

自然语言处理  
文本聚类

软件学院

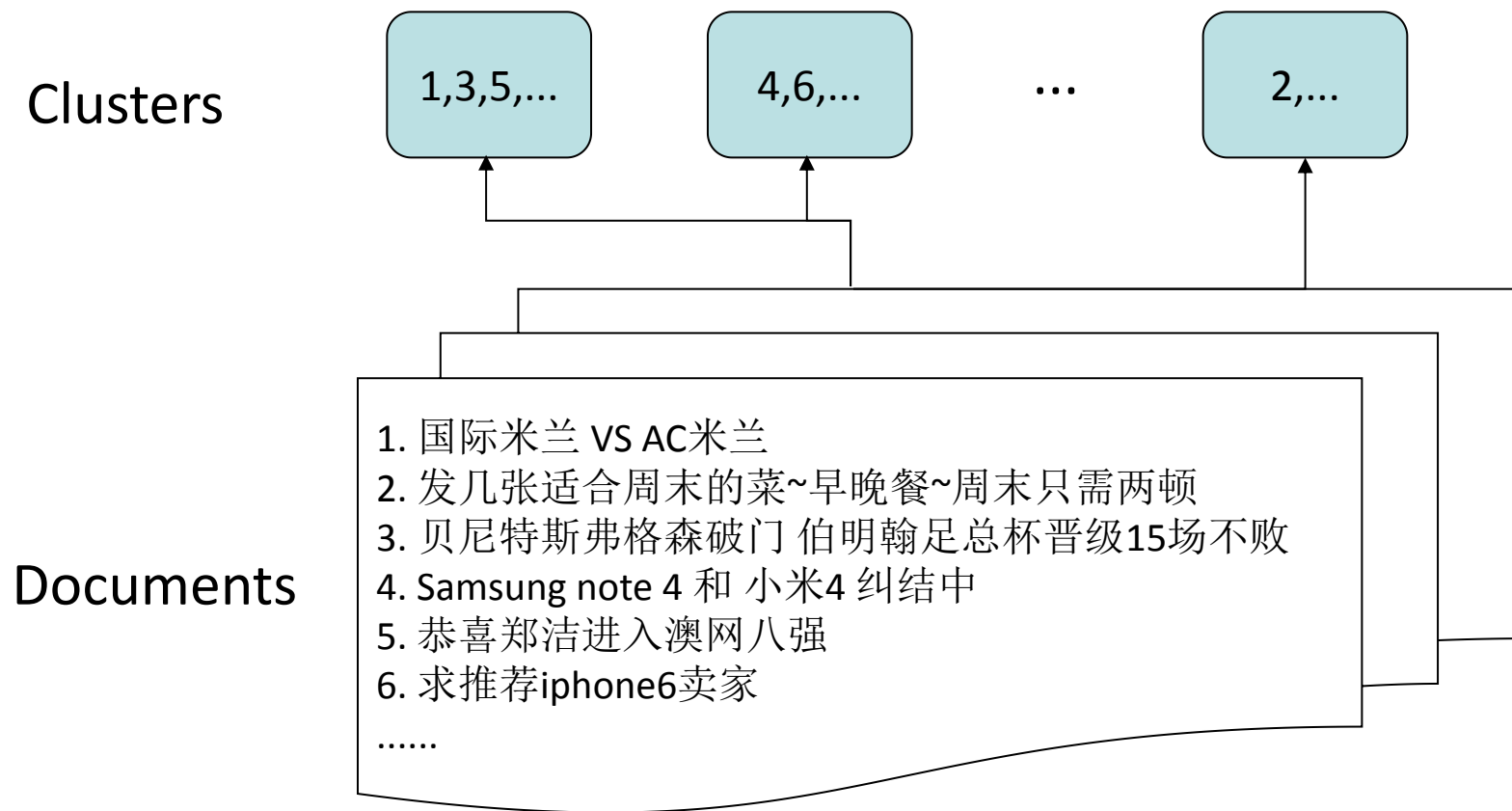
# Overview

- 目前为止：分类问题
  - 应用：文本主题分类、语言识别、词义消歧等
  - “有监督”学习框架：学习时需要类别标签
    - 生成模型： e.g., Naïve Bayes
      - 模型=一组分布  $P(x, y; \theta)$ ,  $\theta \in \Theta$
      - 选择一种最可能从中抽样得到训练数据的分布  $P(x, y; \hat{\theta})$
      - 根据  $\hat{y} = \arg \max_y P(x, y; \hat{\theta})$  对新样本  $x$  进行分类
    - 判别模型： e.g., maximum entropy models (a.k.a. logistic regression)
      - 模型=一组分类函数  $F$
      - 选择一种可以对训练样本正确分类的函数  $\hat{f} \in F$
      - 根据  $\hat{y} = \text{sign}(\hat{f}(x))$  对新样本  $x$  进行分类

# Overview

- 接下来：聚类问题
  - “无监督”学习框架：学习时不需要类别标签
  - **Magic:** 挖掘数据内部的隐藏模式（结构、类别等）
  - 很多NLP任务都关乎聚类问题：IR, recommendation system, exploratory data analysis

# 例子：标题聚类



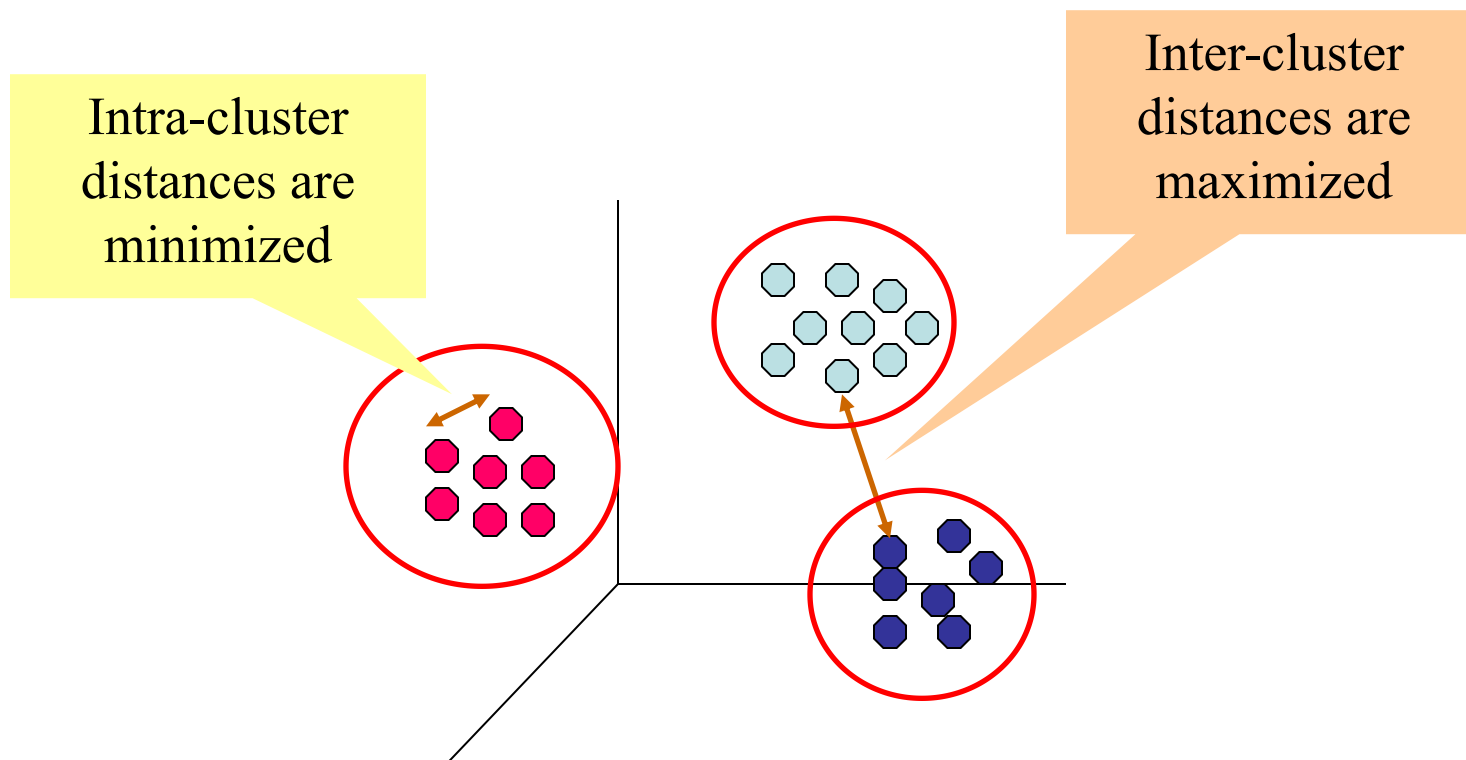
# 例子：web搜索优化

- 用户提交的query本身往往包含歧义
  - E.g., Jaguar, NLP, Paris Hilton



- 前10页几乎没有关于“animal”词义下的结果
- 仅靠词义消歧通常不能很好的解决web搜索的歧义
  - 需要交由用户选择相关的搜索结果
  - 将搜索结果按照query的语义进行分组
  - Jaguar or Jaguar car

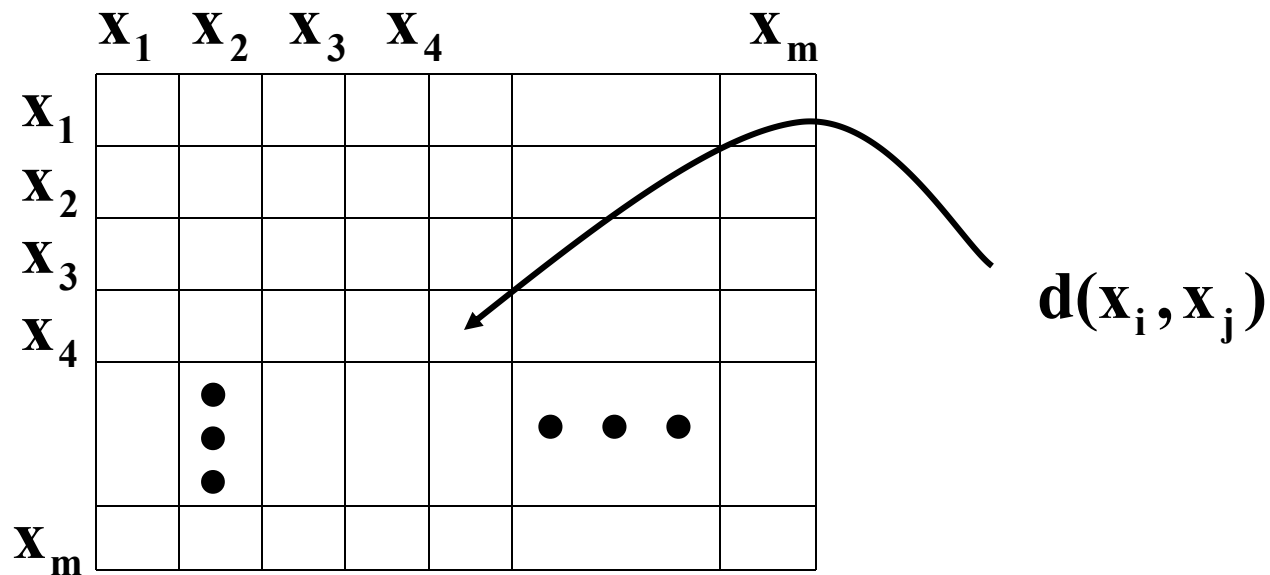
# 聚类



- 如何定义样本间的距离？
- 如何找到这样一组划分？

# 距离测度

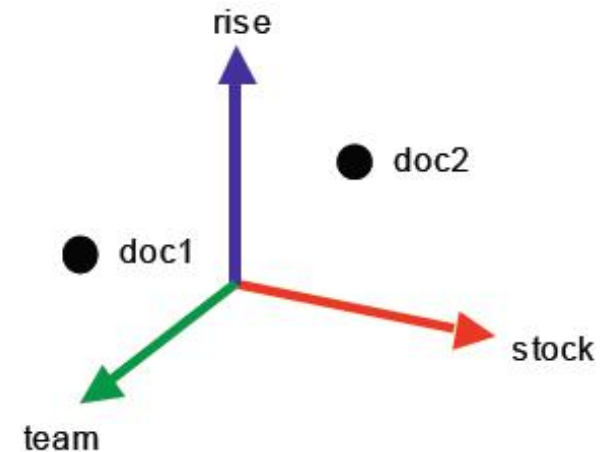
- 假设每个样本可以看成高维空间内的一个点，则可以采用距离测度表示样本间的距离
- 则聚类问题的输入可以看成由样本距离张成的一个矩阵：



# 距离测度

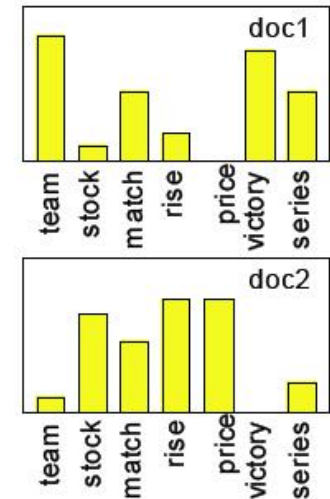
- 向量空间表示

- 一个文档可以看成高维向量空间内的一个点
- Distance: Euclidean / cosine



- 概率表示

- 一个文档可以看做从某个经验分布抽样得到的词汇的集合
- Distance: divergences





# 距离测度

- 传统意义上，距离测度可以表示为一个函数，满足：

1.  $d(x,y) \geq 0, d(x,y) = 0 \Leftrightarrow x = y$

2.  $d(x,y) + d(y,z) \geq d(x,z)$

3.  $d(x,y) = d(y,x)$

- 但对于聚类问题，采用的距离测度（或相似度）可以不严格满足以上三个约束

# 距离测度

## Examples:

- Euclidean Distance:

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})} = \sqrt{\sum_{i=1}^d (\mathbf{x}_i - \mathbf{y}_i)^2} \quad \mathbf{L}_2$$

- Manhattan Distance:

$$d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_1 = \sum_{i=1}^d |\mathbf{x}_i - \mathbf{y}_i| \quad \mathbf{L}_1$$

- Chebyshev Distance:

$$d(\mathbf{x}, \mathbf{y}) = \max_{1 \leq i \leq d} |\mathbf{x}_i - \mathbf{y}_i| \quad \mathbf{L}_\infty$$

- 注意: 对于欧氏距离  $d(\mathbf{x}, \mathbf{y})$  ,  $d^2(\mathbf{x}, \mathbf{y})$  不是一个测度 (metric) 但可用于距离度量 (measure) (不满足三角不等式)

# 距离测度

Examples:

- Mahalanobis Distance:

$$\mathbf{d}(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T \Sigma^{-1} (\mathbf{x} - \mathbf{y})}$$

其中  $\Sigma$  表示协方差矩阵

# 距离测度

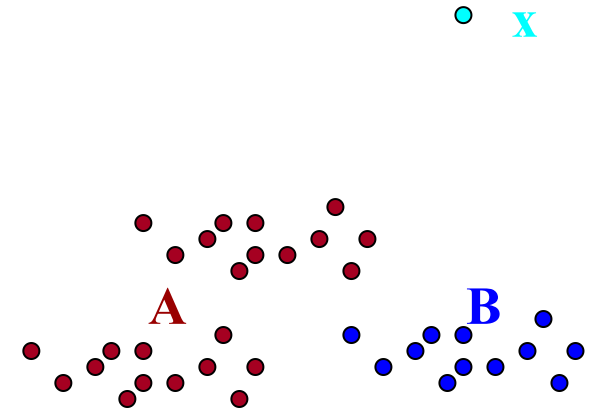
- 样本点  $x$  和样本点集合  $A$  之间的距离:

$$d(x, A) = \frac{1}{|A|} \sum_{y \in A} d(x, y)$$

- 两个样本点集合  $A, B$  之间的距离:

$$d(A, B) = \frac{1}{|A| |B|} \sum_{x \in A, y \in B} d(x, y)$$

- 注意, 还有其它多种不同的刻画方式



# 聚类模型

- K-means
- Hierarchical Clustering
- Density-based Clustering
- Gaussian Mixture Model
- Spectral Clustering
- .....

# K-means

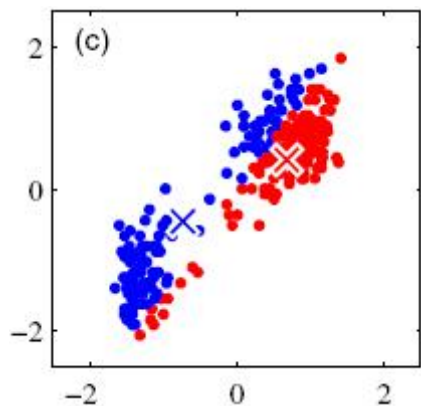
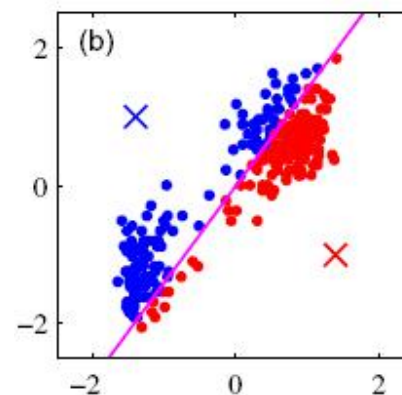
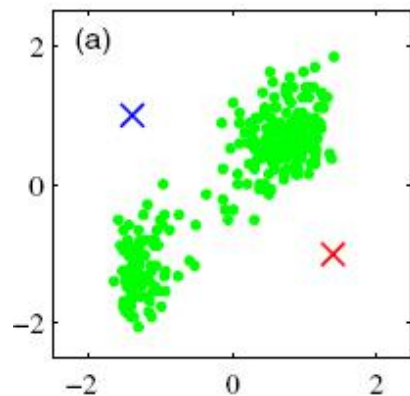
- 先确定簇的个数，**K**
- 假设每个簇都有一个中心点 (**centroid**)
- 将每个样本点划分到距离它最近的中心点所属的簇中
- **迭代过程:**

## **The basic algorithm:**

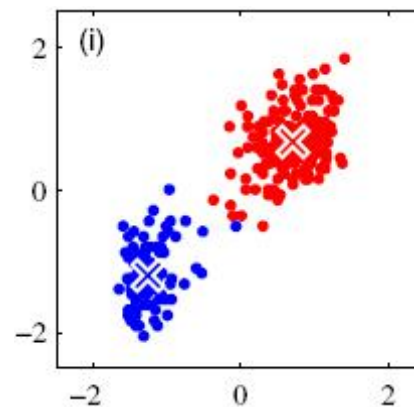
---

- 1: select K points as the initial centroids
  - 2: **repeat**
  - 3: Form K clusters by assigning all points to the *closest centroid*
  - 4: Recompute the centroid of each cluster
  - 5: **until** the centroids don't change
-

# K-means Clustering



...



# K-means

- 目标函数：定义为每个样本与其簇中心点的距离的平方和（the Sum of Squared Error, SSE）

$$SSE = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \text{dist}(x_n - \mu_k) \quad \text{e.g., } \text{dist}(x_n - \mu_k) = ||x_n - \mu_k||^2$$

- $\mu_k$  表示簇  $C_k$  的中心点（或其它能代表  $C_k$  的点）
  - 若  $x_n$  被划分到簇  $C_k$  则  $r_{nk}=1$ ，否则  $r_{nk}=0$
- 目标：找到簇的中心点  $\mu_k$  及簇的划分  $r_{nk}$  使得目标函数 SSE 最小



# K-means

- 1: choose some initial values for the  $\mu_k$  ( $k=1,\dots,K$ )
- 2: **repeat**
- 3: minimize SSE with respect to the  $r_{nk}$  ( $n=1,\dots,N$ )

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|x_n - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}$$

- 4: minimize SSE with respect to the  $\mu_k$

$$\mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$

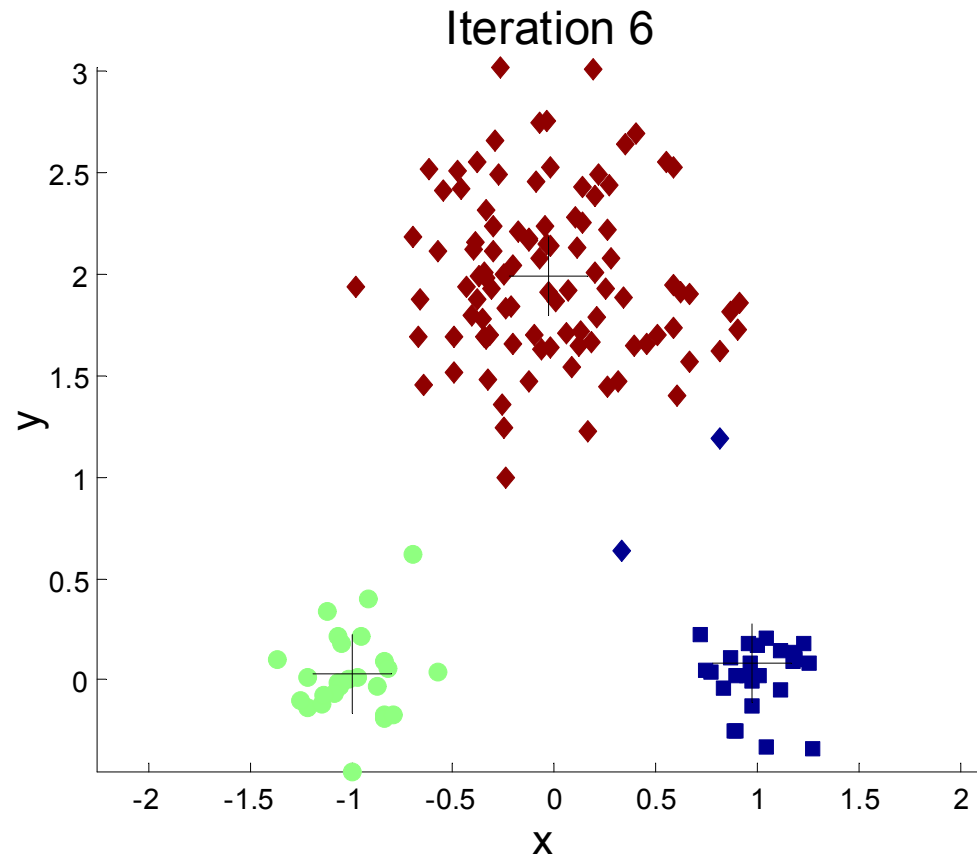
- 5: **until** convergence

# K-means – Details

- 初始中心点通常是随机选取的
  - 产生的簇可能和上一次迭代相差很大
- 中心点通常是当前簇中所有样本点的均值
  - K-medoids: 中位数
- 算法复杂度:  $O(n \times d \times K \times I)$ 
  - $n$ =样本点个数,  $d$ =样本特征维度
  - $K$ =类的个数,  $I$ =迭代次数
- 需要预先确定K!

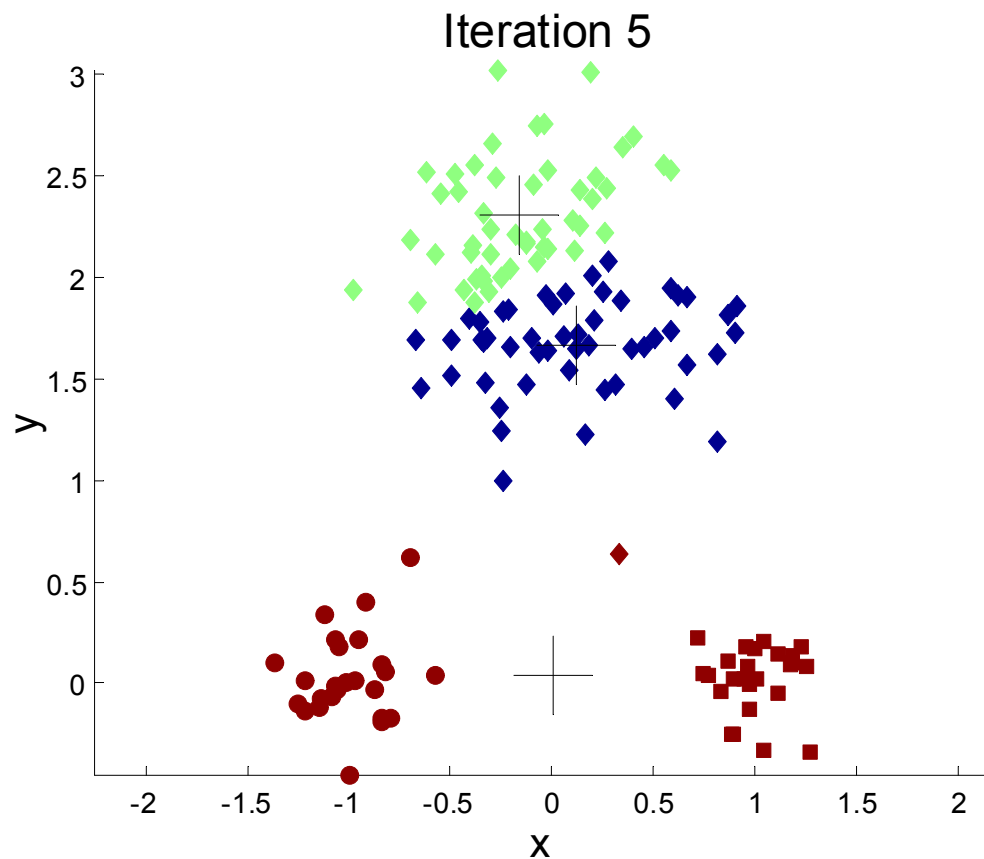
# K-means – Details

不同初始中心点的选取对聚类结果的影响:



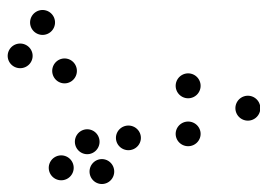
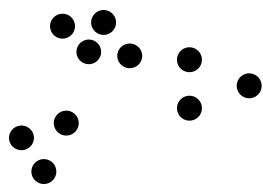
# K-means – Details

不同初始中心点的选取对聚类结果的影响:

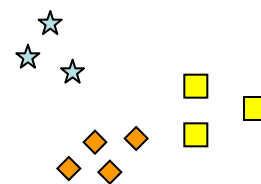
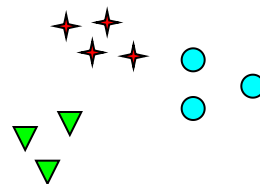


# K-means – Details

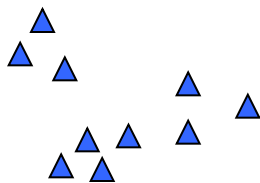
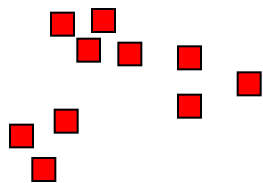
类别本身可能存在歧义：



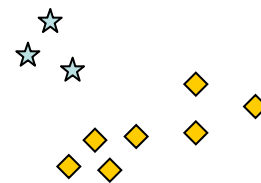
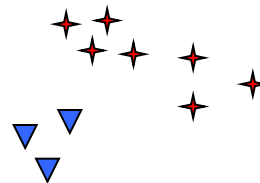
有多少个类？



Six Clusters



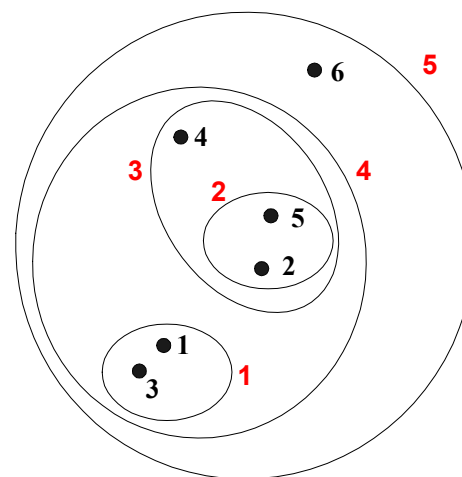
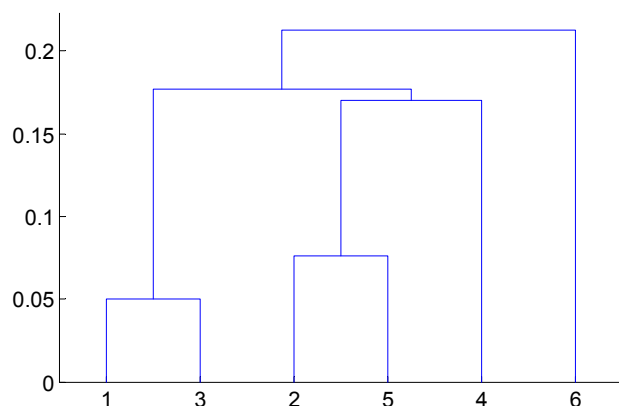
Two Clusters



Four Clusters

# Hierarchical Clustering

- 为数据集输出一个嵌套的、层次化的类别树：  
dendrogram
  - 树结构记录了簇的合并或拆分
  - 自底向上（agglomerative）
  - 自顶向下（divisive）



# Hierarchical Agglomerative Clustering (HAC)

- 层次凝聚聚类：一种很常用的聚类模型

## Basic algorithm:

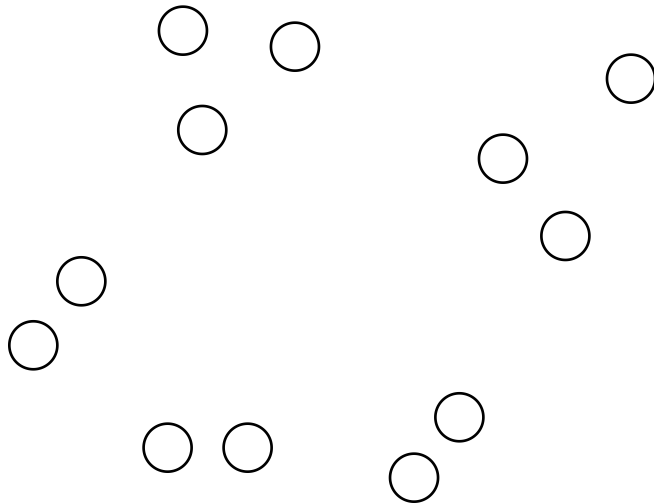
---

1. Compute the proximity matrix
  2. Let each data point be a cluster
  3. **Repeat**
  4.   **Merge** the two closest clusters
  5.   Update the proximity matrix
  6. **Until** only a single cluster remains
- 

- 关键步骤：计算两个簇的相似度
  - 不同的度量两个簇的相似度方法，区分了不同的聚类算法

# 初始状态

- 初始状态：每个点代表一个簇以及一个相似度矩阵(**proximity matrix**)



|    | p1 | p2 | p3 | p4 | p5 | ... |
|----|----|----|----|----|----|-----|
| p1 |    |    |    |    |    |     |
| p2 |    |    |    |    |    |     |
| p3 |    |    |    |    |    |     |
| p4 |    |    |    |    |    |     |
| p5 |    |    |    |    |    |     |
| .  |    |    |    |    |    |     |
| .  |    |    |    |    |    |     |
| .  |    |    |    |    |    |     |

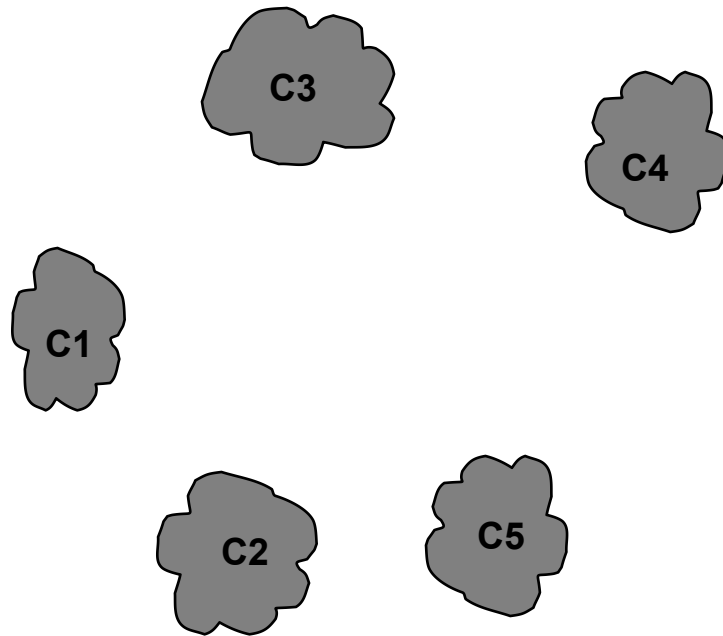
**Proximity Matrix**





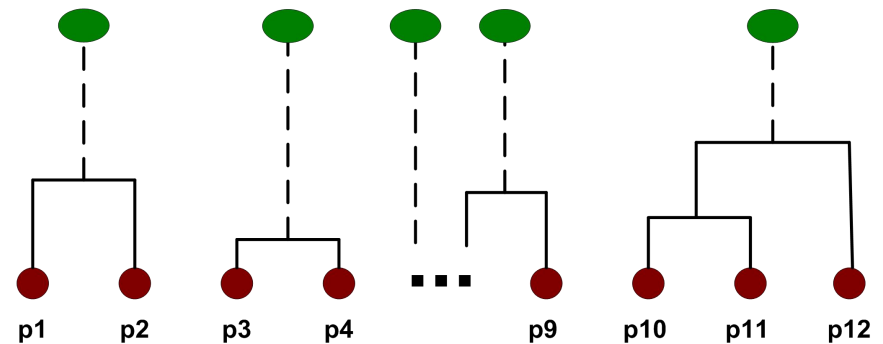
# 中间状态

- 经过几次合并操作，得到一些簇以及一个相似度矩阵



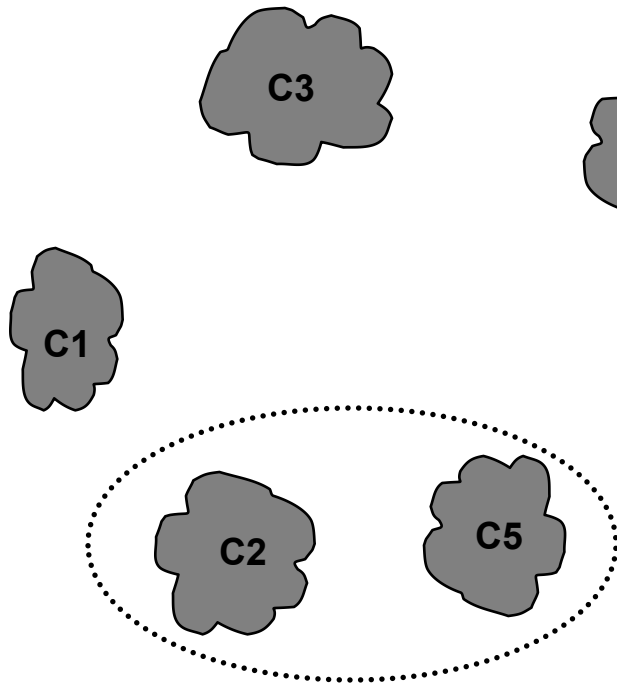
|    | C1 | C2 | C3 | C4 | C5 |
|----|----|----|----|----|----|
| C1 |    |    |    |    |    |
| C2 |    |    |    |    |    |
| C3 |    |    |    |    |    |
| C4 |    |    |    |    |    |
| C5 |    |    |    |    |    |

Proximity Matrix



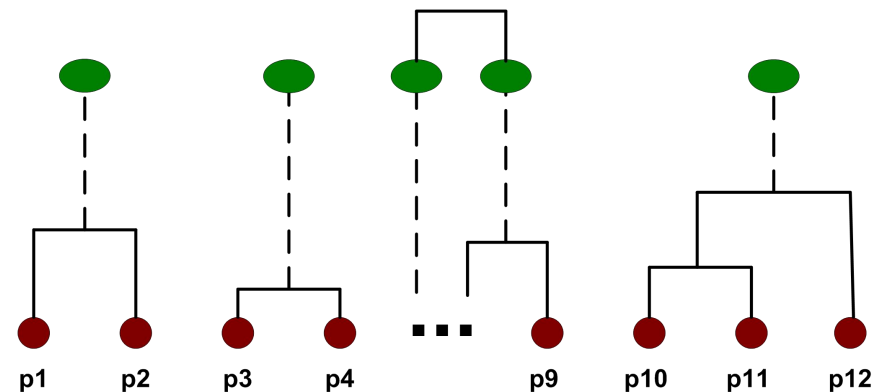
# 合并操作

- 合并操作：合并最近的两个簇(C2 and C5)，同时更新相似度矩阵



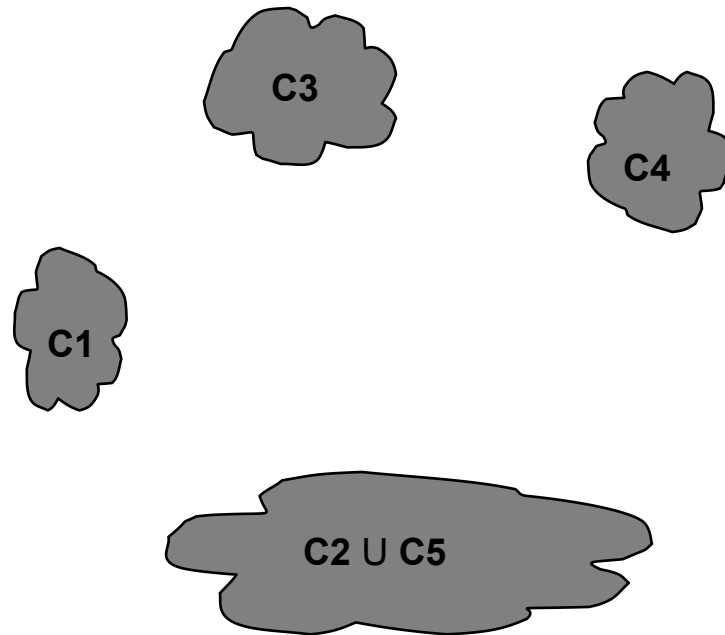
|    | C1 | C2 | C3 | C4 | C5 |
|----|----|----|----|----|----|
| C1 |    |    |    |    |    |
| C2 |    |    |    |    |    |
| C3 |    |    |    |    |    |
| C4 |    |    |    |    |    |
| C5 |    |    |    |    |    |

Proximity Matrix



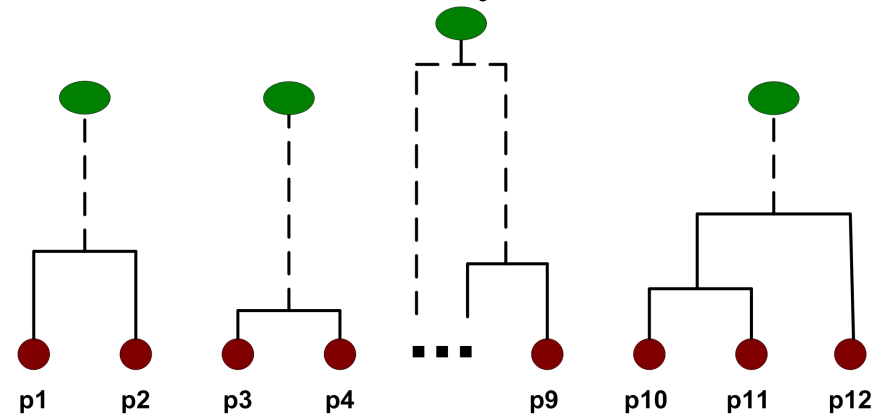
# 合并操作

- 问题：如何更新相似度矩阵？

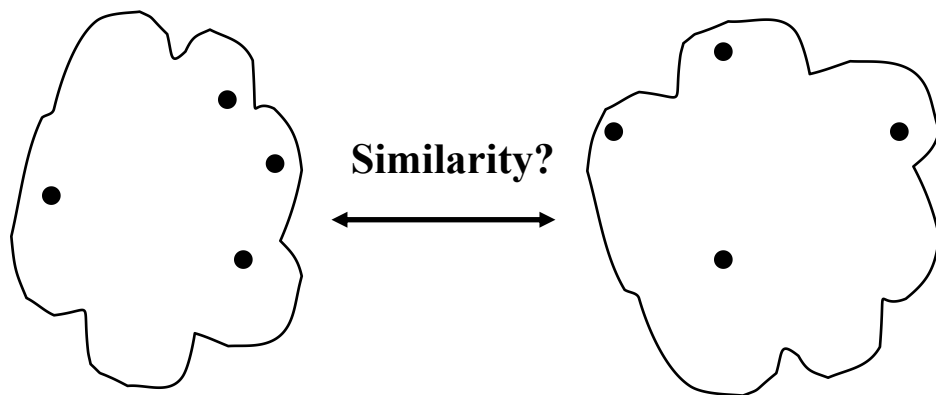


|                | C1 | <b>C2<br/>U<br/>C5</b> | C3 | C4 |
|----------------|----|------------------------|----|----|
| C1             |    | ?                      |    |    |
| <b>C2 U C5</b> | ?  | ?                      | ?  | ?  |
| C3             |    | ?                      |    |    |
| C4             |    | ?                      |    |    |

Proximity Matrix



# Inter-Cluster Similarity

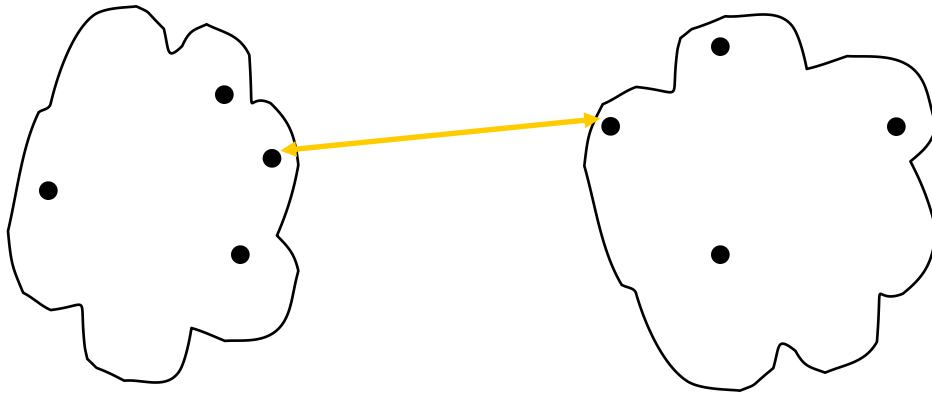


- 两个不相交的簇G和H，其间的相似度  $D(G, H)$  可以通过点点间的相似度 (pairwise similarities)  $D(i, j)$ ,  $i \in G$ ,  $j \in H$ , 计算得到

|    | p1 | p2 | p3 | p4 | p5 | ... |
|----|----|----|----|----|----|-----|
| p1 |    |    |    |    |    |     |
| p2 |    |    |    |    |    |     |
| p3 |    |    |    |    |    |     |
| p4 |    |    |    |    |    |     |
| p5 |    |    |    |    |    |     |
| .  |    |    |    |    |    |     |
| .  |    |    |    |    |    |     |
| .  |    |    |    |    |    |     |

**Proximity Matrix**

# Inter-Cluster Similarity

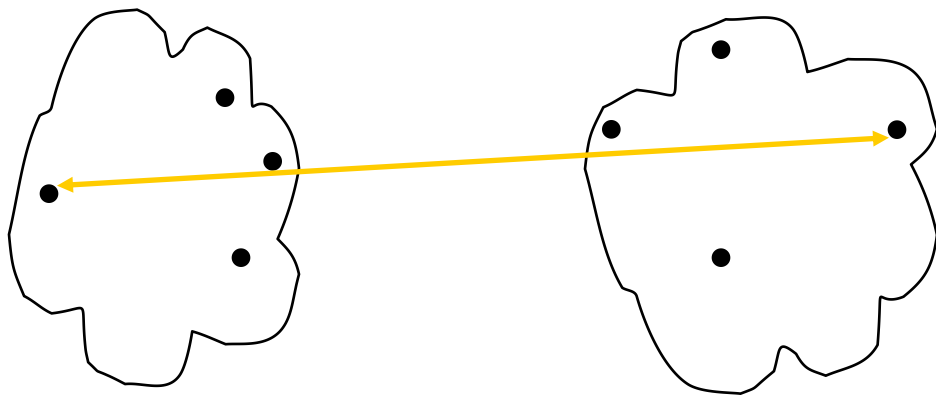


- MIN (single linkage)
- MAX (complete linkage)
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

|    | p1 | p2 | p3 | p4 | p5 | ... |
|----|----|----|----|----|----|-----|
| p1 |    |    |    |    |    |     |
| p2 |    |    |    |    |    |     |
| p3 |    |    |    |    |    |     |
| p4 |    |    |    |    |    |     |
| p5 |    |    |    |    |    |     |
| .  |    |    |    |    |    |     |

· **Proximity Matrix**  
·

# Inter-Cluster Similarity

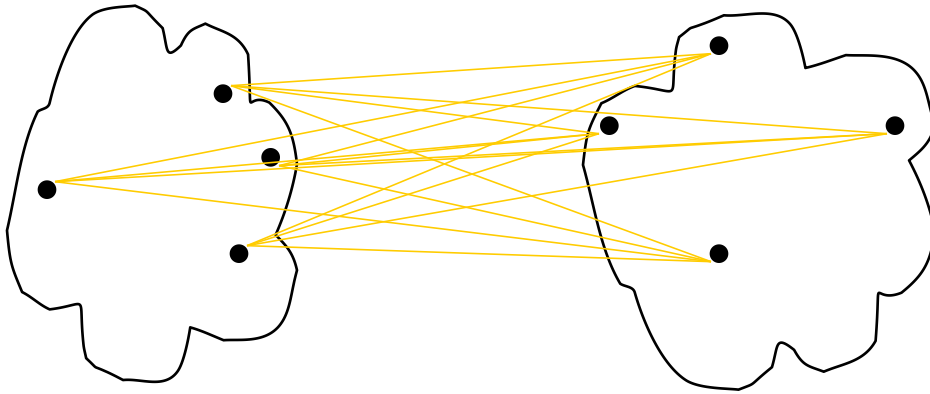


- MIN (single linkage)
- **MAX (complete linkage)**
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

|    | p1 | p2 | p3 | p4 | p5 | ... |
|----|----|----|----|----|----|-----|
| p1 |    |    |    |    |    |     |
| p2 |    |    |    |    |    |     |
| p3 |    |    |    |    |    |     |
| p4 |    |    |    |    |    |     |
| p5 |    |    |    |    |    |     |
| .  |    |    |    |    |    |     |

**Proximity Matrix**

# How to Define Inter-Cluster Similarity



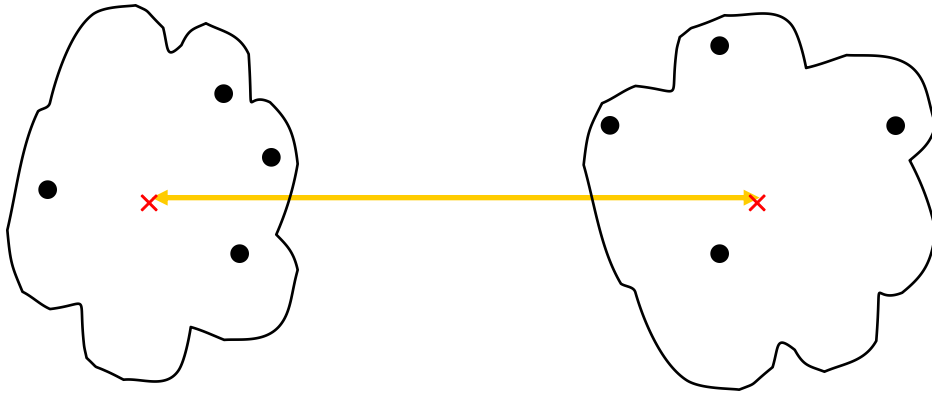
- MIN (single linkage)
- MAX (complete linkage)
- **Group Average**
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

|    | p1 | p2 | p3 | p4 | p5 | ... |
|----|----|----|----|----|----|-----|
| p1 |    |    |    |    |    |     |
| p2 |    |    |    |    |    |     |
| p3 |    |    |    |    |    |     |
| p4 |    |    |    |    |    |     |
| p5 |    |    |    |    |    |     |
| .  |    |    |    |    |    |     |

• **Proximity Matrix**

•

# Inter-Cluster Similarity



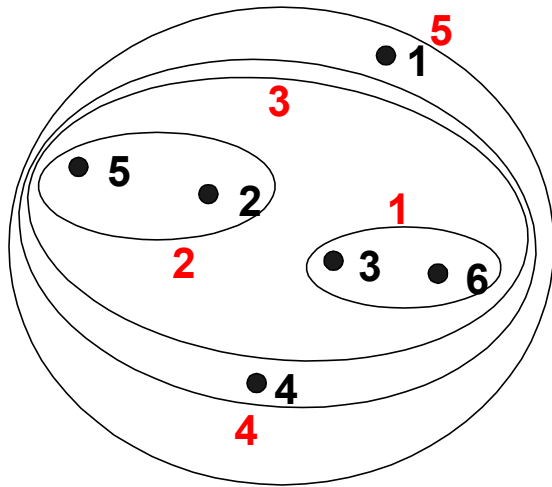
- MIN (single linkage)
- MAX (complete linkage)
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

|    | p1 | p2 | p3 | p4 | p5 | ... |
|----|----|----|----|----|----|-----|
| p1 |    |    |    |    |    |     |
| p2 |    |    |    |    |    |     |
| p3 |    |    |    |    |    |     |
| p4 |    |    |    |    |    |     |
| p5 |    |    |    |    |    |     |
| .  |    |    |    |    |    |     |

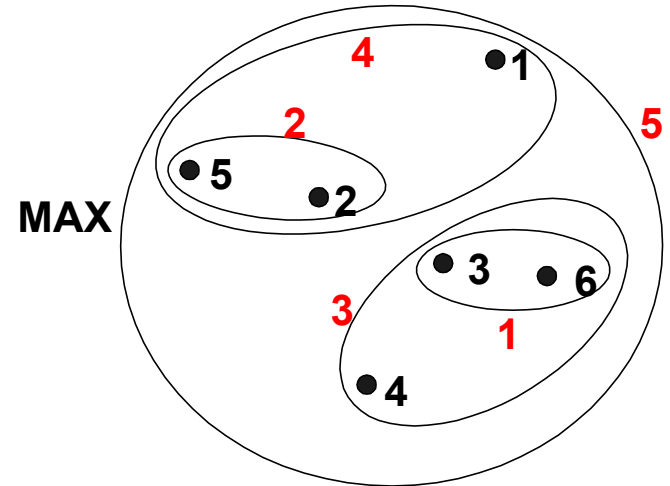
• **Proximity Matrix**



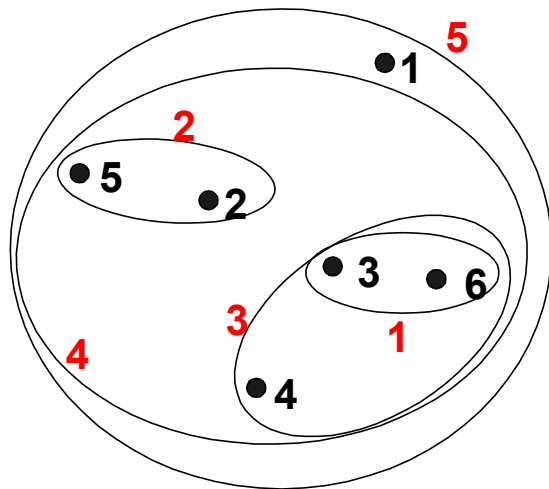
# Hierarchical Clustering: Comparison



MIN

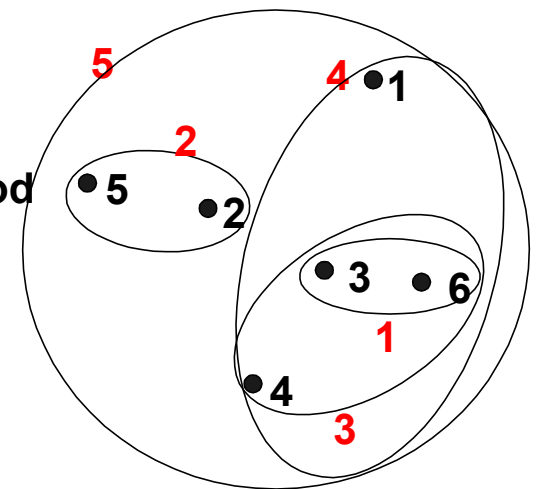


MAX



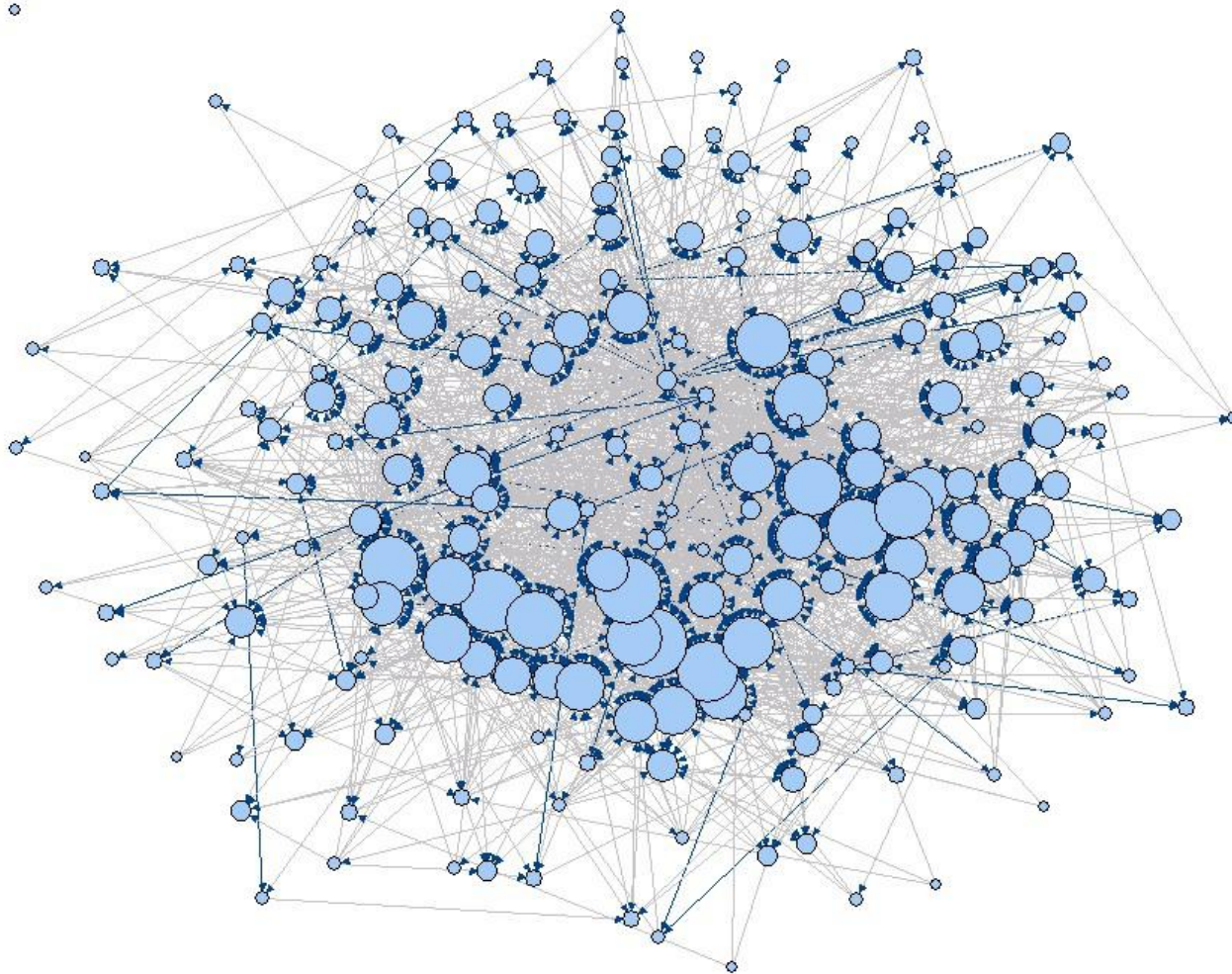
Group Average

Ward's Method



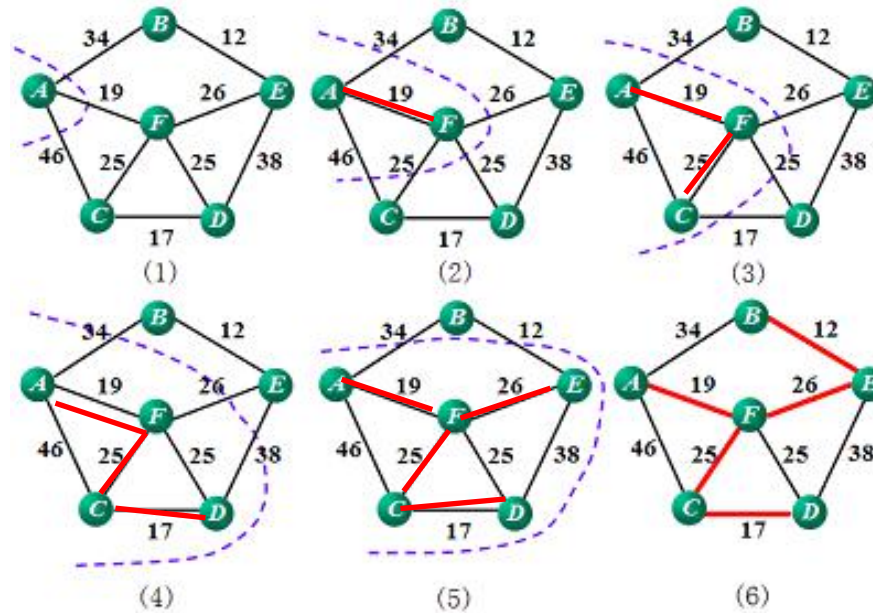
# Divisive Hierarchical Clustering

## Social Network Graphs



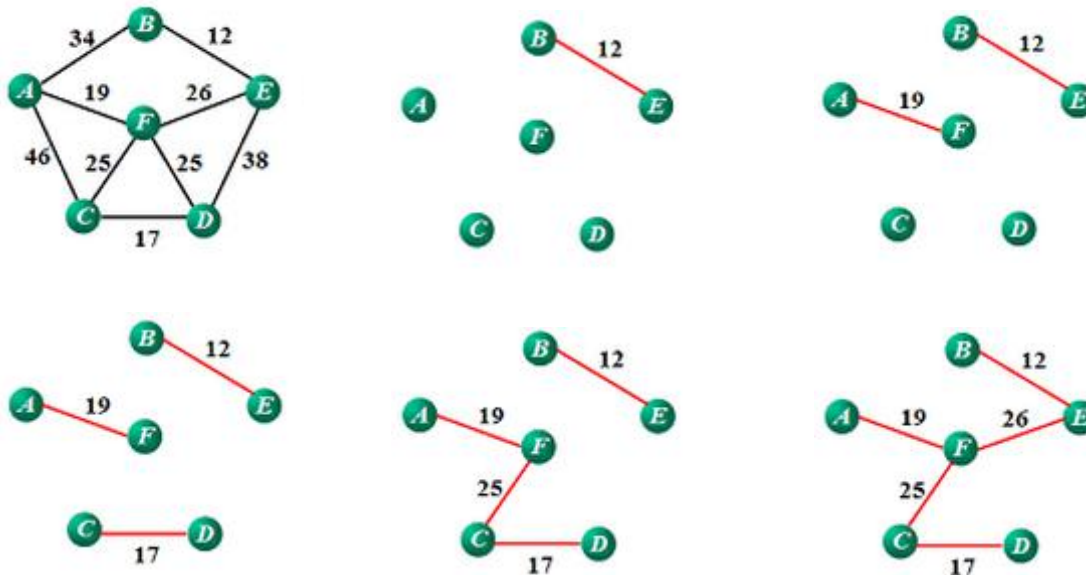
# Divisive Hierarchical Clustering

- E.g., in a MST approach
- 构建最小生成树 (Minimum Spanning Tree)
  - Prime算法



# Divisive Hierarchical Clustering

- E.g., in a MST approach
- 构建最小生成树 (Minimum Spanning Tree)
  - Kruskal 算法



# Hierarchical Clustering: Strengths

- Do not have to assume any particular number of clusters
  - Any desired number of clusters can be obtained by 'cutting' the dendrogram at the proper level
- They may correspond to meaningful taxonomies
  - Example in biological sciences (e.g., animal kingdom, phylogeny reconstruction, ...)

# Hierarchical Clustering: Problems and Limitations

- Once a decision is made to combine two clusters, it cannot be undone
- No objective function is directly minimized
- Different schemes have problems with one or more of the following:
  - Sensitivity to noise and outliers
  - Difficulty handling different sized clusters and convex shapes
  - Breaking large clusters

# Probabilistic Clustering

- Represent the probability distribution of the data as a *mixture model*
  - captures uncertainty in cluster assignments
  - gives model for data distribution
  - *Bayesian mixture model allows us to determine  $K$*
- Consider mixtures of *Gaussians*

# The Gaussian Distribution

- 单高斯模型(Gaussian Single Model)

- 一个随机变量 $x$ 服从高斯分布时，概率密度函数为：

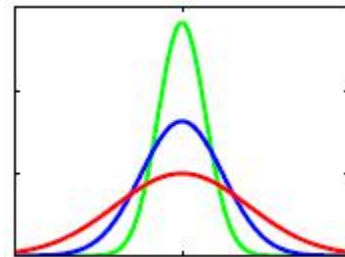
$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2}(x - \mu)^2 \right\}$$

- $\mu$ : 模型均值,  $\sigma^2$ 为模型方差

- 多维变量 $\mathbf{x}$ 服从高斯分布时，概率密度函数为：

$$\mathcal{N}(\mathbf{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu) \right\}$$

- $\mathbf{x}$ 是维度为 $D$ 的列向量,  $\mu$ : 模型均值,  $\Sigma$ 为 $D \times D$ 的协方差矩阵





# 似然函数

- 数据集:

$$D=\{x_i\}, i=1, \dots, N$$

- 考虑单个Gaussian模型
- 假设观测样本点由单个Gaussian独立等分布地抽样得到:

$$p(D|\mu, \Sigma) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \Sigma)$$

- 可以被看做模型参数的函数，因此被称为似然函数

# 极大似然估计

- 求解使得似然函数取最大值时对应的模型参数
- 等价地，极大化log似然函数：

$$\begin{aligned}\ln p(D|\mu, \Sigma) = & -\frac{N}{2} \ln |\Sigma| - \frac{Nd}{2} \ln(2\pi) \\ & - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \mu)^\top \Sigma^{-1} (\mathbf{x}_n - \mu)\end{aligned}$$

# 极大似然估计

- 相对于均值求似然最大化，得到样本均值：

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

- 相对于方差求似然最大化，得到样本方差：

$$\Sigma_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mu_{\text{ML}})(\mathbf{x}_n - \mu_{\text{ML}})^{\top}$$

# Gaussian混合模型

- Gaussian Mixture Model (GMM): 多个Gaussian模型的线性混合:

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

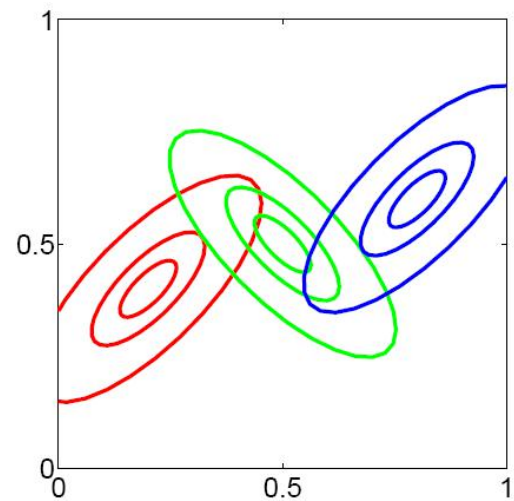
- 混合系数:  $\sum_{k=1}^K \pi_k = 1 \quad 0 \leq \pi_k \leq 1$

- 可以看做一种先验概率:

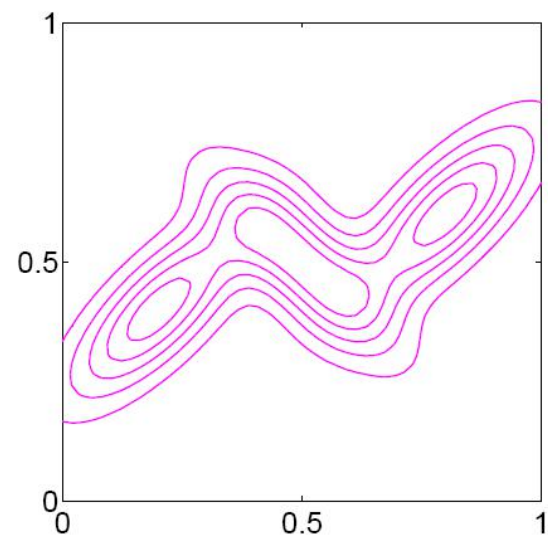
$$p(\mathbf{x}) = \sum_{k=1}^K p(k) p(\mathbf{x} | k)$$

# 举例

3个Gaussian的混合：

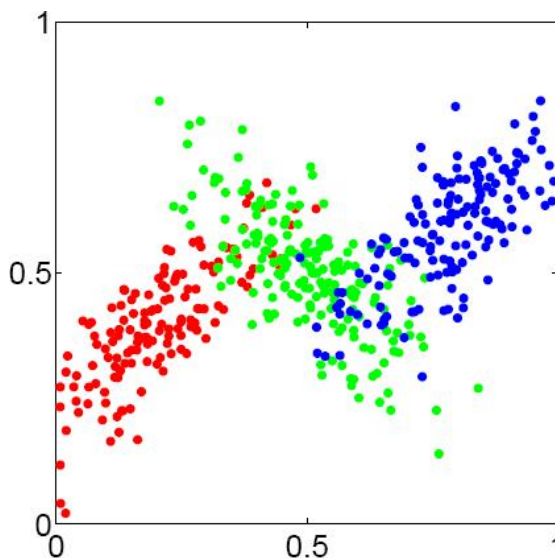


概率分布的等高线：



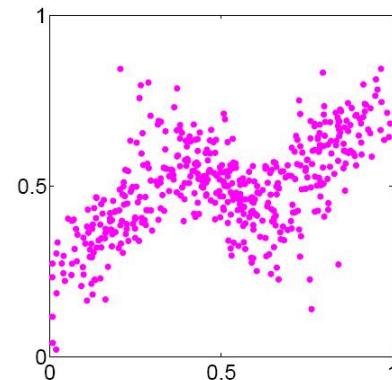
# 从Gaussian混合分布中抽样得到数据

- 样本点 $x_n$ 的生成过程：
  - 首先，以概率 $\pi_k$ 选择一个混合成分；
  - 接着，从该混合成分中抽样得到样本点 $x_n$
- 对于每个样本点，重复以上两个步骤



# 从数据中估计Gaussian混合分布的参数

- 求解上述过程的逆过程 – 给定样本点，估计相应的
  - 混合系数
  - 均值
  - 协方差

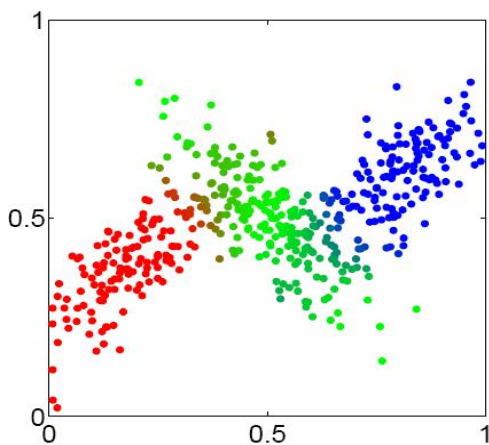


- 如果知道每个样本点由哪个成分抽样得到，则通过极大似然方法可以得到每个类对应的Gaussian模型的参数
- 问题：数据集缺少类别标注
- 因此类别labels可以看做是隐变量 (latent/hidden variable)

# 后验概率

- 可以把混合系数看做每个成分的先验概率
- 给定一个类别标记 $k$ ，可以估计相应的后验概率（posterior probabilities, 或*responsibilities*）
- 可以通过Bayes' theorem得到：

$$\begin{aligned}\gamma_k(\mathbf{x}) \equiv p(k|\mathbf{x}) &= \frac{p(k)p(\mathbf{x}|k)}{p(\mathbf{x})} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}\end{aligned}$$





# GMM的极大似然估计

- Log似然函数：

$$\ln p(D|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

- 注意：sum出现在log内
- 没有封闭解！

# GMM的极大似然估计

- 对log似然函数简单求导
- 令  $\ln p(X/\pi, \mu, \Sigma)$  相对于第k个Gaussian的均值  $\mu_k$  的倒数为零, 得:

$$0 = - \sum_{n=1}^N \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)}}_{\gamma(z_{nk})} \Sigma_k (\mathbf{x}_n - \mu_k)$$

$$\longrightarrow \mu_k = \frac{1}{N_k} \sum_{i=1}^N \gamma_k(x_i) x_i$$

$$N_k = \sum_{i=1}^N \gamma_k(x_i) \text{ 类别 } k \text{ 中所属的有效样本个数}$$

# GMM的极大似然估计

- 类似的，求解协方差：

$$\Sigma_k = \frac{1}{N_k} \sum_{i=1}^N \gamma_k(x_i) (x_i - \mu_k)(x_i - \mu_k)^T$$

- 根据Lagrange multiplier求解混合系数：

$$\pi_k = \frac{1}{N} \sum_{i=1}^N \gamma_k(x_i)$$

# EM Algorithm

- 上述解构不成封闭形式，因为变量之间互为耦合
- 采用一种迭代的方式求解：
  - 给参数一个初始值
  - 通过下述两个步骤更新参数：
    - E-step: 估计后验概率或responsibilities
    - M-step: 根据MLE的结果更新参数
- 每一次 EM 循环都能保证likelihood值增大

# EM Algorithm

- E-step: 估计responsibilities

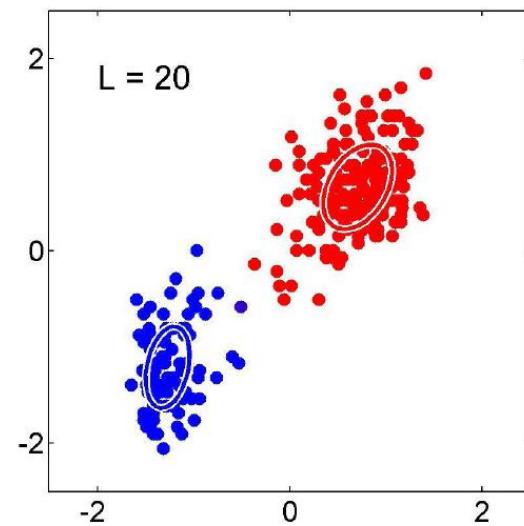
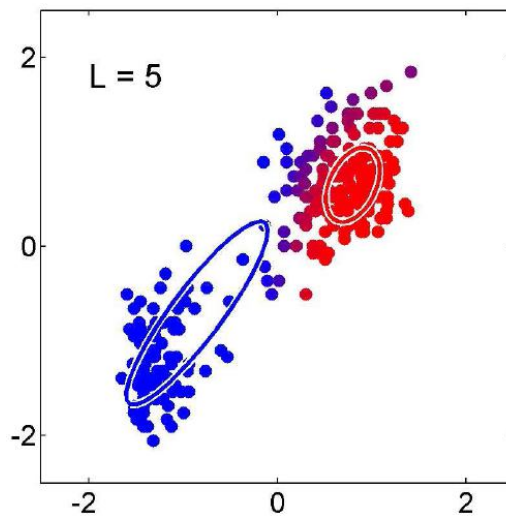
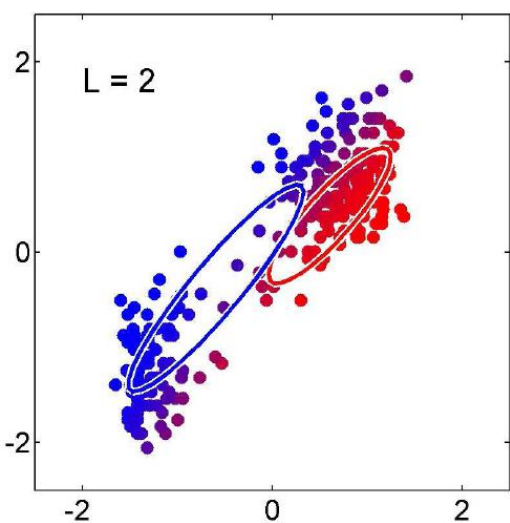
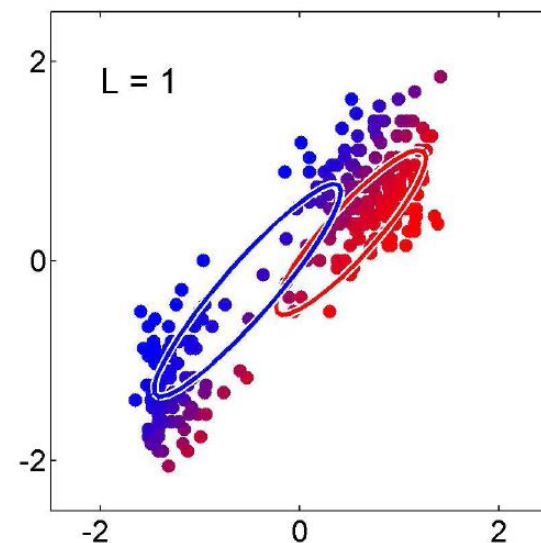
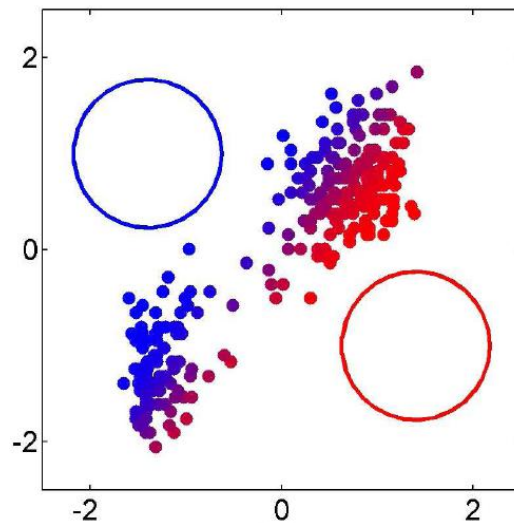
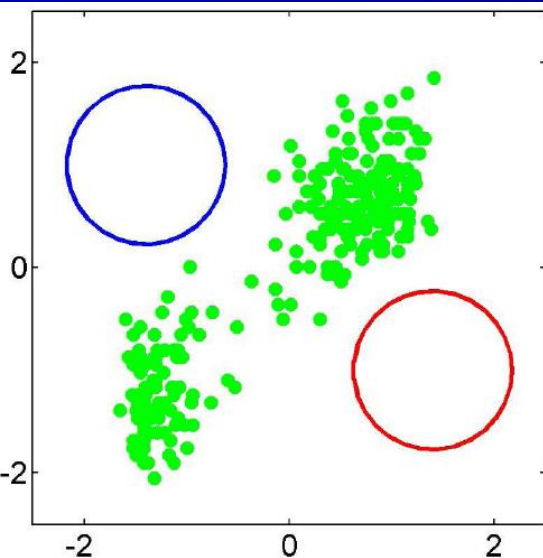
$$\gamma_k(x_i) = \frac{\pi_k N(x_i | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_i | \mu_j, \Sigma_j)}$$

- M-step: 采用MLE估计更新参数

$$\mu_k = \frac{1}{N_k} \sum_{i=1}^N \gamma_k(x_i)$$

$$\Sigma_k = \frac{1}{N_k} \sum_{i=1}^N \gamma_k(x_i) (x_i - \mu_k)(x_i - \mu_k)^T$$

# Example



# GMM应用于分类

- 分类时：
  - 每个混合成分（类）的参数 $\mu$ 和 $\Sigma$ 已知
  - 把数据点 $x$ 带入到每个混合成分 $C_k$ 中
$$N(x | \mu_k, \Sigma_k)$$
  - 当概率大于一定阈值时便认为 $x$ 属于 $C_k$ 类

# 与 K-means之间的关系

- 考虑GMM的协方差为一个常数  $\epsilon$
- 令极限  $\epsilon \rightarrow 0$
- Responsibilities取两个值:

$$\gamma_i(\mathbf{x}_n) = \frac{\pi_i \exp \left\{ -\|\mathbf{x}_n - \boldsymbol{\mu}_i\|^2 / 2\epsilon \right\}}{\sum_j \pi_j \exp \left\{ -\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2 / 2\epsilon \right\}} \rightarrow r_{ni} \in \{0, 1\}$$

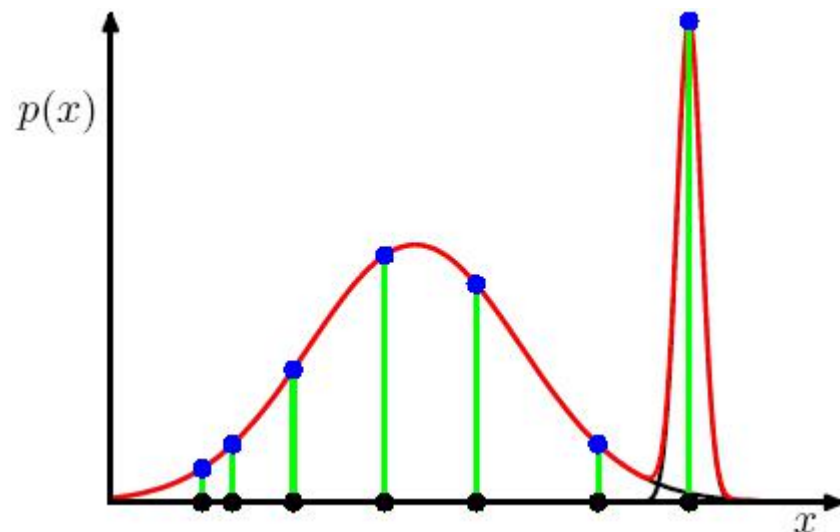
- 此时, EM algorithm与K-means等价



# Other issues: Over-fitting in Gaussian Mixture Models

- 假设：某个混合成分仅包括一个样本点

$$\Sigma_j = \sigma^2 \mathbf{I}, \mu_j = \mathbf{x}_n$$
$$\mathcal{N}(\mathbf{x}_n | \mathbf{x}_n, \sigma_j^2 \mathbf{I}) = \frac{1}{(2\pi)^{1/2}} \frac{1}{\sigma_j}$$



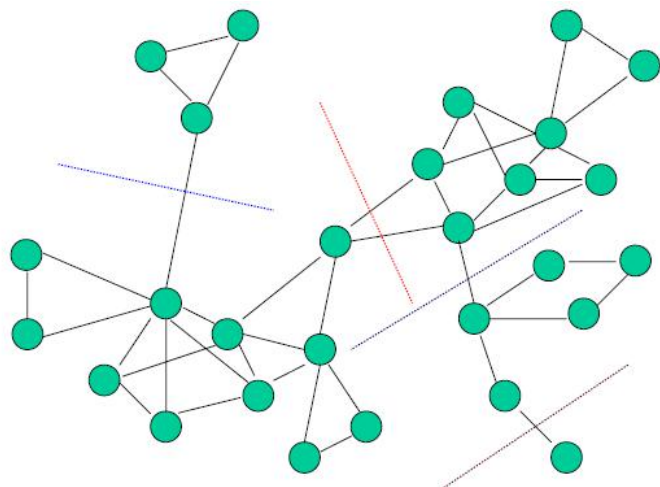
- 考虑情形  $\sigma^2 \rightarrow 0$
- 此时  $\log$  似然函数趋于无穷大，则极大化似然估计不是一个良态问题

# Spectral Clustering

- Recall that the intuitive goal of clustering is
  - to divide the data points into several groups such that points in the same group are **similar** and points in different groups are **dissimilar** to each other
- Suppose the data points  $\{x_1, \dots, x_n\}$  are organized by **similarity graph  $G = (V, E, W)$** 
  - Each vertex  $v_i$  in this graph represents a data point  $x_i$
  - Each edge  $e_{ij}$  is weighted by similarity  $s_{ij}$  (or  $w_{ij}$ ) between  $v_i$  and  $v_j$
- The graph construction depends on the application

# Spectral Clustering

- The problem of clustering can now be reformulated using the **similarity graph**:
  - to find a partition of the graph such that the edges between different groups have very low weights and the edges within a group have high weights



# Graph partitioning

- Clustering partitions the vertices of the graph.
- A good clustering places dissimilar vertices in different partitions.
- The **loss function for a partition of  $(A, \bar{A})$**  (or the weighted adjacency of  $(A, \bar{A})$ ) is given by the cut:

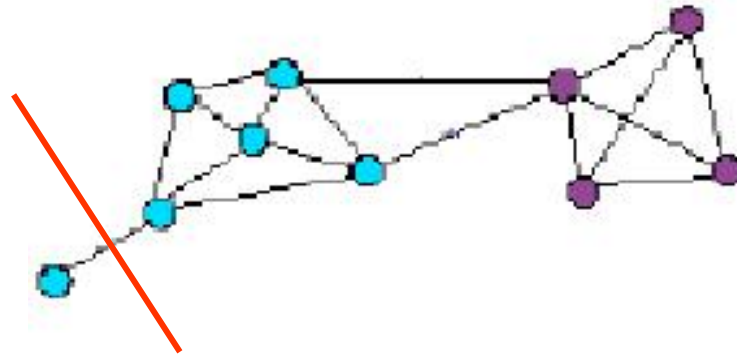
$$cut(A, \bar{A}) = \sum_{i \in A, j \in \bar{A}} w_{ij}$$

- For a given number  $k$  of subsets, the mincut approach simply consists in choosing a partition  $A_1, \dots, A_k$  which minimizes

$$cut(A_1, \dots, A_k) = \sum_{i \in A_i, j \in \bar{A}_i} w_{ij}$$

- Find a partition **that minimizes the cut (Mincut criterion) -- HOW?**

- In many cases, the solution of mincut simply separates one individual vertex from the rest of the graph.



e.g., 2-way Partitioning...

- Of course this is not what we want to achieve in clustering
- Explicitly request that the sets  $A_1, \dots, A_k$  are "reasonably large"

Two ways of measuring the "size" of a subset  $A$ :

- $|A| :=$  the number of vertices in  $A$
- $\text{vol}(A) := \sum d_i$ , for all  $i$  in  $A$

# Graph partitioning

- A loss function that favors such clusters is Normalized cut

$$Ncut = \frac{1}{2} \sum_{i=1}^k \frac{W(A_i, \bar{A}_i)}{vol(A_i)} = \sum_{i=1}^k \frac{cut(A_i, \bar{A}_i)}{vol(A_i)}$$

$$- vol(A) = \sum_{i \in A} d_i, d_i = \sum_{j=1, \dots, n} w_{ij}$$

- **A good partition** should separate **dissimilar** vertices and should produce **balanced** clusters
- **Minimizing normalized cut is NP-hard.**
- One way of approximately optimizing normalized cuts leads to **spectral clustering**

# 2-way Spectral Graph Partitioning

- Partition membership indicator:  $q_i = \begin{cases} 1 & \text{if } i \in A \\ -1 & \text{if } i \in B \end{cases}$

$$\begin{aligned} J = \text{CutSize} &= \frac{1}{4} \sum_{i,j} w_{ij} [q_i - q_j]^2 \\ &= \frac{1}{4} \sum_{i,j} w_{ij} [q_i^2 + q_j^2 - 2q_i q_j] = \frac{1}{2} \sum_{i,j} q_i [d_i \delta_{ij} - w_{ij}] q_j \\ &= \frac{1}{2} q^T (D - W) q \end{aligned}$$

$$D = \text{diag}(d_1, \dots, d_n), d_i = \sum_{j=1, \dots, n} w_{ij}$$

- Relax indicators  $q_i$  from discrete values to continuous values, the solution for  $\min J(q)$  is given by the eigenvectors of

$$(D - W)q = \lambda q$$

# Properties of Graph Laplacian (Proof omitted)

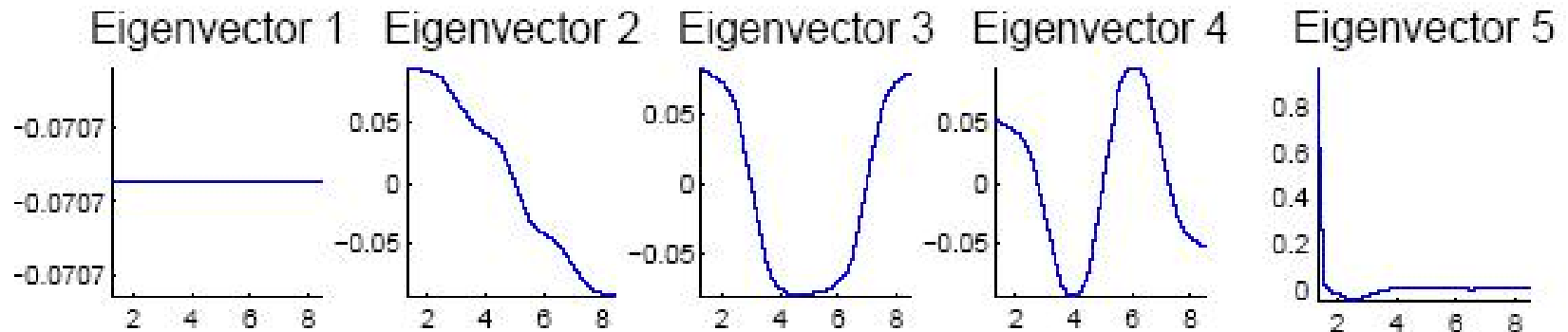
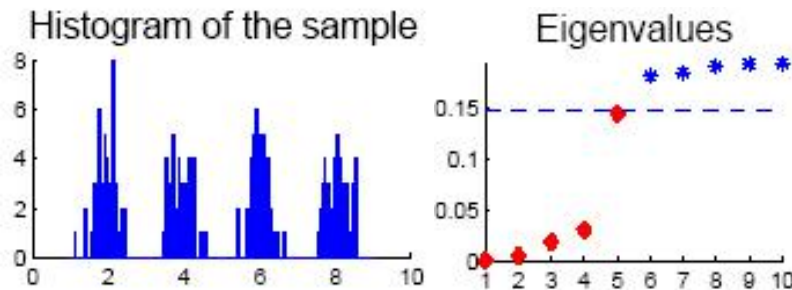
- Laplacian matrix of the Graph:  $L = D - W$
- $L$  is semi-positive definite  $x^T L x \geq 0$  for any  $x$ .
- Assume eigenvalues will always be ordered increasingly
  - First eigenvector is  $q_1 = (1, \dots, 1)^T$  with  $\lambda_1 = 0$ .
  - The eigenvectors with small eigenvalues provide important information for clustering
    - Second eigenvector  $q_2$  can serve as the desired solution
    - Higher eigenvectors are also useful



# Recovering Partitions

- From the definition of cluster indicators, partitions A, B are determined by:

$$A = \{i \mid q_2(i) < 0\}, B = \{i \mid q_2(i) \geq 0\}$$



# Spectral Clustering

- In most common view:
  - Build a weighted graph  $G = (V, E, W)$ .
  - Construct a Laplacian matrix  $L=f(W)$  (different variants of spectral clustering result from different functions  $f$ ).
  - Compute the eigenvectors of  $k$  smallest eigenvalues of  $L$ . These provide a new representation of the original data points.
  - Cluster the points in this new representation (e.g. using k-means).

# Spectral Clustering

## Spectral clustering algorithm

---

**Input:** Similarity matrix  $S \in \mathbb{R}^{n \times n}$ , number  $k$  of clusters to construct.

1: **Construct a similarity graph**  $G=(V, E)$ . Let  $W$  be its weighted adjacency matrix.

2: **Compute the Laplacian matrix**  $L$ .

3: **Compute the  $k$  eigenvectors**  $u_1, \dots, u_k$  **with  $k$  smallest eigenvalues** of  $L$ .

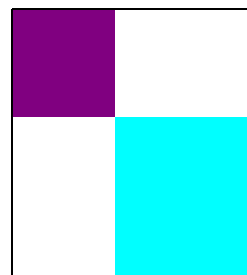
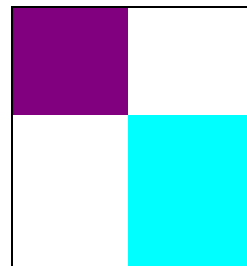
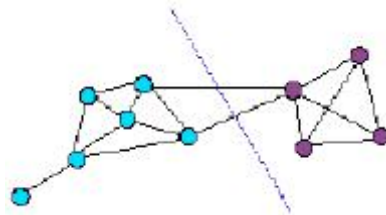
3: Let  $U \in \mathbb{R}^{n \times k}$  be the matrix containing the vectors  $u_1, \dots, u_k$  as columns.

4: For  $i = 1, \dots, n$ , let  $y_i \in \mathbb{R}^k$  be the vector corresponding to the  $i$ -th row of  $U$ .  
(for  $i = 1, \dots, n$ ,  **$x_i$  can be represented by  $y_i = [U_{i1}, \dots, U_{ik}]$** ).

5: **Cluster the points  $(y_i)_{i=1, \dots, n}$  in  $\mathbb{R}^k$**  with the  $k$ -means algorithm into clusters  $C_1, \dots, C_k$ .

**Output:** Clusters  $A_1, \dots, A_k$  with  $A_i = \{j \mid y_j \in C_i\}$ .

---



$l_1, \dots, l_k$

similarity matrix  
 $W$

transition matrix  
 $L = D - W$

first  $k$  eigenvalue of  $L$   
 $l_1, \dots, l_k$

cosine similarity

normalize rows

spectral mapping

clustering in  $R^k$

# Summary

- Clustering is cool
- It's easy to find the most salient pattern
- It's quite hard to find the pattern you want
- It's hard to know how to fix when broken
- EM is a useful optimization technique you should understand well if you don't already