $\epsilon$  和  $m$  密度相连的。
- 簇：一个基于密度的簇是最大的密度相连对象的集合。
- 噪声：不包含在任何簇中的对象称为噪声。



# DBSCAN算法

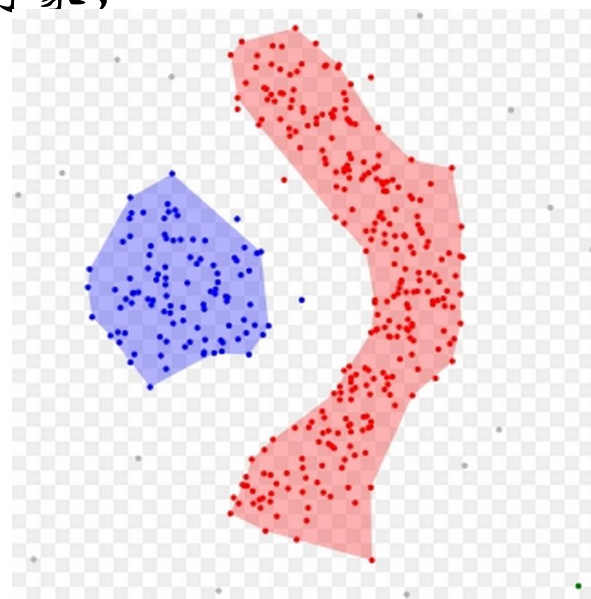


## □ DBSCAN算法流程：

- 如果一个点 $p$ 的 $\epsilon$ -邻域包含多于 $m$ 个对象，则创建一个 $p$ 作为核心对象的新簇；
- 寻找并合并核心对象直接密度可达的对象；
- 没有新点可以更新簇时，算法结束。

## □ 有上述算法可知：

- 每个簇至少包含一个核心对象；
- 非核心对象可以是簇的一部分，构成了簇的边缘(edge)；
- 包含过少对象的簇被认为是噪声。





# Code

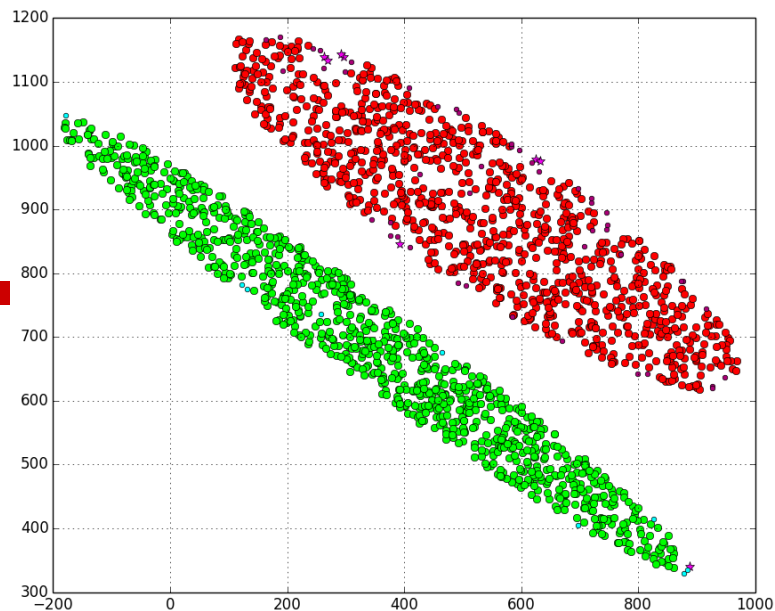
```
def density_cluster(data):
    # 计算核心对象
    m = len(data)          # 样本数目
    sim = [[] for i in range(m)] * m    # sim[i]:第i个样本的近邻
    r2 = r * r
    for i in range(m):     # 距离
        print i
        for j in range(m):
            if distance2(data[i], data[j]) < r2:
                sim[i].append(j)

    # 根据核心对象合并
    uf = uf_init(m)        # 并查集初始化
    for i in range(m):
        if len(sim[i]) > kernel_num: # i是核心对象
            for j in sim[i]:
                uf_union(uf, i, j)

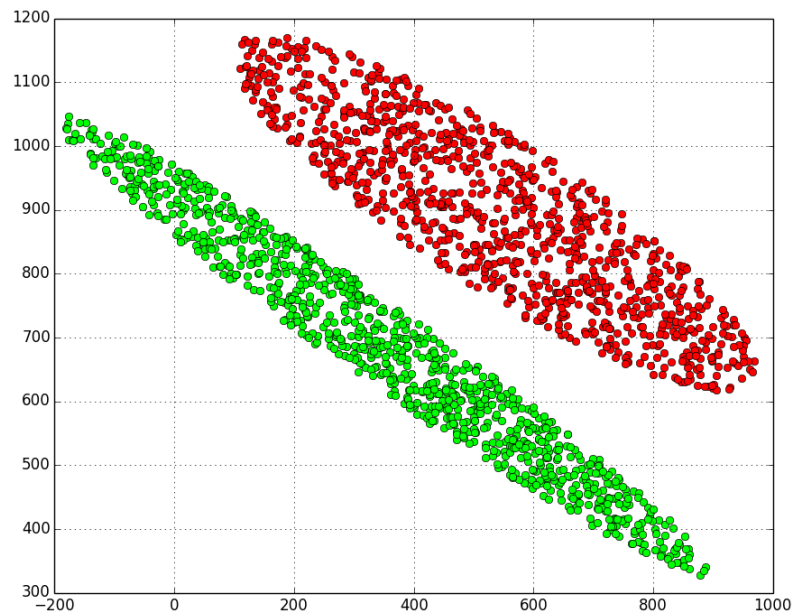
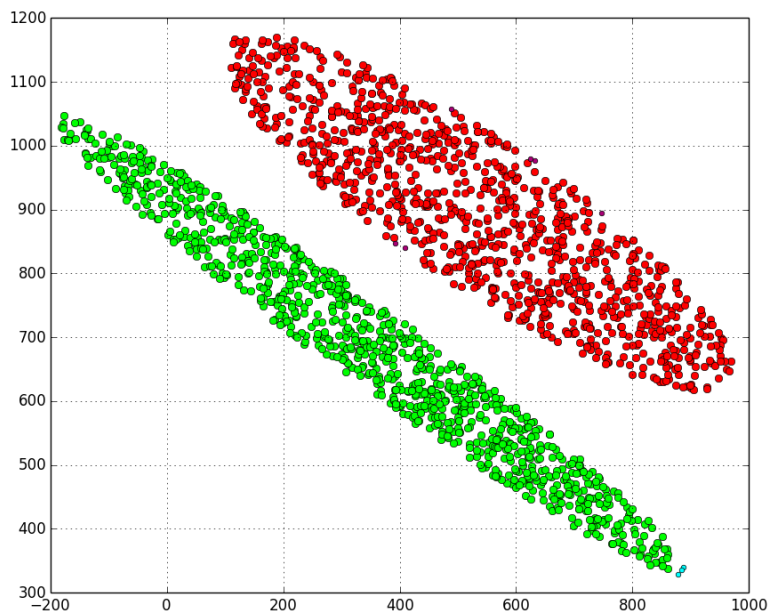
    # 计算并查集的每个对象所属类型
    types = {}
    t = 0
    for i in range(m):
        c = uf_find(uf, i)
        if not types.has_key(c):
            if c != i or len(sim[i]) > kernel_num: # 去除孤立点
                t += 1
            types[c] = t

    # 根据字典映射关系, 计算每个样本所属的簇
    cluster = [0] * m
    for i in range(m):
        c = uf_find(uf, i)
        if (c != i) or (len(sim[i]) > kernel_num):
            if len(sim[i]) > kernel_num: # 是核心对象
                cluster[i] = types[c]
            else:
                cluster[i] = k+types[c]
    return cluster
```

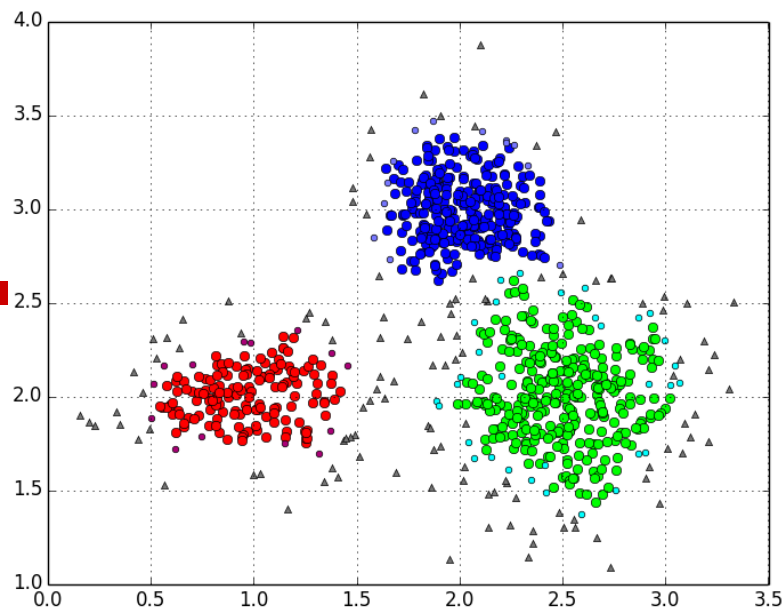
参数:  $n=5$



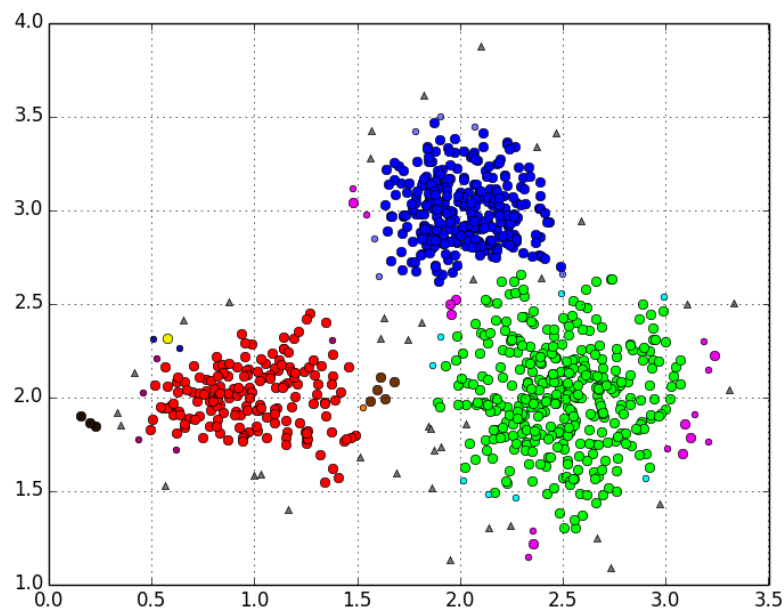
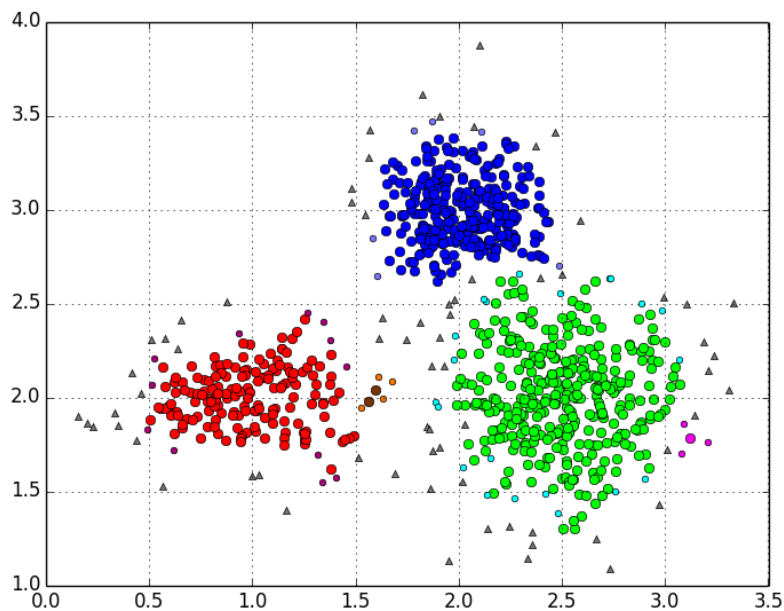
25  
30 | 40



参数:  $r=0.1$



5  
3 | 2



# 密度最大值聚类

□ 密度最大值聚类是一种简洁优美的聚类算法,可以识别各种形状类簇,并且参数很容易确定。

□ 定义: 局部密度  $\rho_i$

$$\rho_i = \sum_j \chi(d_{ij} - d_c), \quad \text{其中, } \chi(x) = \begin{cases} 1 & x < 0 \\ 0 & \text{otherwise} \end{cases}$$

■  $d_c$  是一个截断距离,  $\rho_i$  即到对象  $i$  的距离小于  $d_c$  的对象的个数。由于该算法只对  $\rho_i$  的相对值敏感, 所以对  $d_c$  的选择是稳健的, 一种推荐做法是选择  $d_c$ , 使得平均每个点的邻居数为所有点的 1%-2%

□ 定义: 高局部密度点距离  $\delta_i$

# 局部密度的其他定义

□ 截断值:

$$\rho_i = \sum_j \chi(d_{ij} - d_c), \quad \text{其中, } \chi(x) = \begin{cases} 1 & x < 0 \\ 0 & \text{otherwise} \end{cases}$$

□ 高斯核相似度:

$$\rho_i = \sum_{j \in I_S \setminus \{i\}} \exp\left(-\left(\frac{d_{ij}}{d_c}\right)^2\right)$$

□ K近邻均值:

$$\rho_i = \frac{1}{K} \sum_{j=1}^K d_{ij}, \quad \text{其中, } d_{i1} > d_{i2} > \dots > d_{iK} > d_{i,K+1} > \dots$$

# 高局部密度点距离

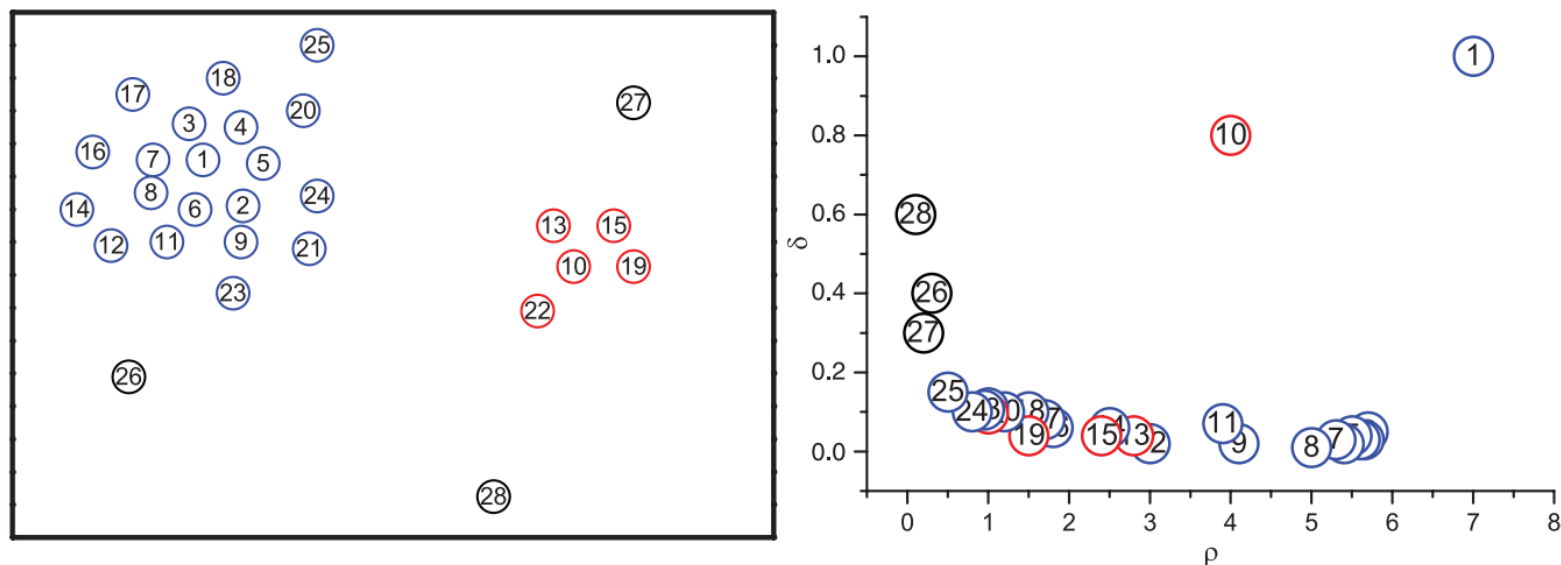
- 高局部密度点距离  $\delta_i = \min_{j: \rho_j > \rho_i} (d_{ij})$
- 在密度高于对象i的所有对象中，到对象i最近的距离，即高局部密度点距离。
- 对于密度最大的对象，设置 $\delta_i = \max(d_{ij})$ (即：该问题中的无穷大)。
  - 只有那些密度是局部或者全局最大的点才会有远大于正常值的高局部密度点距离。

# 簇中心的识别

- 那些有着比较大的局部密度 $\rho_i$ 和很大的高密距离 $\delta_i$ 的点被认为是簇的中心；
- 高密距离 $\delta_i$ 较大但局部密度 $\rho_i$ 较小的点是异常点；
  - 确定簇中心之后，其他点按照距离已知簇的中心最近进行分类
  - 注：也可按照密度可达的方法进行分类。

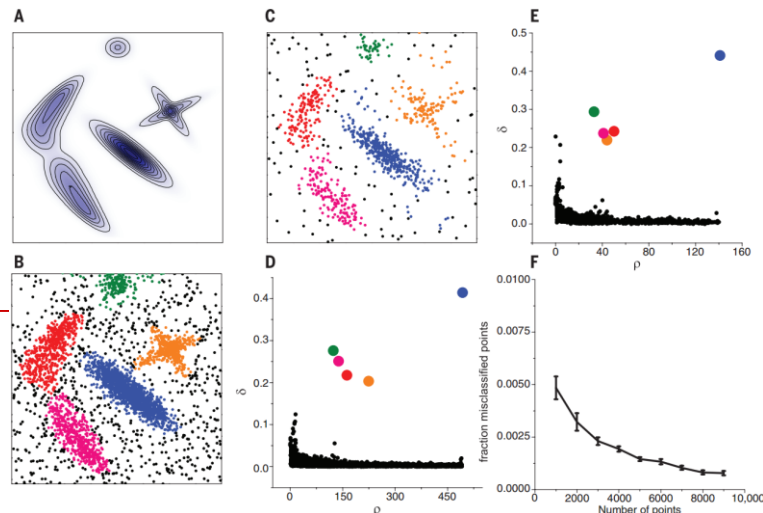
# DensityPeak与决策图Decision Graph

□ 左图是所有点在二维空间的分布, 右图是以 $\rho$ 为横坐标, 以 $\delta$ 为纵坐标绘制的决策图。可以看到, 1和10两个点的 $\rho_i$ 和 $\delta_i$ 都比较大, 作为簇的中心点。26、27、28三个点的 $\delta_i$ 也比较大, 但是 $\rho_i$ 较小, 所以是异常点。

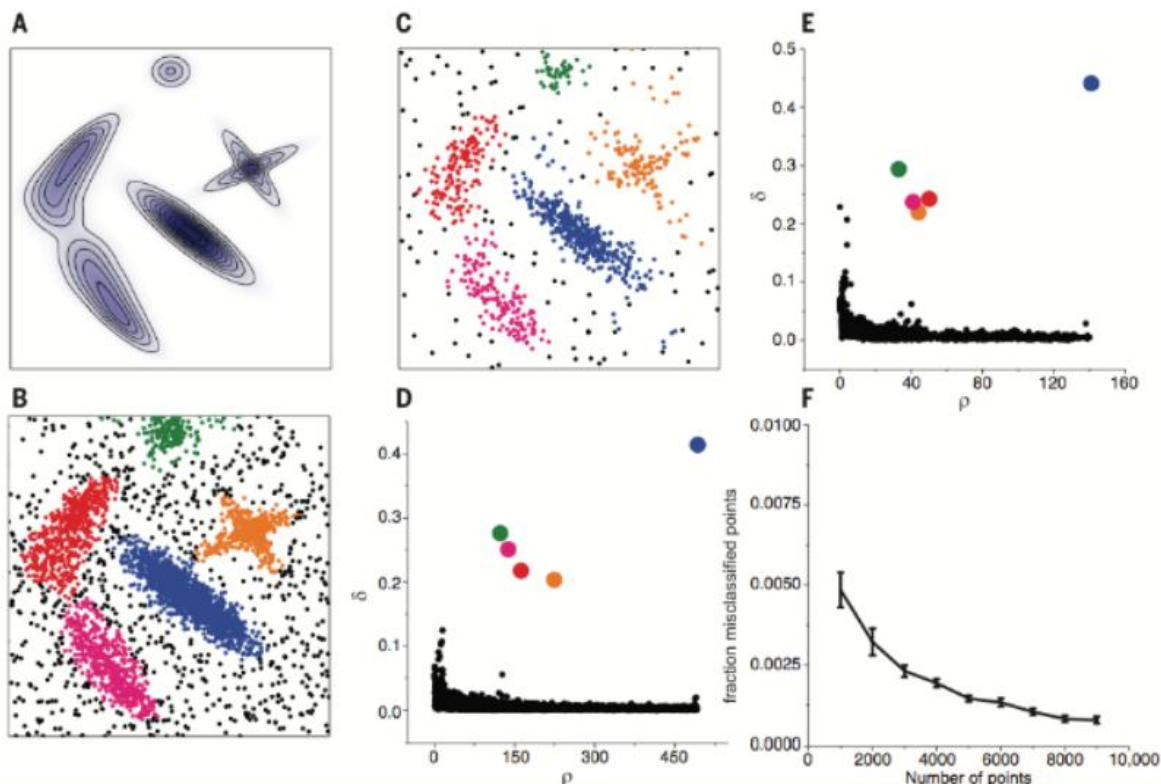




# 边界和噪声的重认识

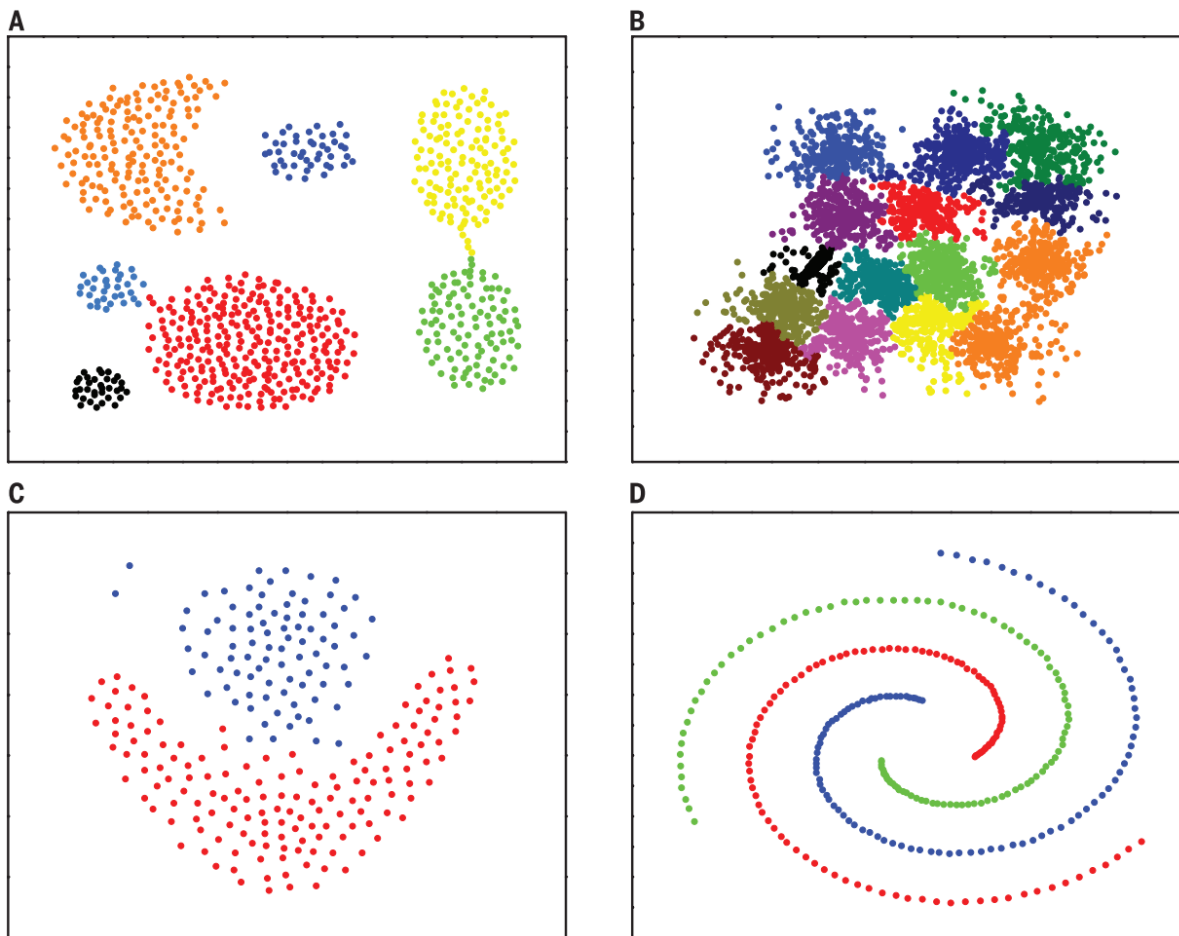


- 在聚类分析中，通常需要确定每个点划分给某个簇的可靠性：
  - 在该算法中，可以首先为每个簇定义一个边界区域 (border region)，亦即划分给该簇但是距离其他簇的点的距离小于  $d_c$  的点的集合。然后为每个簇找到其边界区域的局部密度最大的点，令其局部密度为  $\rho_h$ 。
  - 该簇中所有局部密度大于  $\rho_h$  的点被认为是簇核心的一部分(亦即将该点划分给该类簇的可靠性很大)，其余的点被认为是该类簇的光晕(halo)，亦即可以认为是噪声。
- 注：关于可靠性问题，在EM算法中仍然会有相关涉及。



A图为生成数据的概率分布, B, C二图为分别从该分布中生成了4000, 1000个点. D, E分别是B, C两组数据的决策图(decision tree), 可以看到两组数据都只有五个点有比较大的 $\rho_i$ 和很大的 $\delta_i$ . 这些点作为类簇的中心, 在确定了类簇的中心之后, 每个点被划分到各个类簇(彩色点), 或者是划分到类簇光晕(黑色点). F图展示的是随着抽样点数量的增多, 聚类的错误率在逐渐下降, 说明该算法是鲁棒的.

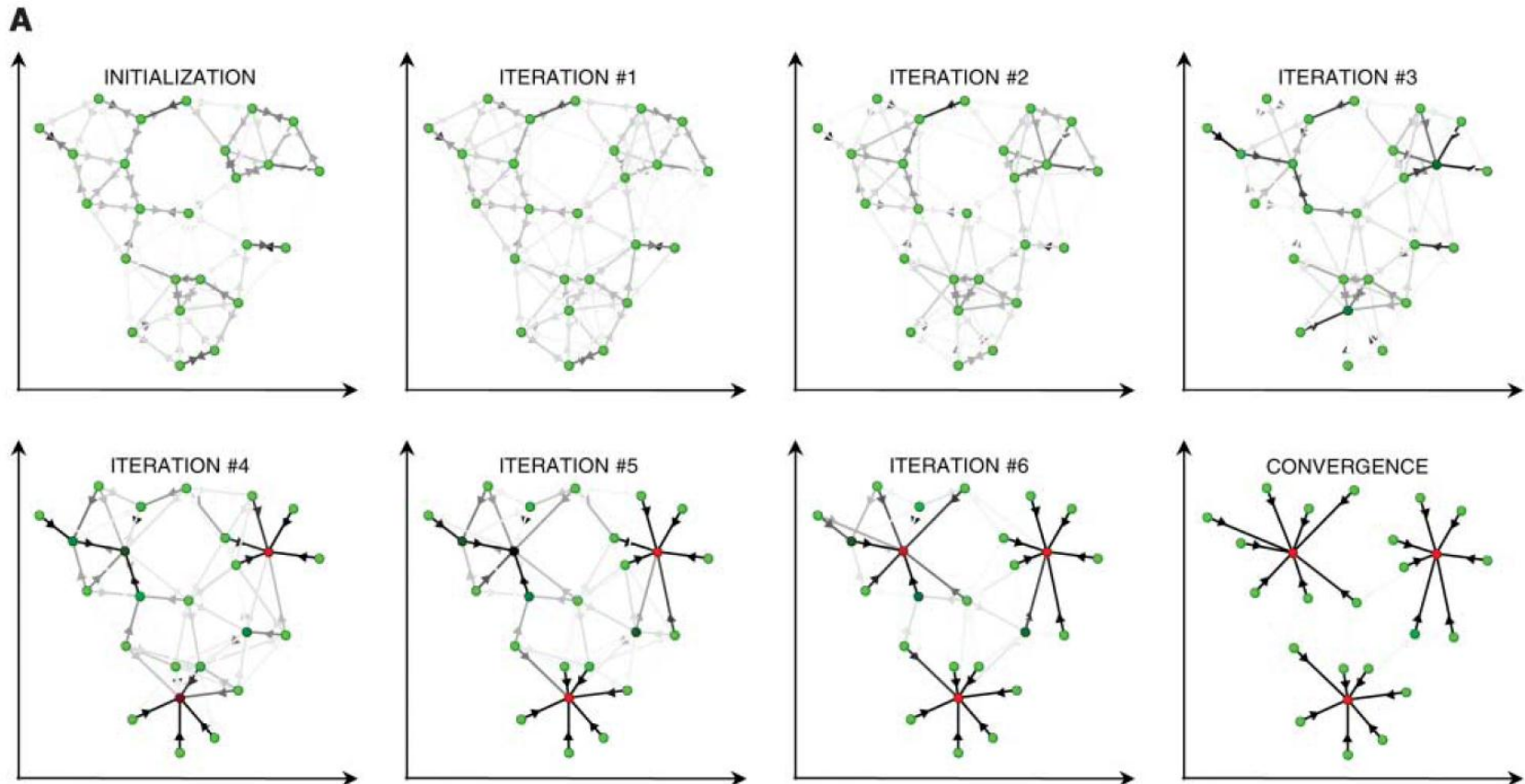
# 不同数据下密度最大值聚类效果



$$r(i,k) \leftarrow s(i,k) - \max_{k' \text{ s.t. } k' \neq k} \{a(i,k') + s(i,k')\}$$

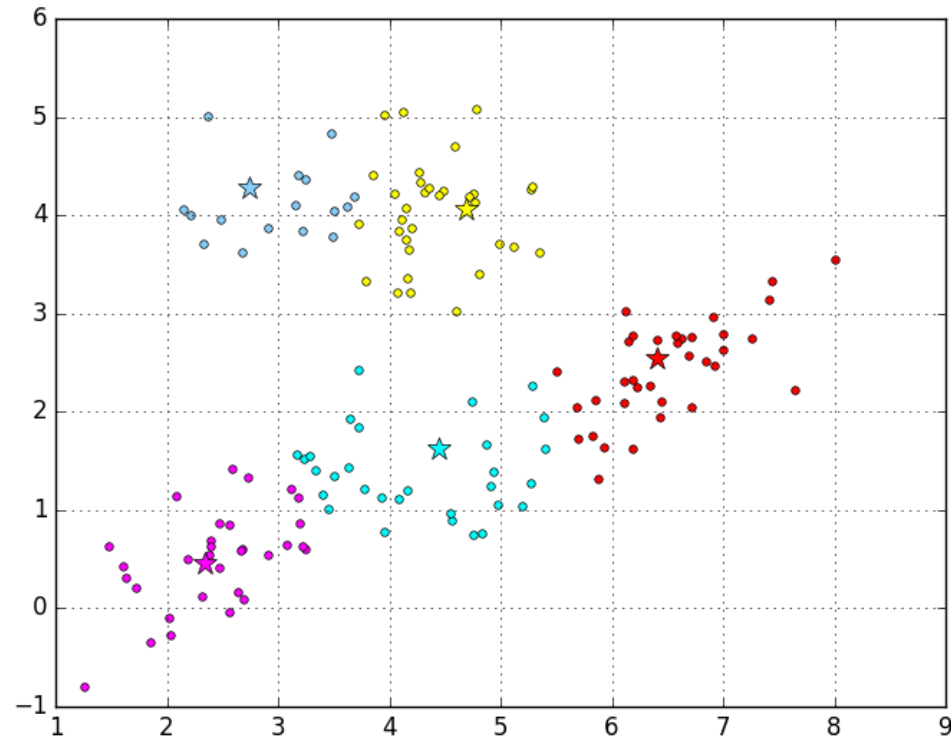
# Affinity Propagation

$$a(i,k) \leftarrow \min \left\{ 0, r(k,k) + \sum_{i' \text{ s.t. } i' \neq i, k} \max \{ 0, r(i',k) \} \right\}$$



# Code

# Cluster

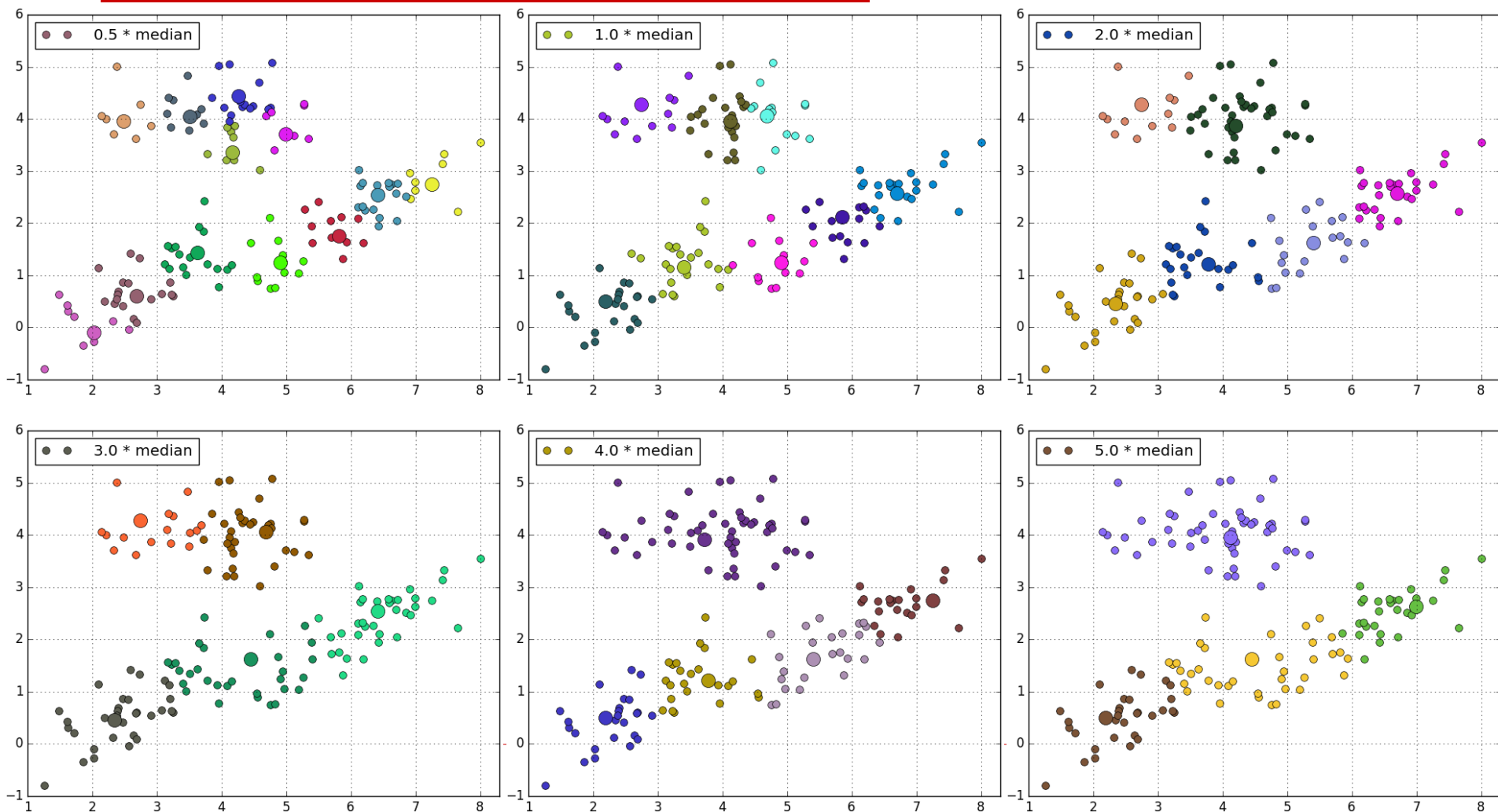


```
def affinity_propagation(data, factor):
    m, n = data.shape  # m样本个数, n样本维度
    s = np.zeros((m, m))
    for i in np.arange(m):
        for j in np.arange(i+1, m):
            s[i][j] = s[j][i] = -((data[i] - data[j]) ** 2).sum()
    p = factor * np.median(s)  # 自聚类因子
    print p
    for i in np.arange(m):
        s[i][i] = p
    r = np.zeros((m, m))  # r(i,k):i对k的依赖度
    a = np.zeros((m, m))  # a(i,k):k度i的适合度
    times = 100  # 迭代次数
    lamda = 0.5  # 阻尼因子
    cluster = np.zeros(m, dtype=np.int)
    center = {}
    for tt in np.arange(times):
        # r
        for i in np.arange(m):
            for k in np.arange(m):
                a_s = None  # a_s: a+s
                for t in np.arange(m):
                    if t != k:
                        if (a_s is None) or (a_s < a[i][t] + s[i][t]):
                            a_s = a[i][t] + s[i][t]
                r[i][k] = lamda * (s[i][k] - a_s) + (1-lamda)*r[i][k]

        # a
        for i in np.arange(m):
            for k in np.arange(m):
                if k == i:  # a[i][i]单独计算
                    continue
                r_s = r[k][k]  # sum(r[:k])
                for t in np.arange(m):
                    if (t != i) and (t != k):
                        r_s += max(r[t][k], 0)
                        # r_s += r[t][k]
                a[i][k] = lamda * min(0, r_s) + (1-lamda)*a[i][k]
        # 计算a[i][i]
        r_s = 0
        for t in np.arange(m):
            if t != i:
                r_s += max(r[t][i], 0)
                # r_s += r[t][i]
        a[i][i] = lamda * r_s + (1-lamda)*a[i][i]
```



# AP算法调参



# 复习：实对称阵的特征值是实数

□ 首先  $A\bar{x} = \overline{A\bar{x}} = \overline{Ax} = \overline{\lambda x} = \bar{\lambda}\bar{x}$

□ 因为  $\bar{x}^T(Ax) = \overline{\bar{x}^T(Ax)} = \bar{x}^T \lambda x = \lambda \bar{x}^T x$

$$\bar{x}^T(Ax) = (\bar{x}^T A^T)x = (A\bar{x})^T x = (\bar{\lambda}\bar{x})^T x = \bar{\lambda}\bar{x}^T x$$

□ 从而

$$\lambda \bar{x}^T x = \bar{\lambda} \bar{x}^T x \Rightarrow (\lambda - \bar{\lambda}) \bar{x}^T x = 0$$

□ 而

$$\bar{x}^T x = \sum_{i=1}^n \overline{x_i} x_i = \sum_{i=1}^n |x_i|^2 \neq 0$$

□ 所以

$$\lambda - \bar{\lambda} = 0 \Rightarrow \lambda = \bar{\lambda}$$

# 实对称阵不同特征值的特征向量正交

□ 令实对称矩阵为A，其两个不同的特征值 $\lambda_1, \lambda_2$ 对应的特征向量分别是 $\mu_1, \mu_2$ ；

■  $\lambda_1, \lambda_2, \mu_1, \mu_2$ 都是实数或是实向量。

$$\begin{aligned} & \begin{cases} A\mu_1 = \lambda_1\mu_1 \\ A\mu_2 = \lambda_2\mu_2 \Rightarrow \mu_1^T A\mu_2 = \mu_1^T \lambda_2\mu_2 \end{cases} \\ & \Rightarrow (A^T \mu_1)^T \mu_2 = \lambda_2 \mu_1^T \mu_2 \Rightarrow (\underline{A\mu_1})^T \mu_2 = \lambda_2 \mu_1^T \mu_2 \\ & \Rightarrow (\underline{\lambda_1 \mu_1})^T \mu_2 = \lambda_2 \mu_1^T \mu_2 \\ & \Rightarrow \lambda_1 \mu_1^T \mu_2 = \lambda_2 \mu_1^T \mu_2 \\ & \xrightarrow{\lambda_1 \neq \lambda_2} \mu_1^T \mu_2 = 0 \end{aligned}$$



# 谱和谱聚类

- 方阵作为线性算子，它的所有特征值的全体统称方阵的谱。
  - 方阵的谱半径为最大的特征值
  - 矩阵A的谱半径： $(A^T A)$ 的最大特征值
- 谱聚类是一种基于图论的聚类方法，通过对样本数据的拉普拉斯矩阵的特征向量进行聚类，从而达到对样本数据聚类的目的。

# 谱分析的整体过程

- 给定一组数据  $x_1, x_2, \dots, x_n$ ，记任意两个点之间的相似度(“距离”的减函数)为  $s_{ij} = \langle x_i, x_j \rangle$ ，形成相似度图(similarity graph):  $G=(V, E)$ 。如果  $x_i$  和  $x_j$  之间的相似度  $s_{ij}$  大于一定的阈值，那么，两个点是连接的，权值记做  $s_{ij}$ 。
- 接下来，可以用相似度图来解决样本数据的聚类问题：找到图的一个划分，形成若干个组(Group)，使得不同组之间有较低的权值，组内有较高的权值。

# 若干概念

---

□ 无向图  $G=(V,E)$

□ 邻接矩阵  $W = (w_{ij})_{i,j=1,\dots,n}$

□ 顶点的度  $d_i \rightarrow$  度矩阵  $D$  (对角阵)

$$d_i = \sum_{j=1}^n w_{ij}$$

# 若干概念

---

## □ 子图A的指示向量

$$\mathbb{1}_A = (f_1, \dots, f_n)' \in \mathbb{R}^n$$

$$f_i = 1 \text{ if } v_i \in A$$

$$f_i = 0 \text{ otherwise}$$

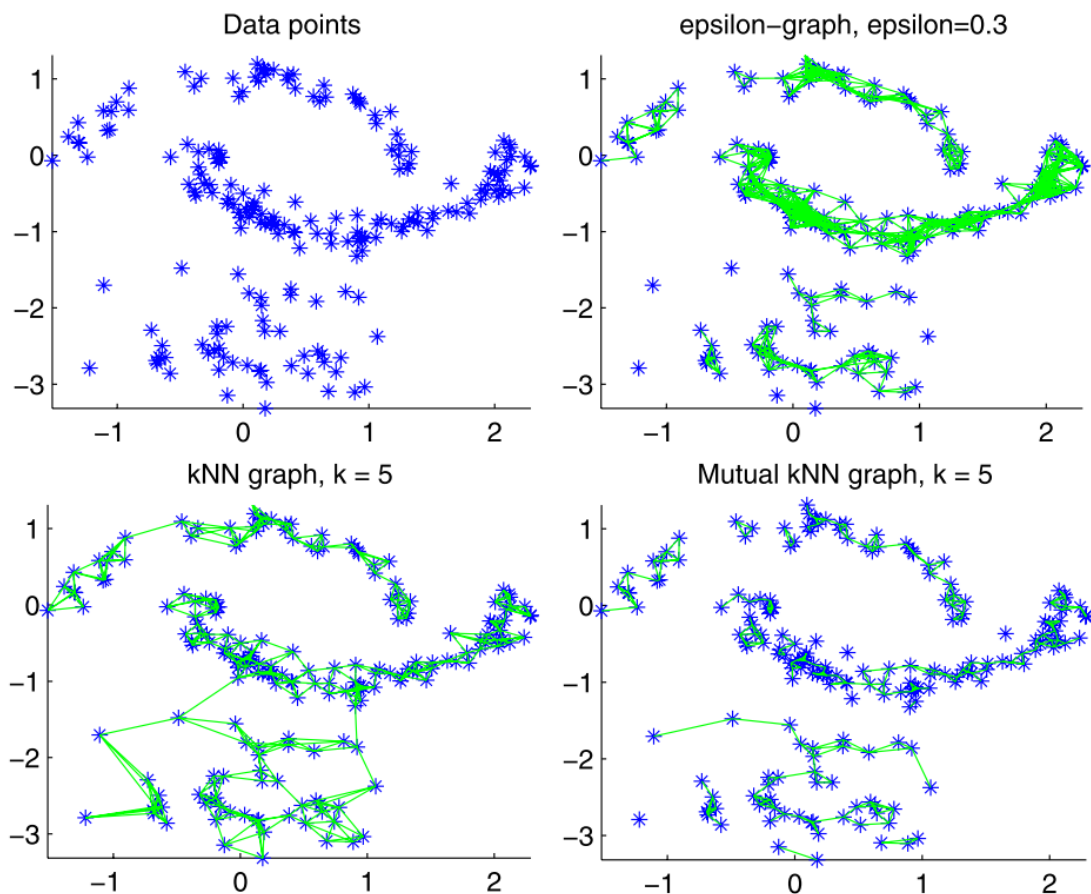
## □ A和B是图G的不相交子图，则定义子图的连接权：

$$W(A, B) := \sum_{i \in A, j \in B} w_{ij}$$

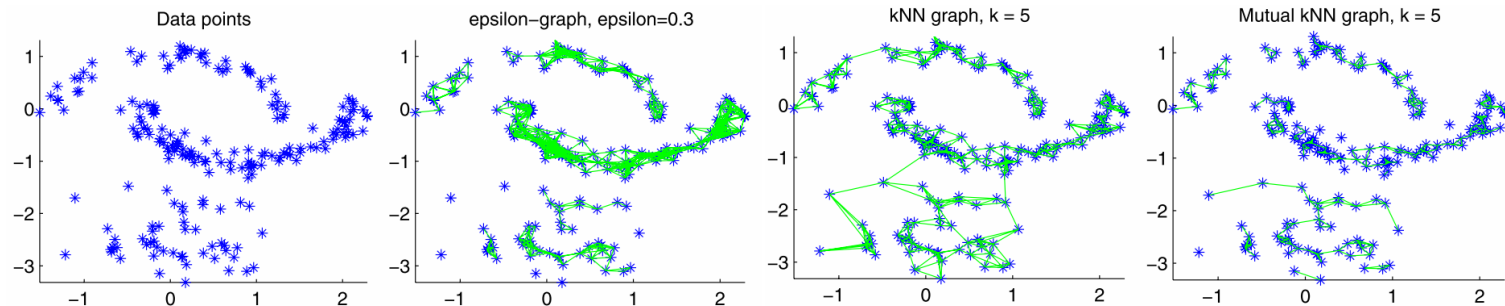
# 相似度图G的建立方法

- 全连接图：距离越大，相似度越小
  - 高斯相似度  $s(x_i, x_j) = \exp\left(\frac{-\|x_i - x_j\|^2}{2\sigma^2}\right)$
- $\epsilon$ 近邻图
  - 给定参数 $\epsilon$ /如何选择 $\epsilon$ ?
    - 图G的权值的均值
    - 图G的最小生成树的最大边
- k近邻图(k-nearest neighbor graph)
  - 若 $v_i$ 的k最近邻包含 $v_j$ ， $v_j$ 的k最近邻不一定包含 $v_i$ ：有向图
  - 忽略方向的图，往往简称“k近邻图”
  - 两者都满足才连接的图，称作“互k近邻图(mutual)”

# 相似度图G的举例



# 权值比较



- $\epsilon$ 近邻图:  $\epsilon=0.3$ , “月牙部分”非常紧的连接了, 但“高斯部分”很多都没连接。当数据有不同的“密度”时, 往往发生这种问题。
- $k$ 近邻图: 可以解决数据存在不同密度时有些无法连接的问题, 甚至低密度的“高斯部分”与高密度的“月牙部分”也能够连接。同时, 虽然两个“月牙部分”的距离比较近, 但 $k$ 近邻还可以把它们区分开。
- 互 $k$ 近邻图: 它趋向于连接相同密度的部分, 而不连接不同密度的部分。这种性质介于 $\epsilon$ 近邻图和 $k$ 近邻图之间。如果需要聚类不同的密度, 这个性质非常有用。
- 全连接图: 使用高斯相似度函数可以很好的建立权值矩阵。但缺点是建立的矩阵不是稀疏的。
- 总结: 首先尝试使用 $k$ 近邻图。

# 拉普拉斯矩阵及其性质

□ 拉普拉斯矩阵:  $L = D - W$

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

□  $L$  是对称半正定矩阵, 最小特征值是0, 相应的特征向量是全1向量。



# 拉普拉斯矩阵的定义

- 计算点之间的邻接相似度矩阵  $W$ 
  - 若两个点的相似度值越大，表示这两个点越相似；
  - 同时，定义  $w_{ij}=0$  表示  $v_i, v_j$  两个点没有任何相似性(无穷远)
- $W$  的第  $i$  行元素的和为  $v_i$  的度。形成顶点度对角阵  $D$ 
  - $d_{ii}$  表示第  $i$  个点的度
  - 除主对角线元素， $D$  其他位置为 0
- 未正则的拉普拉斯矩阵： $L = D - W$
- 正则拉普拉斯矩阵
  - 对称拉普拉斯矩阵  $L_{sym} = D^{-\frac{1}{2}} \cdot L \cdot D^{\frac{1}{2}} = I - D^{-\frac{1}{2}} \cdot W \cdot D^{\frac{1}{2}}$
  - 随机游走拉普拉斯矩阵  $L_{rw} = D^{-1}L = I - D^{-1}W$
  - Random walk

# 谱聚类算法：未正则拉普拉斯矩阵

- 输入：n个点 $\{p_i\}$ ，簇的数目k
  - 计算 $n \times n$ 的相似度矩阵 $W$ 和度矩阵 $D$ ；
  - 计算拉普拉斯矩阵 $L=D-W$ ；
  - 计算 $L$ 的前k个特征向量 $u_1, u_2, \dots, u_k$ ；
  - 将k个列向量 $u_1, u_2, \dots, u_k$ 组成矩阵 $U$ ， $U \in \mathbb{R}^{n \times k}$ ；
  - 对于 $i=1, 2, \dots, n$ ，令 $y_i \in \mathbb{R}^k$ 是 $U$ 的第i行的向量；
  - 使用k-means算法将点 $(y_i)_{i=1, 2, \dots, n}$ 聚类成簇 $C_1, C_2, \dots, C_k$ ；
  - 输出簇 $A_1, A_2, \dots, A_k$ ，其中， $A_i = \{j | y_j \in C_i\}$

# 谱聚类算法：随机游走拉普拉斯矩阵

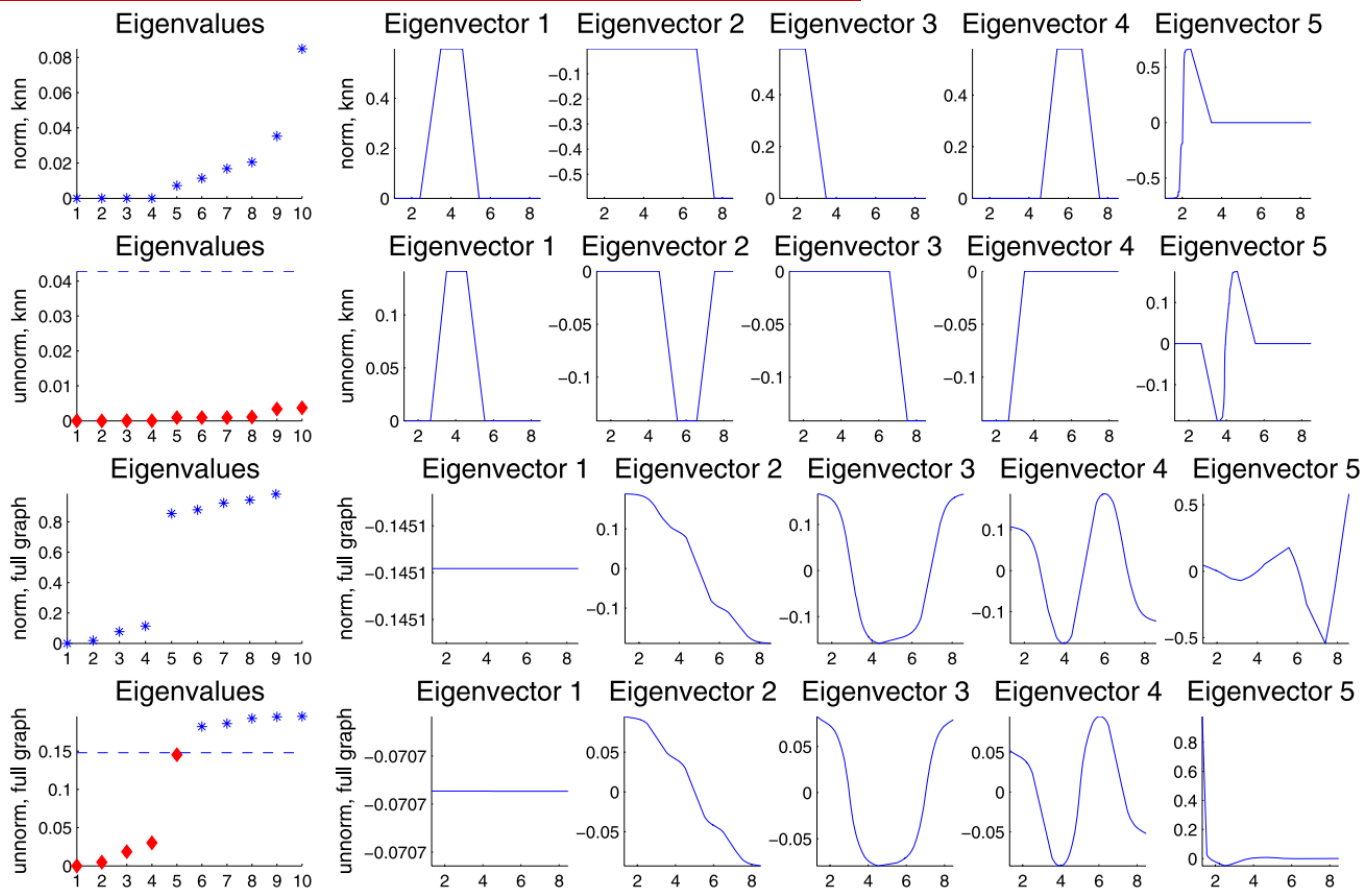
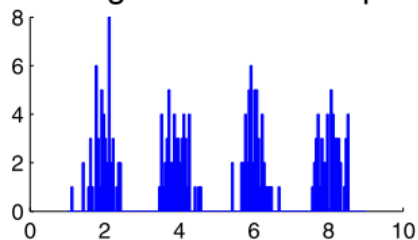
- 输入：n个点 $\{p_i\}$ ，簇的数目k
  - 计算 $n \times n$ 的相似度矩阵 $W$ 和度矩阵 $D$ ；
  - 计算正则拉普拉斯矩阵 $L_{rw} = D^{-1}(D - W)$ ；
  - 计算 $L_{rw}$ 的前k个特征向量 $u_1, u_2, \dots, u_k$ ；
  - 将k个列向量 $u_1, u_2, \dots, u_k$ 组成矩阵 $U$ ， $U \in R^{n \times k}$ ；
  - 对于 $i=1, 2, \dots, n$ ，令 $y_i \in R^k$ 是 $U$ 的第i行的向量；
  - 使用k-means算法将点 $(y_i)_{i=1, 2, \dots, n}$ 聚类成簇 $C_1, C_2, \dots, C_k$ ；
  - 输出簇 $A_1, A_2, \dots, A_k$ ，其中， $A_i = \{j | y_j \in C_i\}$

# 谱聚类算法：对称拉普拉斯矩阵

- 输入：n个点 $\{p_i\}$ ，簇的数目k
  - 计算 $n \times n$ 的相似度矩阵W和度矩阵D；
  - 计算正则拉普拉斯矩阵 $L_{\text{sym}} = D^{-1/2}(D-W)D^{-1/2}$ ；
  - 计算 $L_{\text{sym}}$ 的前k个特征向量 $u_1, u_2, \dots, u_k$ ；
  - 将k个列向量 $u_1, u_2, \dots, u_k$ 组成矩阵U， $U \in \mathbb{R}^{n \times k}$ ；
  - 对于 $i=1, 2, \dots, n$ ，令 $y_i \in \mathbb{R}^k$ 是U的第i行的向量；
  - 对于 $i=1, 2, \dots, n$ ，将 $y_i \in \mathbb{R}^k$ 依次单位化，使得 $|y_i|=1$ ；
  - 使用k-means算法将点 $(y_i)_{i=1, 2, \dots, n}$ 聚类成簇 $C_1, C_2, \dots, C_k$ ；
  - 输出簇 $A_1, A_2, \dots, A_k$ ，其中， $A_i = \{j | y_j \in C_i\}$

# 一个实例

Histogram of the sample



# Code

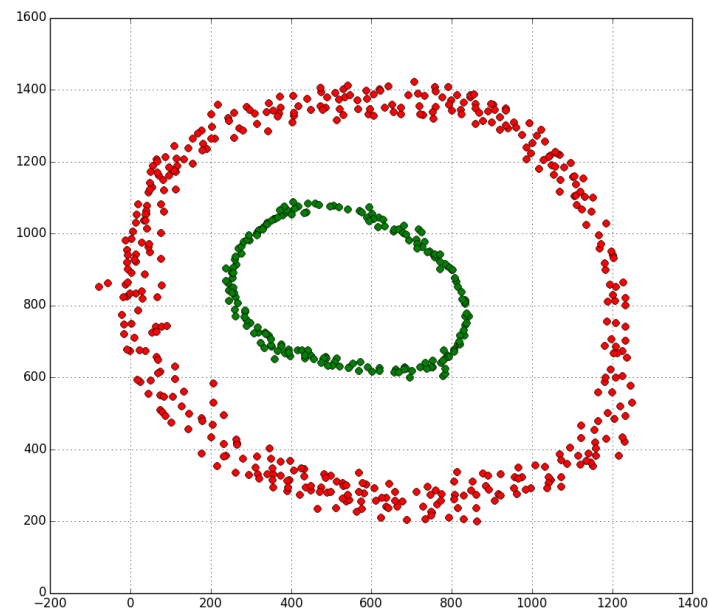
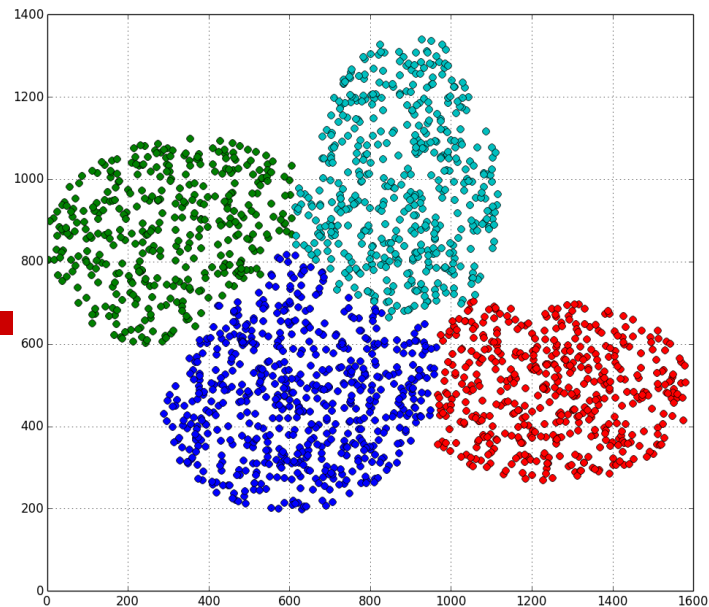
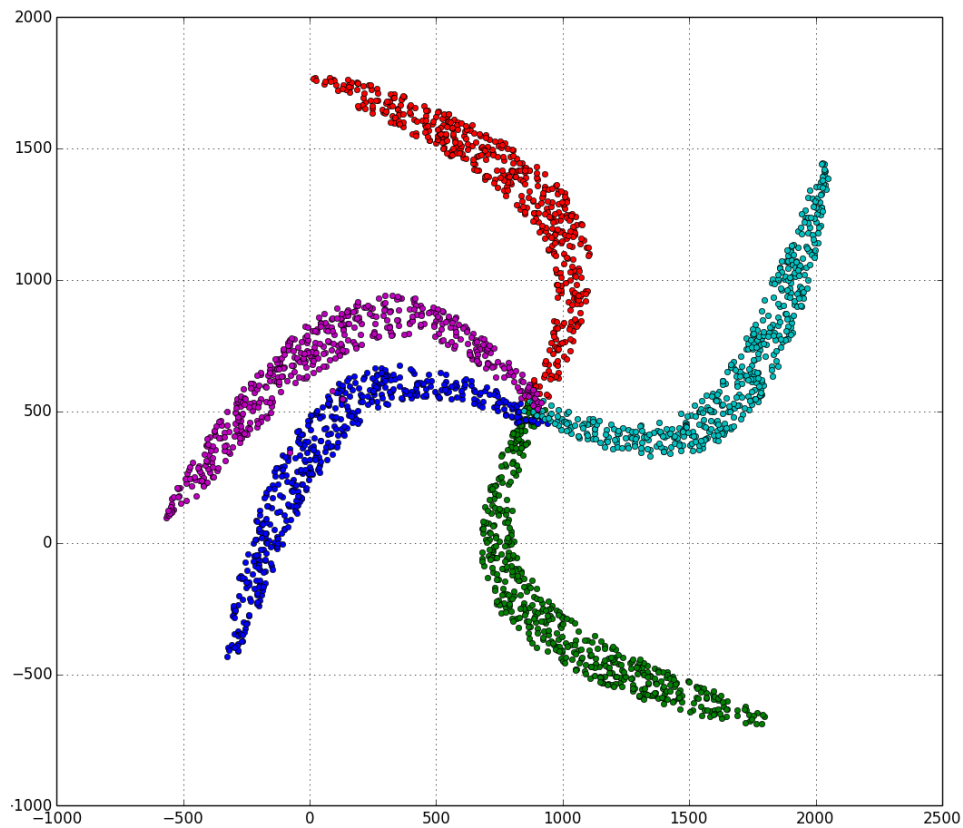
```
def spectral_cluster(data):
    lm = laplace_matrix(data)
    eg_values, eg_vectors = linalg.eig(lm)
    idx = eg_values.argsort()
    eg_vectors = eg_vectors[:, idx]

    m = len(data)
    eg_data = [[] for x in range(m)]
    for i in range(m):
        eg_data[i] = [0 for x in range(k)]
        for j in range(k):
            eg_data[i][j] = eg_vectors[i][j]
    return k_means(eg_data)
```

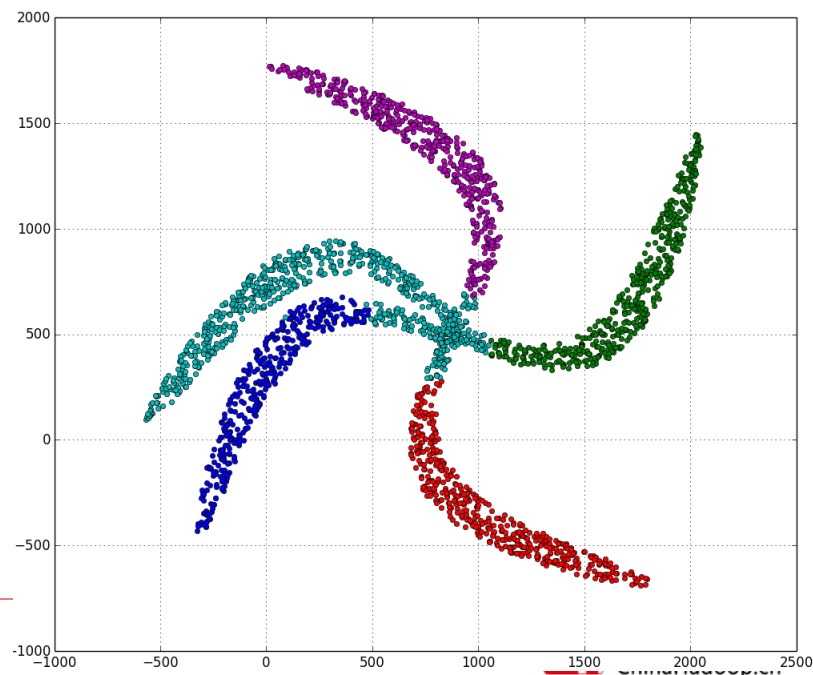
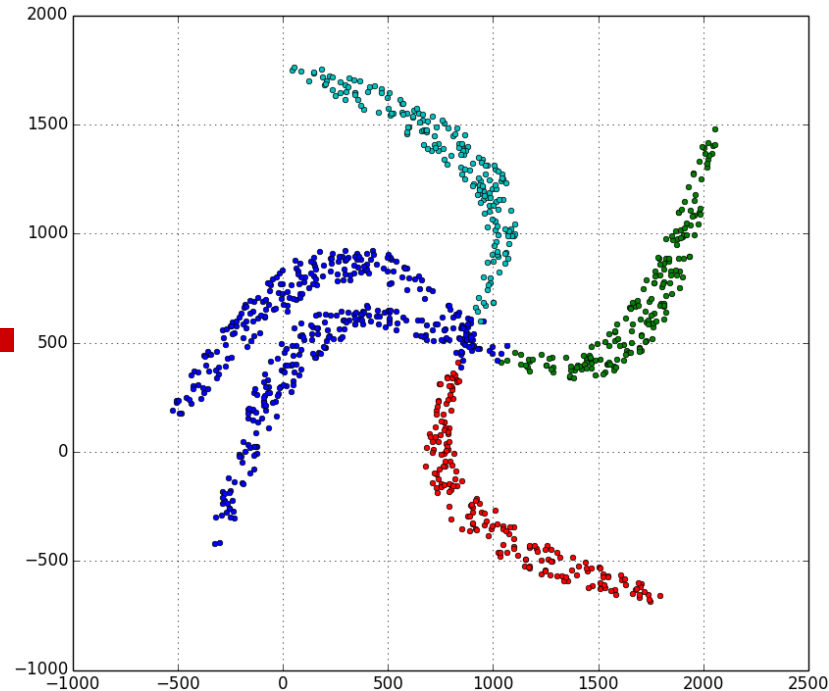
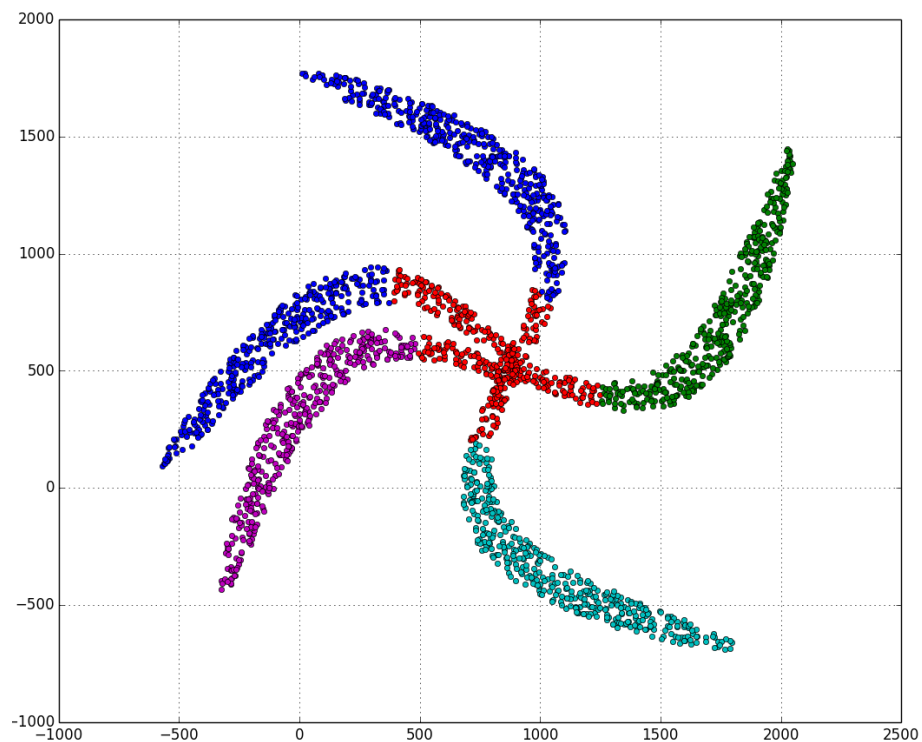
```
def laplace_matrix(data):
    m = len(data)
    w = [[] for x in range(m)]
    for i in range(m):
        w[i] = [0 for x in range(m)]
    nearest = [0 for x in range(neighbor)]

    for i in range(m):
        zero_list(nearest)
        for j in range(i+1, m):
            w[i][j] = similar(data, i, j)
            if not is_neighbor(w[i][j], nearest):
                w[i][j] = 0
            w[j][i] = w[i][j]    #对称
        w[i][i] = 0
    for i in range(m):
        s = 0
        for j in range(m):
            s += w[i][j]
        if s == 0:
            print "矩阵第", i, "行全为0"
            continue
        for j in range(m):
            w[i][j] /= s
            w[i][j] = -w[i][j]
        w[i][i] += 1    #单位阵主对角线为1
    return w
```

# 聚类效果



# 聚类失败的情况





# 进一步思考

- 谱聚类中的K如何确定?  $k^* = \arg \max_k |\lambda_{k+1} - \lambda_k|$
- 最后一步K-Means的作用是什么?
  - 目标函数是关于子图划分指示向量的函数, 该向量的值根据子图划分确定, 是离散的。该问题是NP的, 转换成求连续实数域上的解, 最后用K-Means算法离散化。
- 未正则/对称/随机游走拉普拉斯矩阵, 首选哪个?
  - 随机游走拉普拉斯矩阵
- 谱聚类可以用切割图/随机游走/扰动论等解释。

# 随机游走和拉普拉斯矩阵的关系

- 图论中的随机游走是一个随机过程，它从一个顶点跳转到另外一个顶点。谱聚类即找到图的一个划分，使得随机游走在相同的簇中停留而几乎不会游走到其他簇。
- 转移矩阵：从顶点 $v_i$ 跳转到顶点 $v_j$ 的概率正比于边的权值 $w_{ij}$

$$p_{ij} = w_{ij} / d_i \quad P = D^{-1}W$$

# 标签传递算法

---

- 对于部分样本的标记给定，而大多数样本的标记未知的情形，是半监督学习问题。
- 标签传递算法(Label Propagation Algorithm,LPA)，将标记样本的标记通过一定的概率传递给未标记样本，直到最终收敛。

# Code

```
def label_propagation(data, a):
    p = transition_matrix(data)
    m = len(data)
    n = len(data[0])
    for times in range(100):
        for i in range(a, m):
            j = calc_label(p, i)
            label = data[j][n-1]
            if label > 0:
                data[i][n-1] = label

def calc_label(p, i):
    n = len(p[i])
    k = random.random() #  $k \in [0, 1)$ 
    r = n-1
    for j in range(n):
        if p[i][j] > k:
            r = j
            break
    return r
```

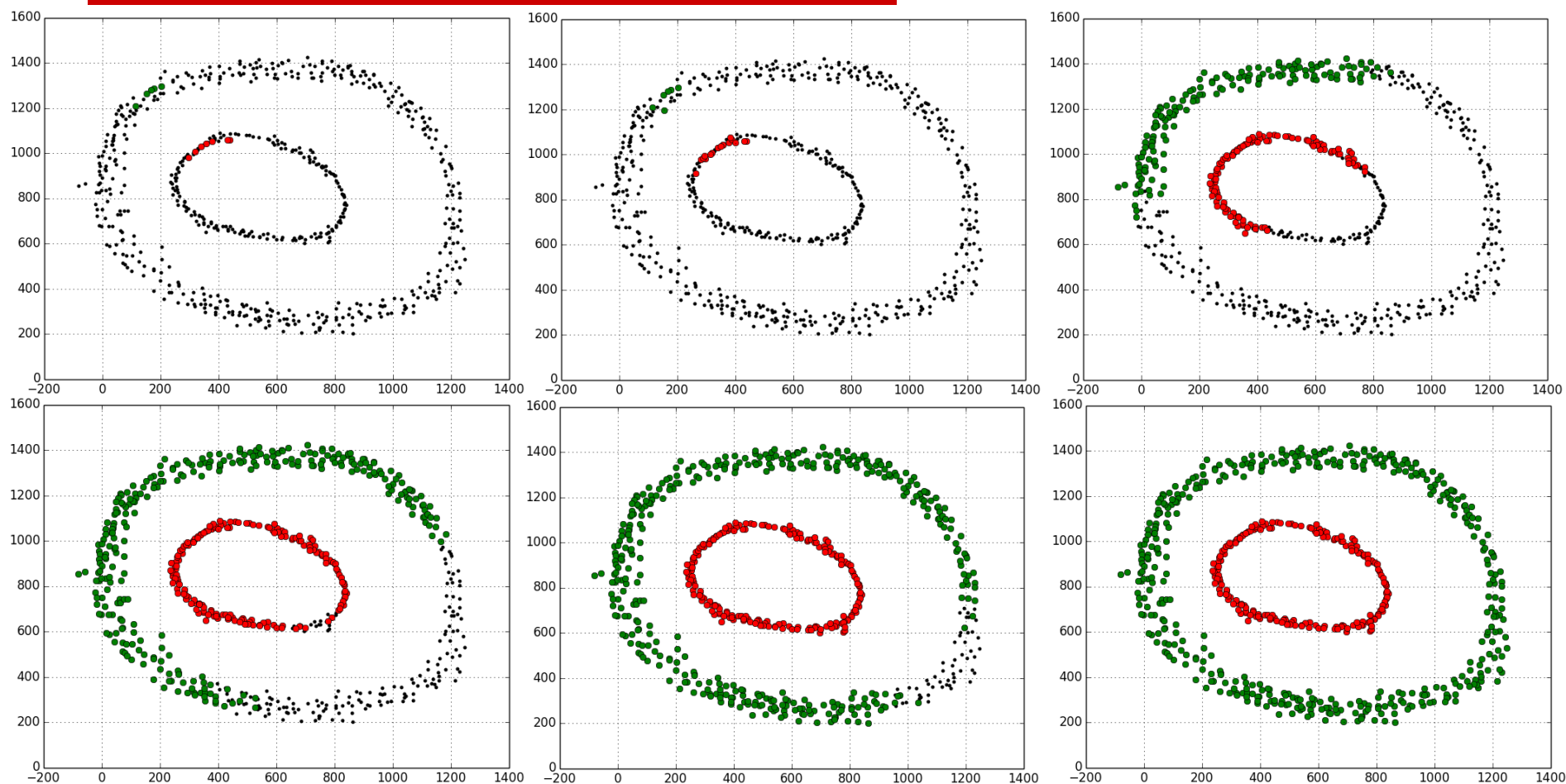
```
def transition_matrix(data):
    m = len(data)
    p = [[] for x in range(m)]
    for i in range(m):
        p[i] = [0 for x in range(m)]
    nearest = [0 for x in range(neighbor)]

    for i in range(m):
        zero_list(nearest)
        for j in range(i+1, m):
            p[i][j] = similar(data, i, j)
            if not is_neighbor(p[i][j], nearest):
                p[i][j] = 0
            p[j][i] = p[i][j] # 对称
        p[i][i] = 0
    for i in range(m):
        s = 0
        for j in range(m):
            s += p[i][j]
        if s == 0:
            print "矩阵第", i, "行全为0"
            continue
        for j in range(m):
            p[i][j] /= s
            if j != 0:
                p[i][j] += p[i][j-1]

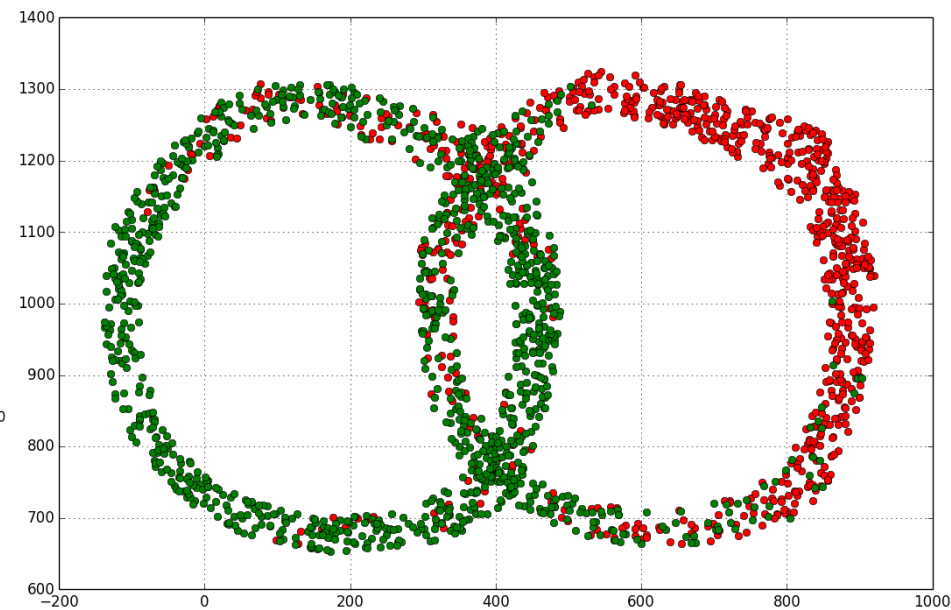
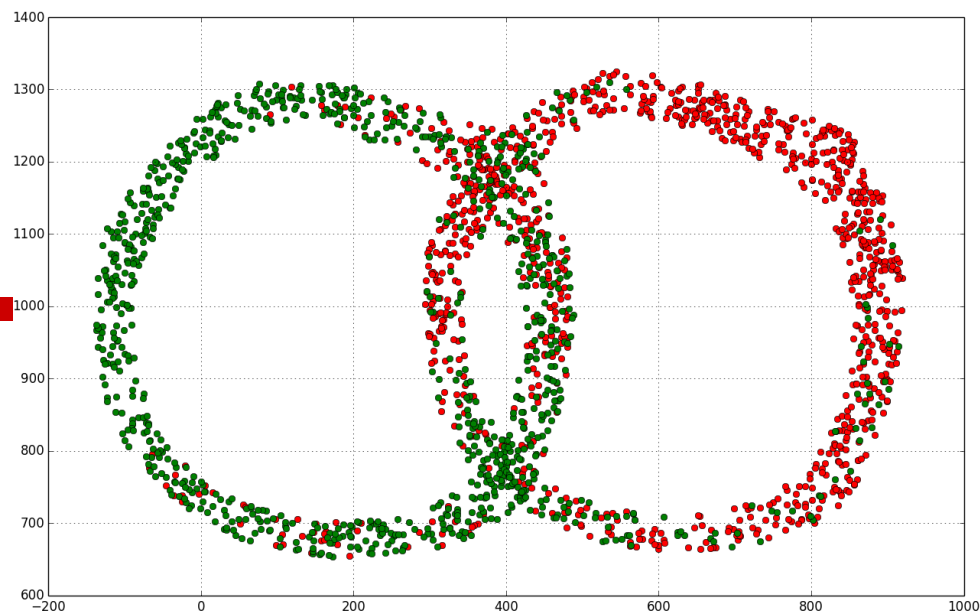
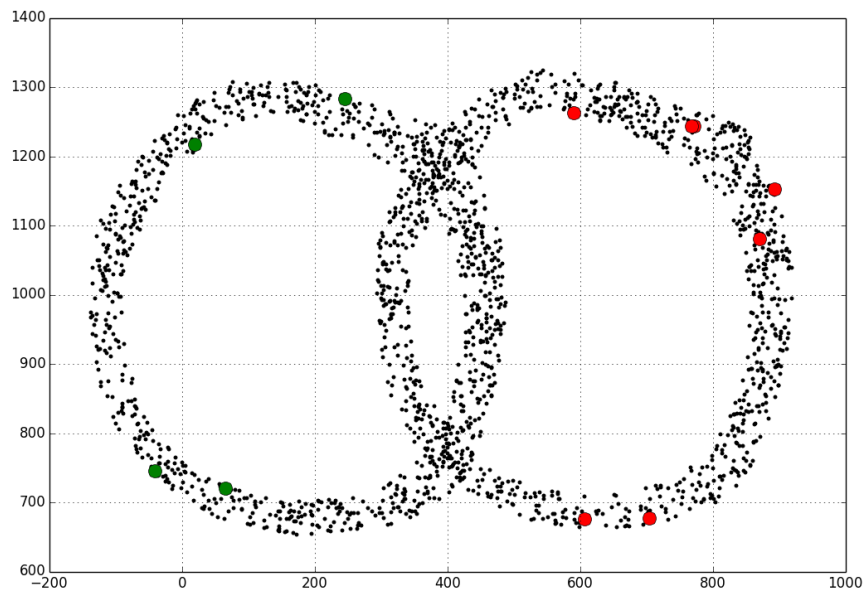
    return p
```

# 标签传递过程

初始	1	10
20	30	40



# 带宽/邻域影响



# 参考文献

---

- ❑ Alex Rodriguez, Alessandro Laio. *Clustering by fast search and find of density peak*. Science. 2014
- ❑ Ulrike von Luxburg. *A tutorial on spectral clustering*. 2007
- ❑ Lang K. *Fixing two weaknesses of the spectral method*. Advances in Neural Information Processing Systems 18, 715–722. MIT Press, Cambridge, 2006
- ❑ Bach F, Jordan M. *Learning spectral clustering*. Advances in Neural Information Processing Systems 16 (NIPS). 305–312. MIT Press, Cambridge, 2004
- ❑ R.J.G.B. Campello, D. Moulavi, A. Zimek and J. Sander  
Hierarchical Density Estimates for Data Clustering, Visualization, and Outlier Detection, ACM Trans. on Knowledge Discovery from Data, Vol 10, 1 (July 2015), 1-51.

# 我们在这里

□ <http://wenda.ChinaHadoop.cn>

■ 视频/课程/社区

□ 微博

■ @ChinaHadoop

■ @邹博\_机器学习

□ 微信公众号

■ 小象学院

■ 大数据分析挖掘

互联网新技术在线教育领航者

小象问答 搜索标题、用户 全站内容搜索 提问 首页 动态 发现 话题 通知

全部 招聘求职 机器学习 大数据平台技术 DCon 大数据行业应用 NoSQL数据库 数据科学 江湖救急

发现 最新 推荐 热门 等待回复

graphviz has no attribute 'write' 贡献  
邹博 回复了问题 • 2 人关注 • 1 个回复 • 3 次浏览 • 2017-04-09 15:48

sklearn中如何理解Pipeline机制 贡献  
数据分析与数据挖掘 邹博 回复了问题 • 2 人关注 • 1 个回复 • 28 次浏览 • 2017-04-09 15:39

关于9.Logistic回归的ppt中第9页的对数线性函数 贡献  
机器学习 邹博 回复了问题 • 3 人关注 • 3 个回复 • 39 次浏览 • 2017-04-09 15:35

关于“贝叶斯估计中，最大后验概率估计就是结构化风险最小化的例子：当模型是条件概率分布，损失函数为对数损失函数，模型的复杂度由模型的先验概率表示，结构化风险最小化就等价于最大后验概率估计” 贡献  
机器学习 邹博 回复了问题 • 2 人关注 • 1 个回复 • 26 次浏览 • 2017-04-09 15:27

关于连续值的预测 贡献  
咨询 邹博 回复了问题 • 2 人关注 • 1 个回复 • 31 次浏览 • 2017-04-09 15:24

拉格朗日对偶函数为什么一定是凸函数 贡献  
数据科学 邹博 回复了问题 • 2 人关注 • 2 个回复 • 26 次浏览 • 2017-04-09 15:20

梯度下降公式中的斯堪J 是 贡献  
机器学习 邹博 回复了问题 • 2 人关注 • 1 个回复 • 29 次浏览 • 2017-04-09 15:17

深度学习适合做预测吗？ 贡献  
深度学习 邹博 回复了问题 • 2 人关注 • 1 个回复 • 27 次浏览 • 2017-04-09 15:15

关于6.4PCA\_FeatureSelection.py中plt.legend的参数疑问 贡献  
机器学习 邹博 回复了问题 • 2 人关注 • 1 个回复 • 28 次浏览 • 2017-04-09 15:04

@邹博 有哪些可以下载数据源的网站？ 贡献  
数据分析与数据挖掘 邹博 回复了问题 • 4 人关注 • 1 个回复 • 31 次浏览 • 2017-04-09 14:53

LDA主题模型 贡献  
机器学习 邹博 回复了问题 • 2 人关注 • 1 个回复 • 29 次浏览 • 2017-04-09 14:45

代码10.6bagging\_ridged老师提到了采样率设为0.2能够使峰值部分的数据被体现出来。这是为什么呢？ 贡献  
机器学习 邹博 回复了问题 • 2 人关注 • 1 个回复 • 22 次浏览 • 2017-04-09 14:26

GraphViz's executables not found 贡献  
机器学习 邹博 回复了问题 • 3 人关注 • 2 个回复 • 23 次浏览 • 2017-04-09 13:47

决策树中关于feature\_importances代码的问题 贡献  
机器学习 邹博 回复了问题 • 2 人关注 • 1 个回复 • 6 次浏览 • 2017-04-09 13:11

专题  
招聘求职  
大数据行业应用  
数据科学  
系统与编程  
云计算技术

热门话题 更多 >  
机器学习 907 个问题, 230 人关注  
spark 387 个问题, 172 人关注  
hadoop 1059 个问题, 155 人关注  
python数据分析 171 个问题, 28 人关注  
数据分析与数据挖掘 54 个问题, 111 人关注

热门用户 更多 >  
小心巴 14 个问题, 0 次赞同  
又又V 45 个问题, 22 次赞同  
铁甲无声 10 个问题, 0 次赞同  
带刀锦衣卫 13 个问题, 0 次赞同



---

感谢大家！

恳请大家批评指正！