

自然语言处理导论

#L9

文本聚类

袁彩霞

yuancx@bupt.edu.cn

智能科学与技术中心

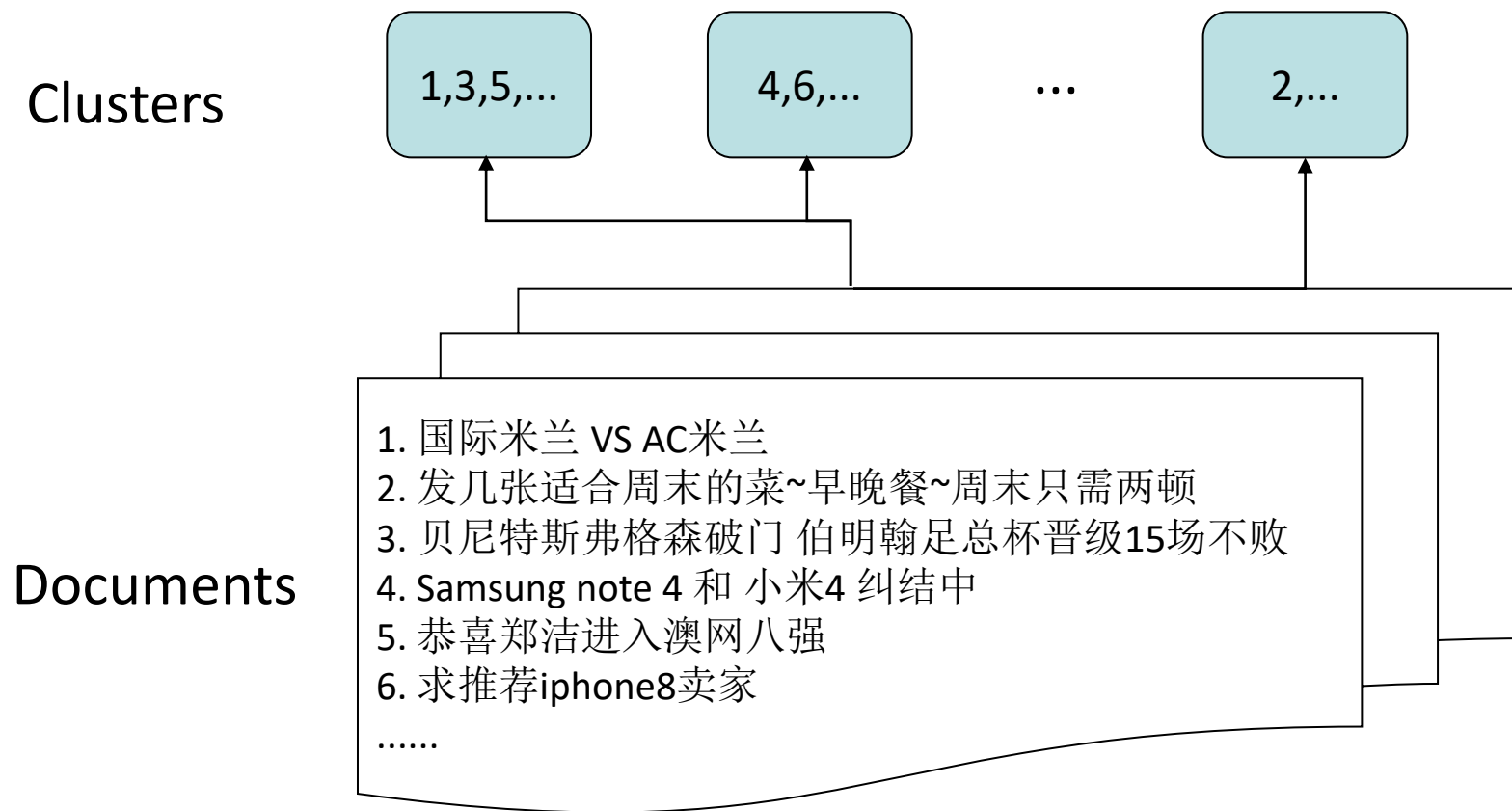
Overview

- 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



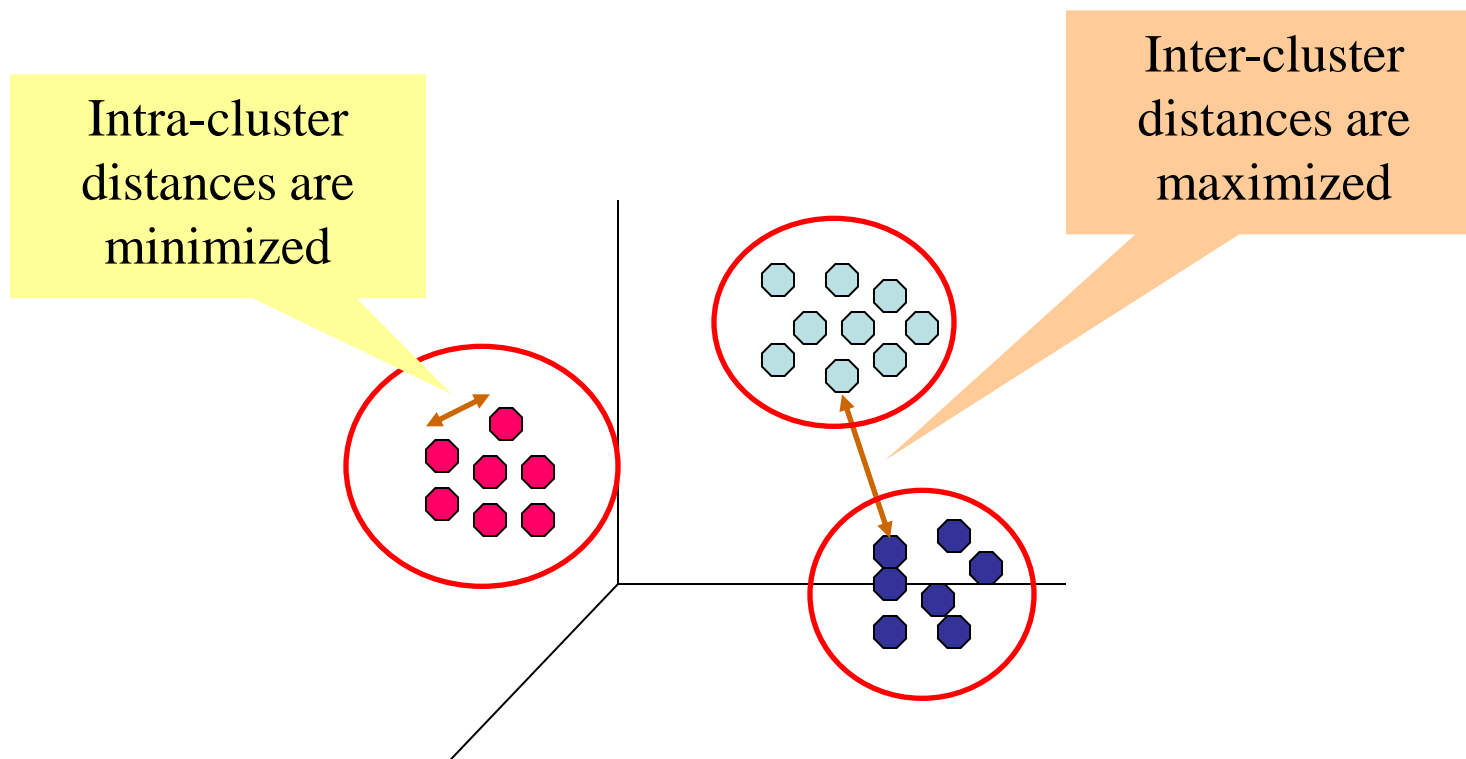
例子：web搜索优化

- 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

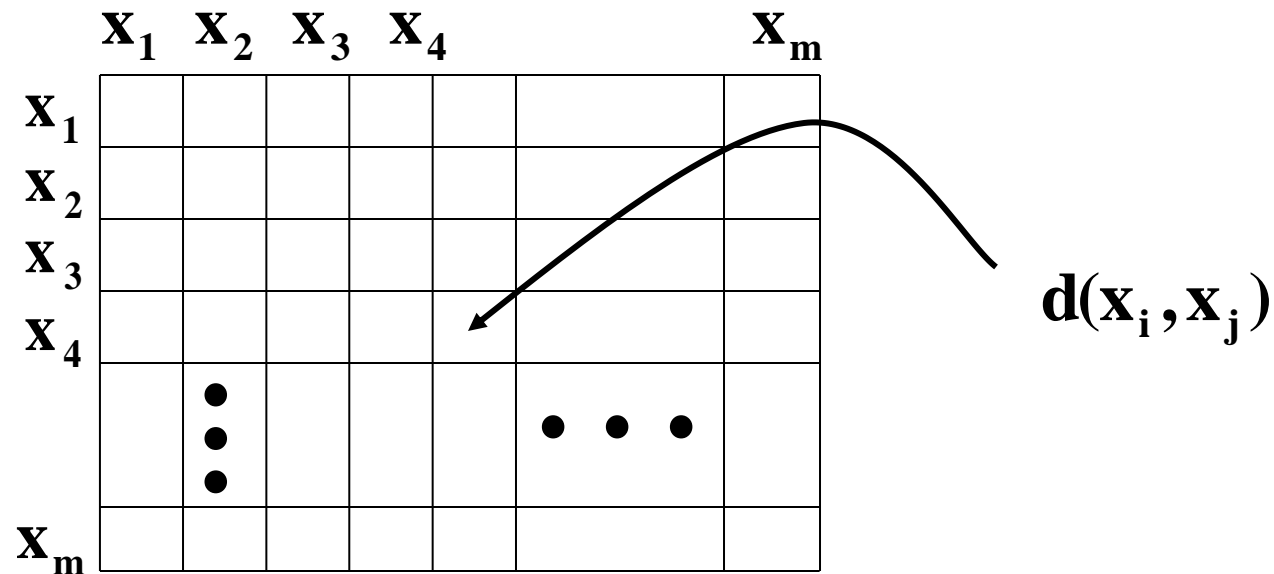
聚类



- 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

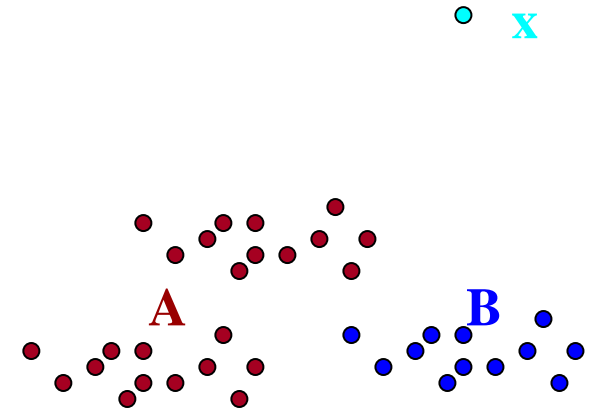
- 样本点 x 和样本点集合 A 之间的距离:

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

- 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
-

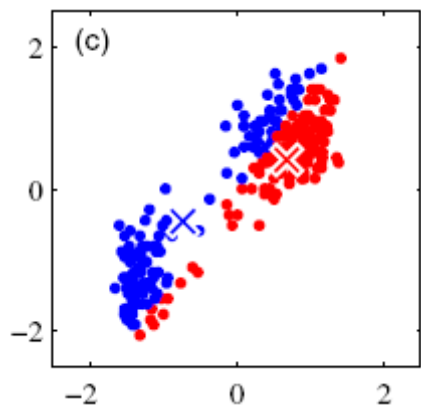
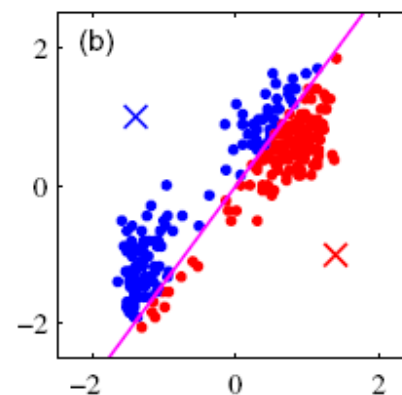
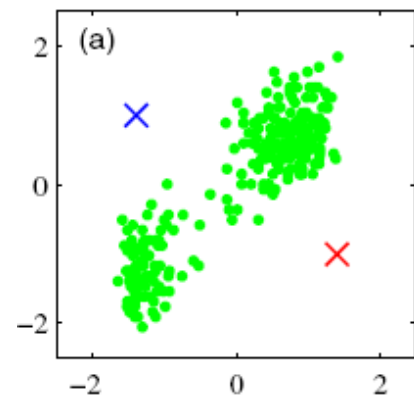
K-means

- 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:**

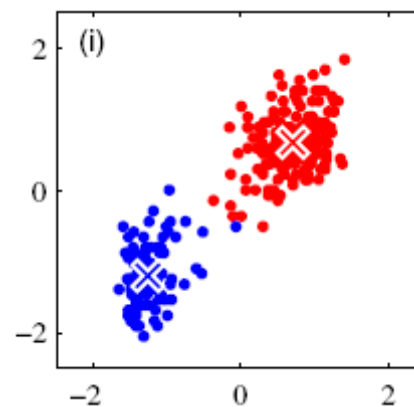
The basic algorithm:

- 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）

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

- μ_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 μ_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 μ_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 μ_k

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

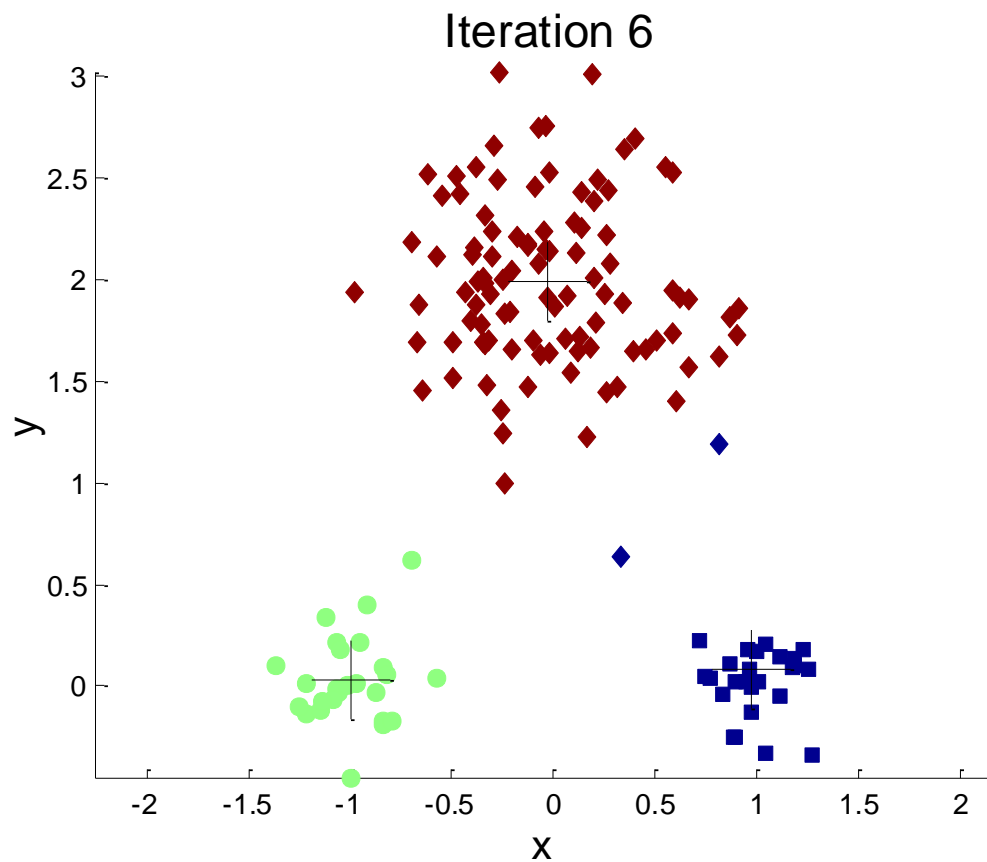
- 5: **until** convergence

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 !

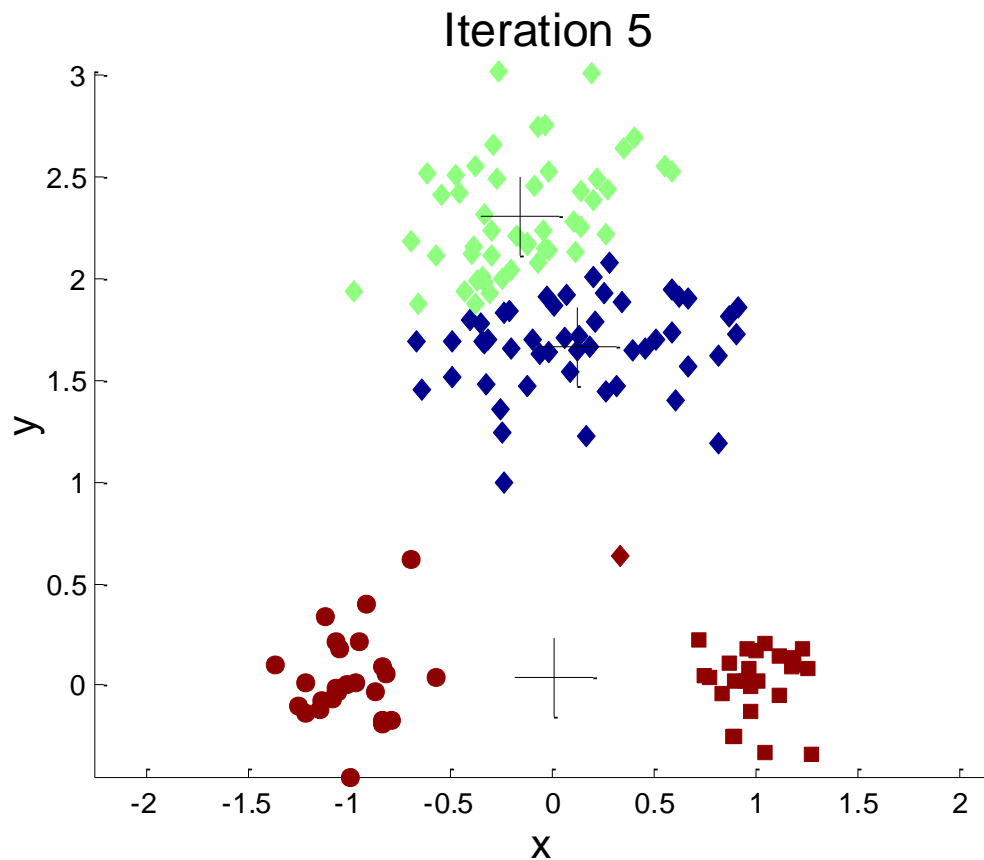
K-means – Details

不同初始中心点的选取对聚类结果的影响:



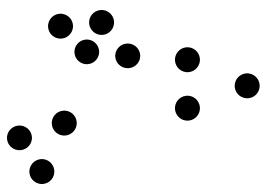
K-means – Details

不同初始中心点的选取对聚类结果的影响:

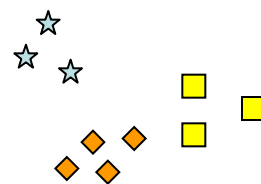
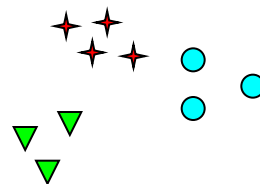
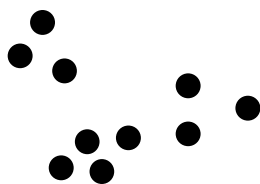


K-means – Details

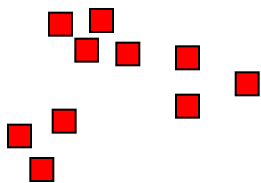
类别本身可能存在歧义：



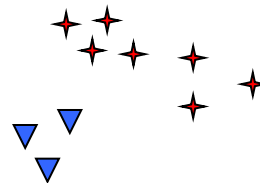
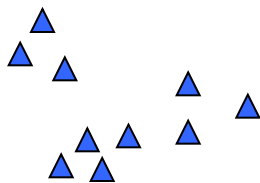
有多少个类？



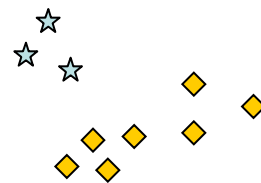
Six Clusters



Two Clusters



Four Clusters

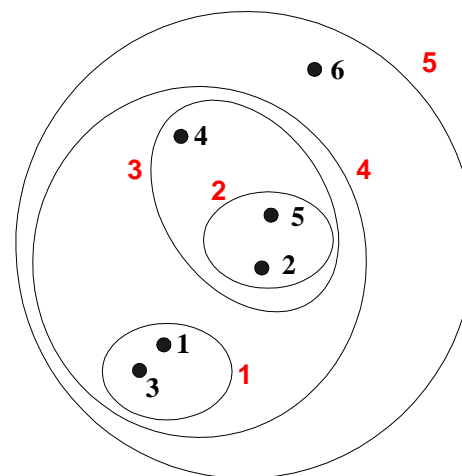
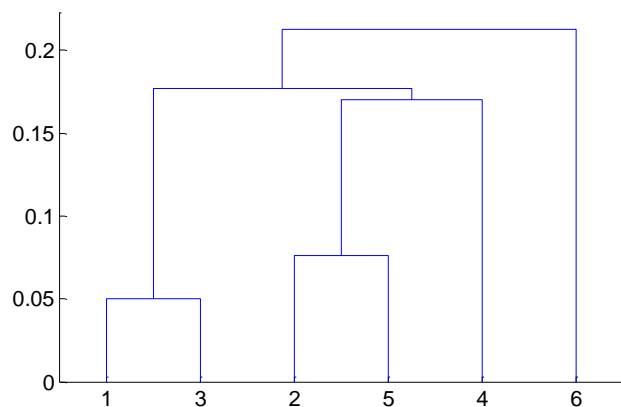


聚类模型

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

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

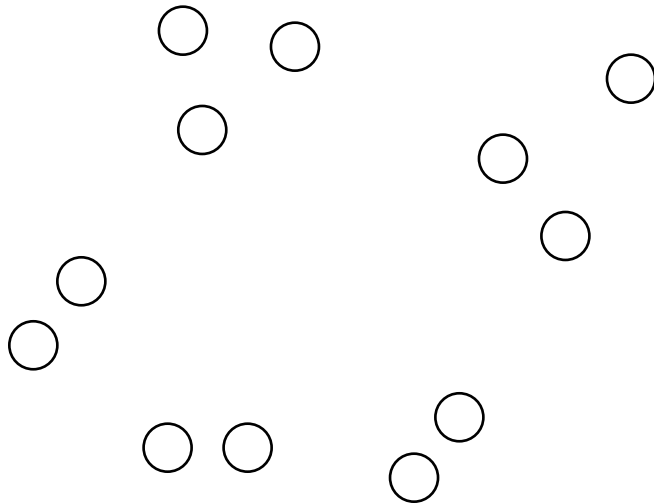
Basic algorithm:

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**)



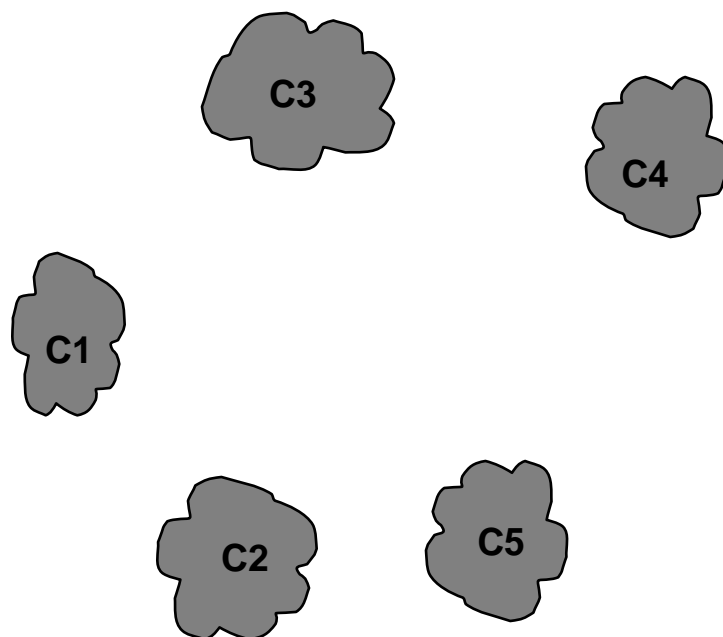
	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						
.						

Proximity Matrix



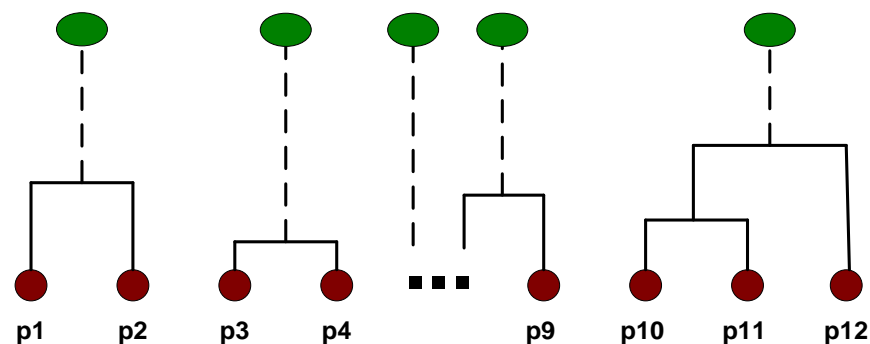
中间状态

- 经过几次合并操作，得到一些簇以及一个相似度矩阵



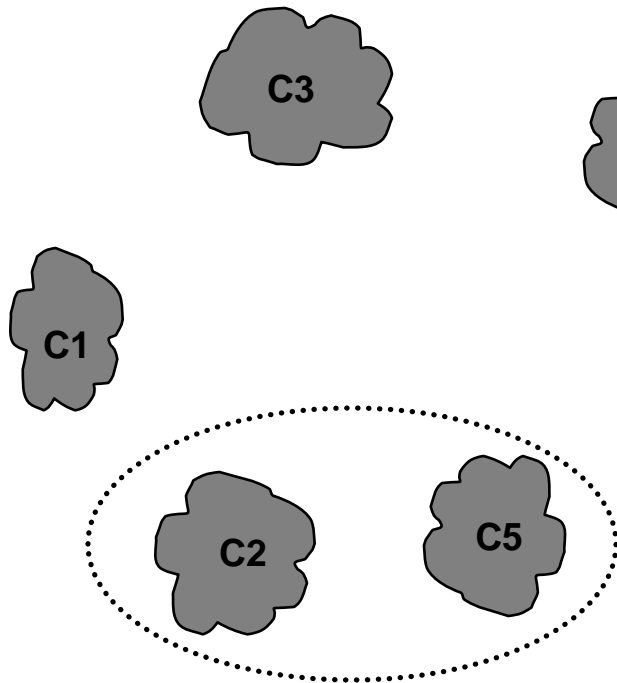
	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Proximity Matrix



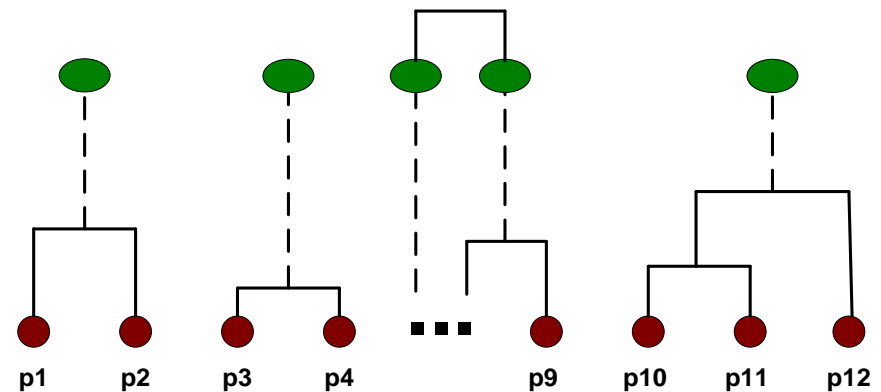
合并操作

- 合并操作：合并最近的两个簇(C2 and C5)，同时更新相似度矩阵



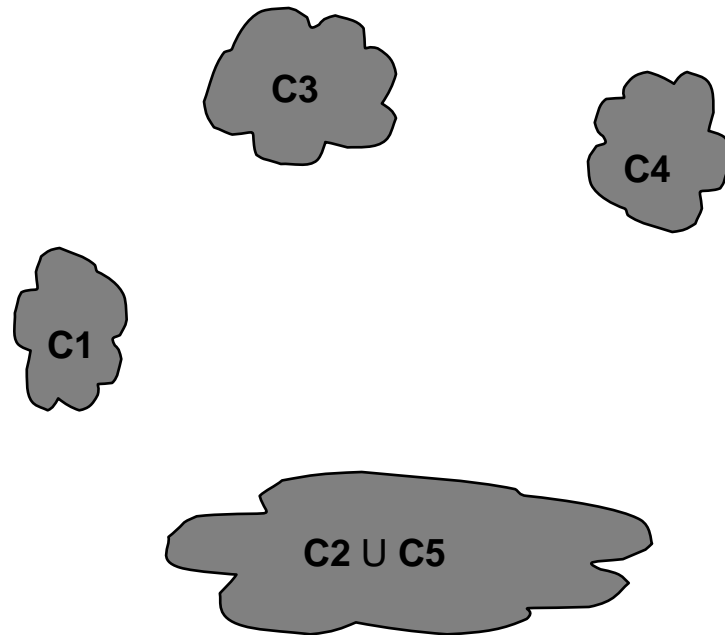
	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Proximity Matrix



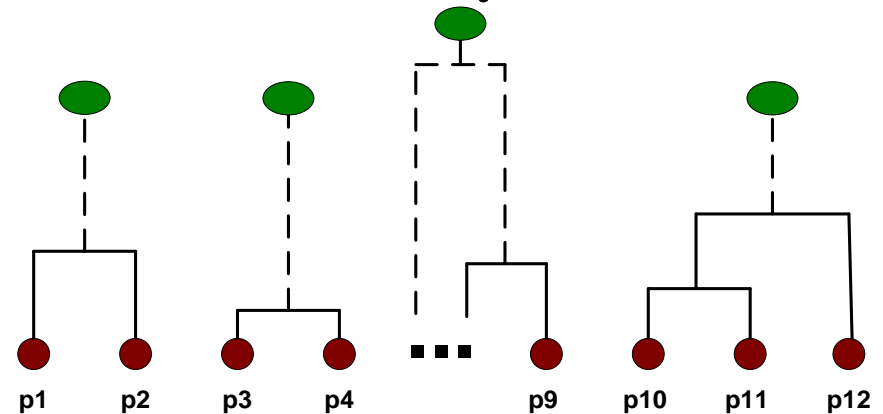
合并操作

- 问题：如何更新相似度矩阵？

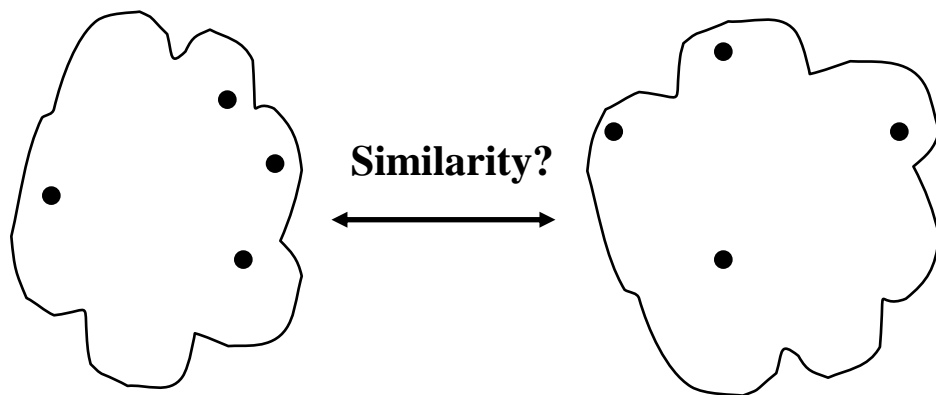


	C1	C2 U C5	C3	C4
C1		?		
C2 U C5	?	?	?	?
C3		?		
C4		?		

Proximity Matrix



Inter-Cluster Similarity

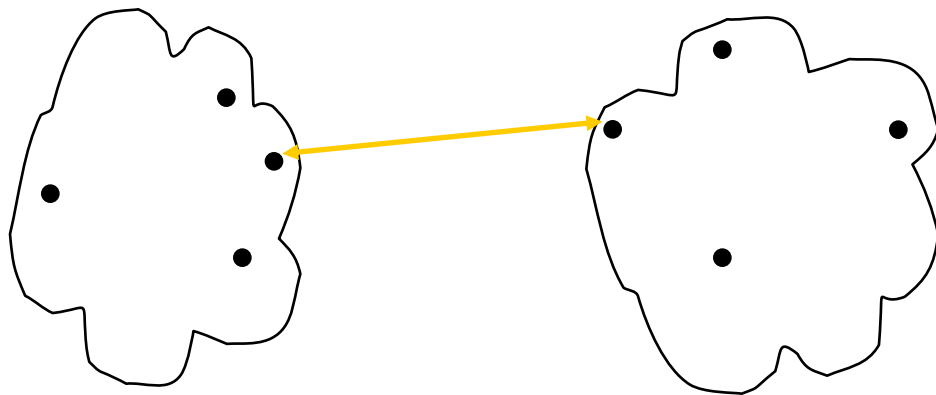


- 两个不相交的簇G和H，其间的相似度 $D(G, H)$ 可以通过点点间的相似度 (pairwise similarities) $D(i, j)$, $i \in G$, $j \in H$, 计算得到

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						

Proximity Matrix

Inter-Cluster Similarity

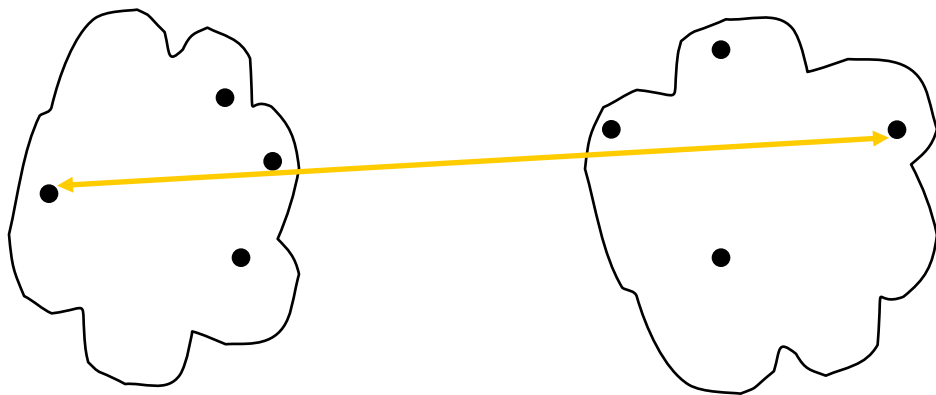


- MIN (single linkage)
- MAX (complete linkage)
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
 - Ward's Method uses squared error

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						

· **Proximity Matrix**
·

Inter-Cluster Similarity

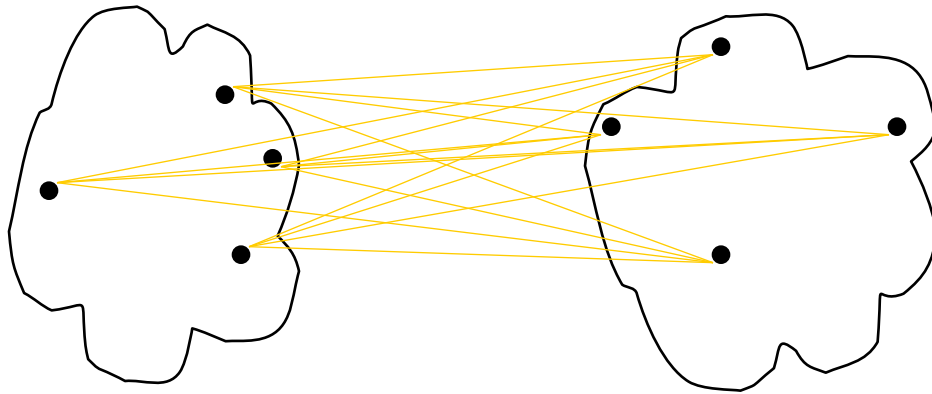


- MIN (single linkage)
- **MAX (complete linkage)**
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
 - Ward's Method uses squared error

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						

Proximity Matrix

How to Define Inter-Cluster Similarity



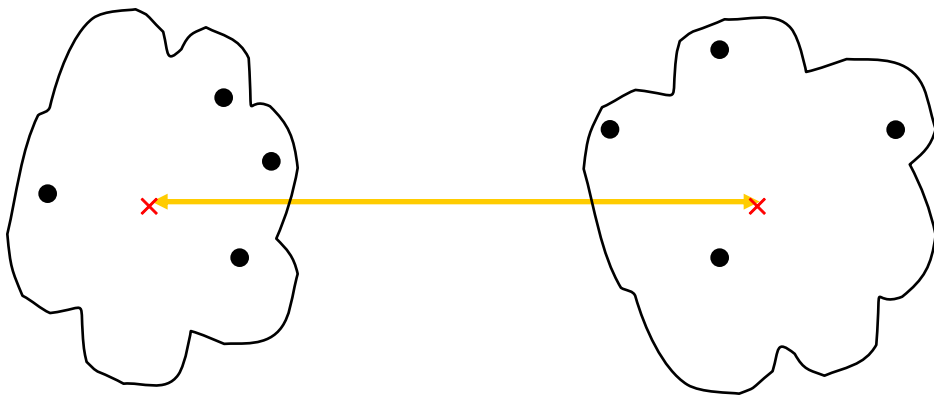
- MIN (single linkage)
- MAX (complete linkage)
- **Group Average**
- Distance Between Centroids
- Other methods driven by an objective function
 - Ward's Method uses squared error

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						

• **Proximity Matrix**

•

Inter-Cluster Similarity



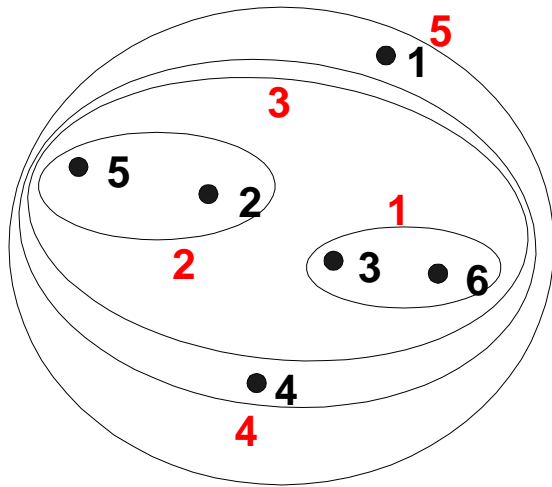
- MIN (single linkage)
- MAX (complete linkage)
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
 - Ward's Method uses squared error

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						

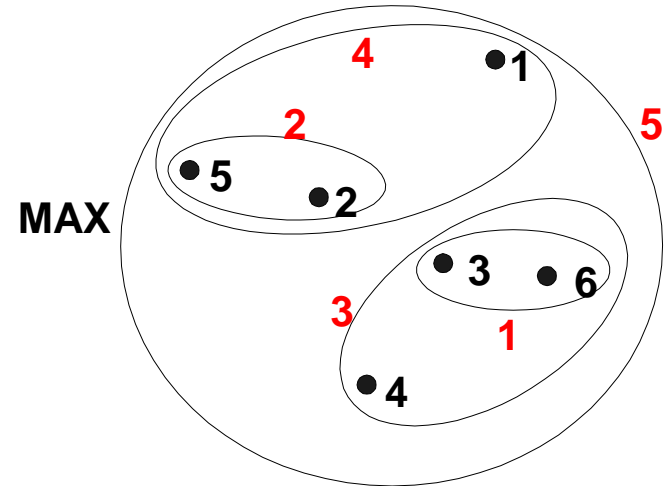
• **Proximity Matrix**

•

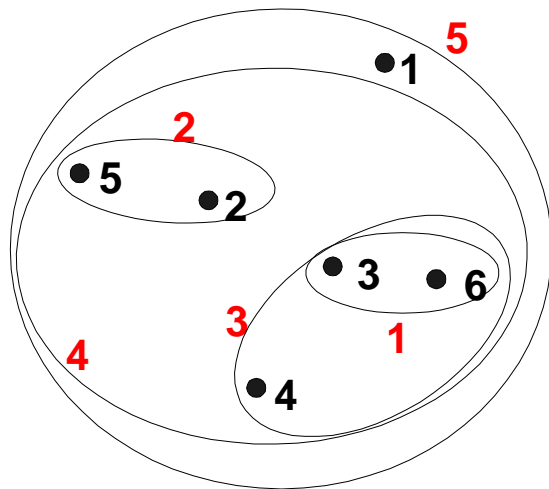
Hierarchical Clustering: Comparison



MIN

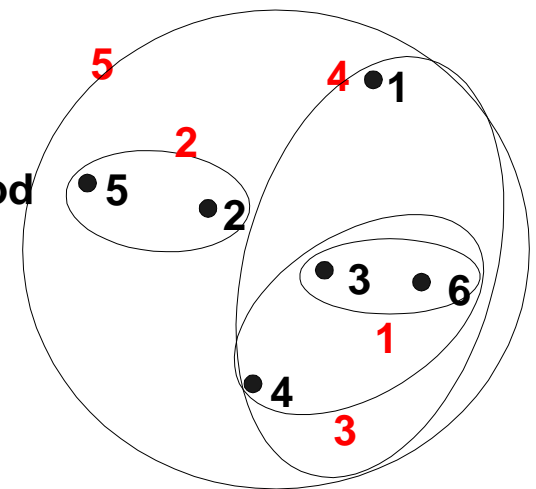


MAX



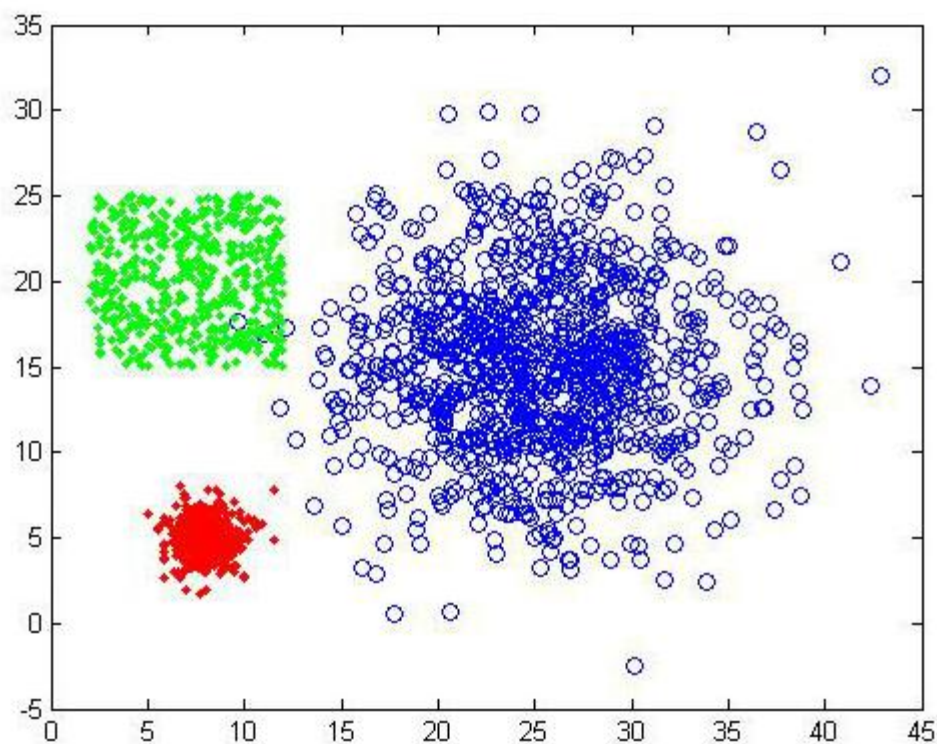
Group Average

Ward's Method



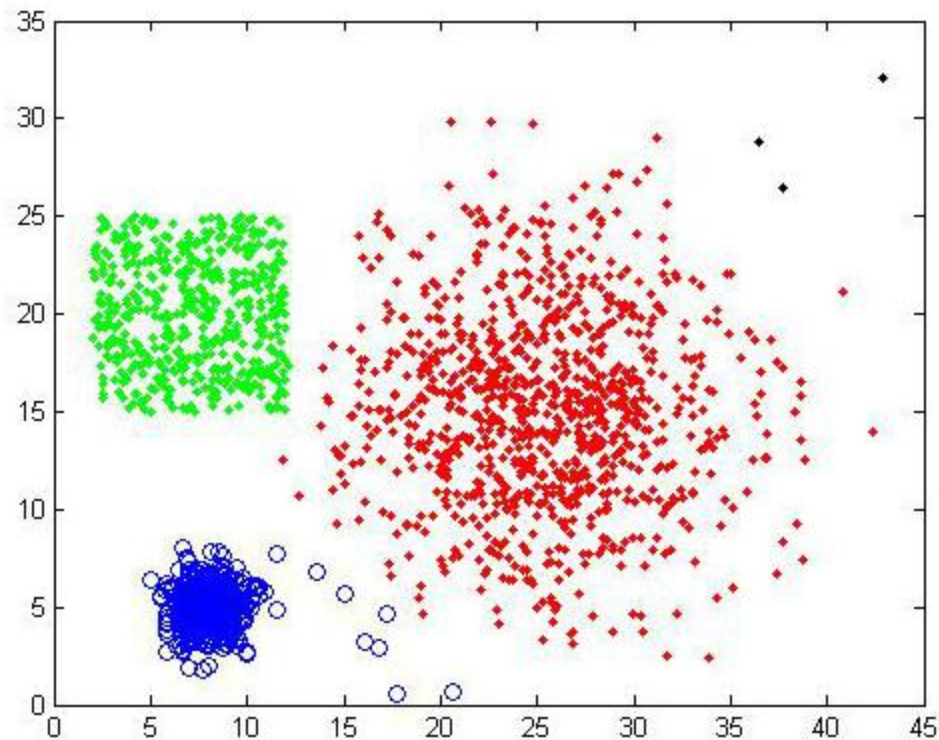
Hierarchical Clustering: Comparison

- A Example: 2000个样本，测试一下算法对cluster size的敏感度



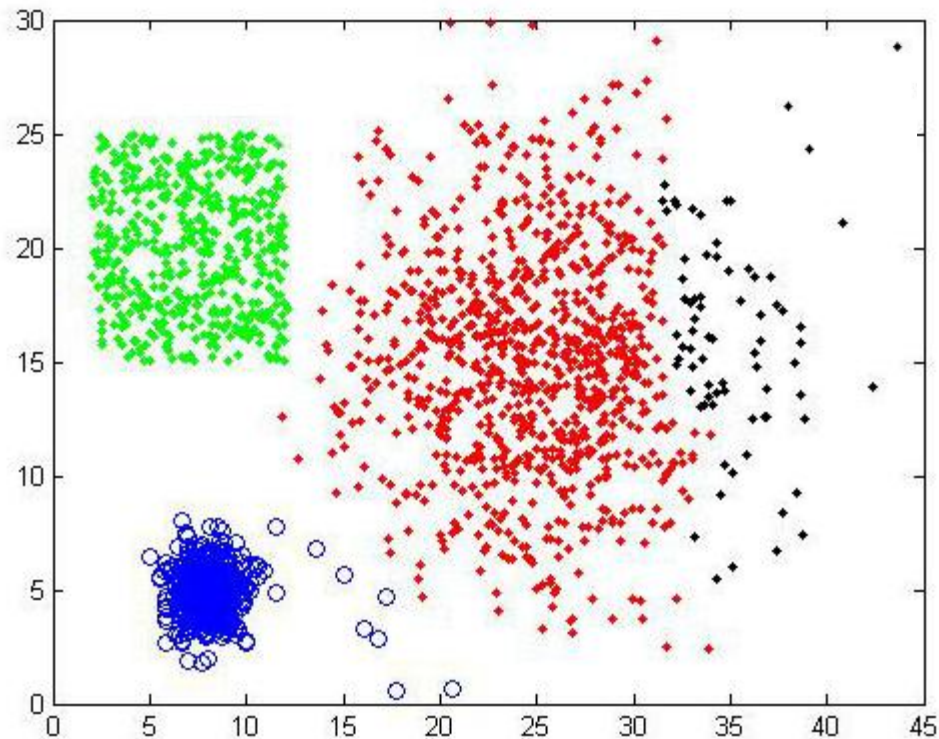
Hierarchical Clustering: Comparison

- A Example: 选择类别数目为4



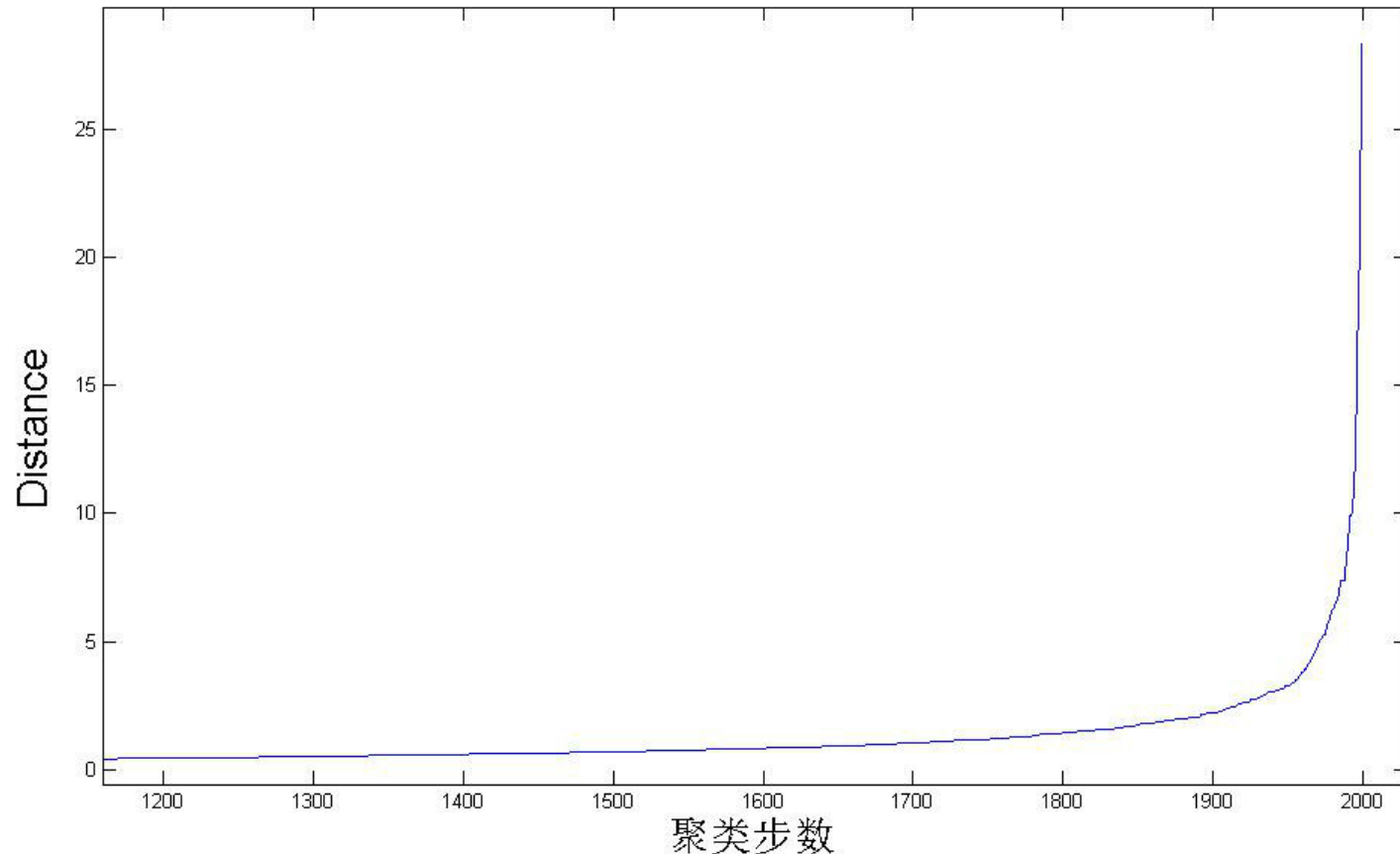
Hierarchical Clustering: Comparison

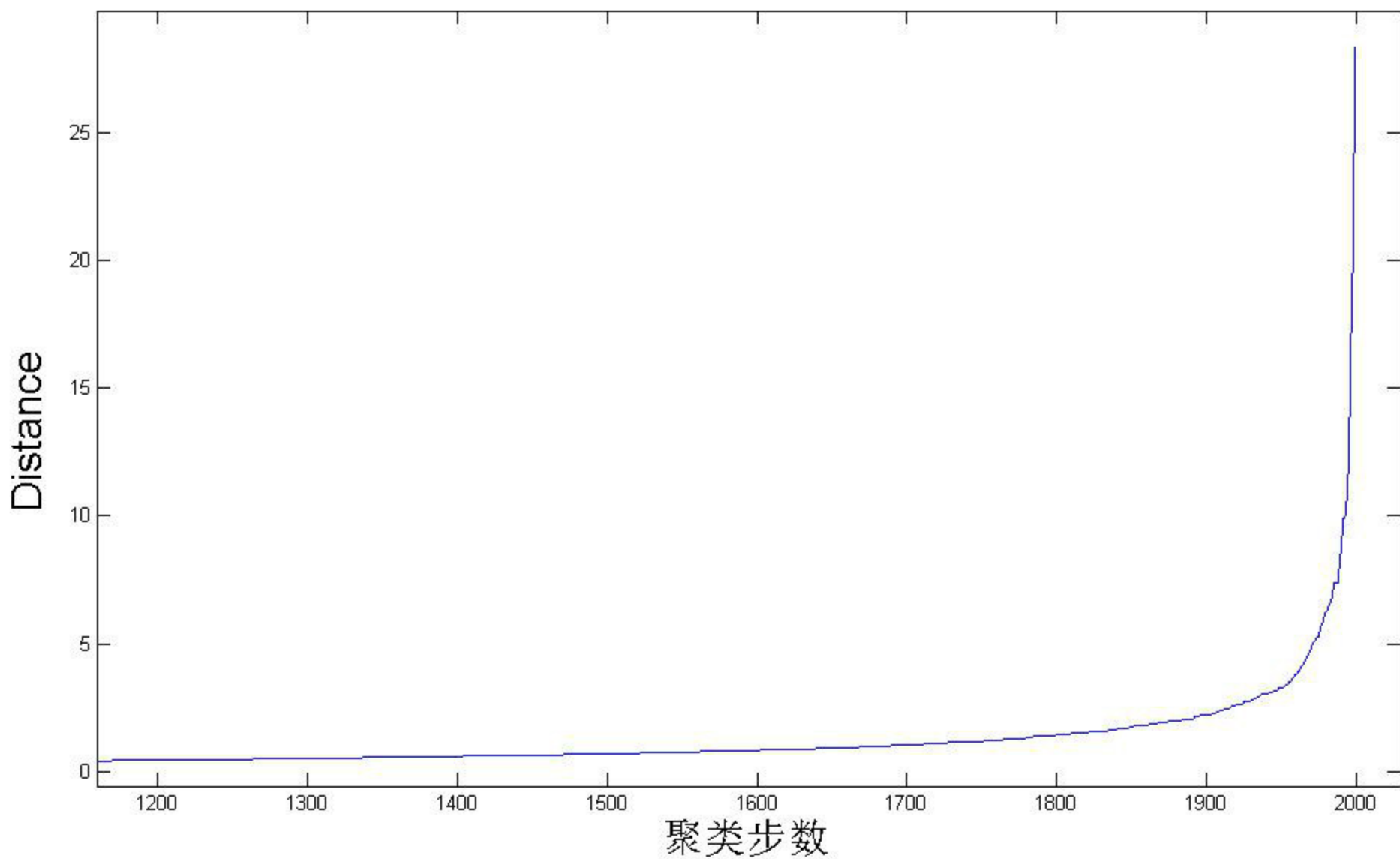
- A Example: 选择类别数目为5



Hierarchical Clustering: Comparison

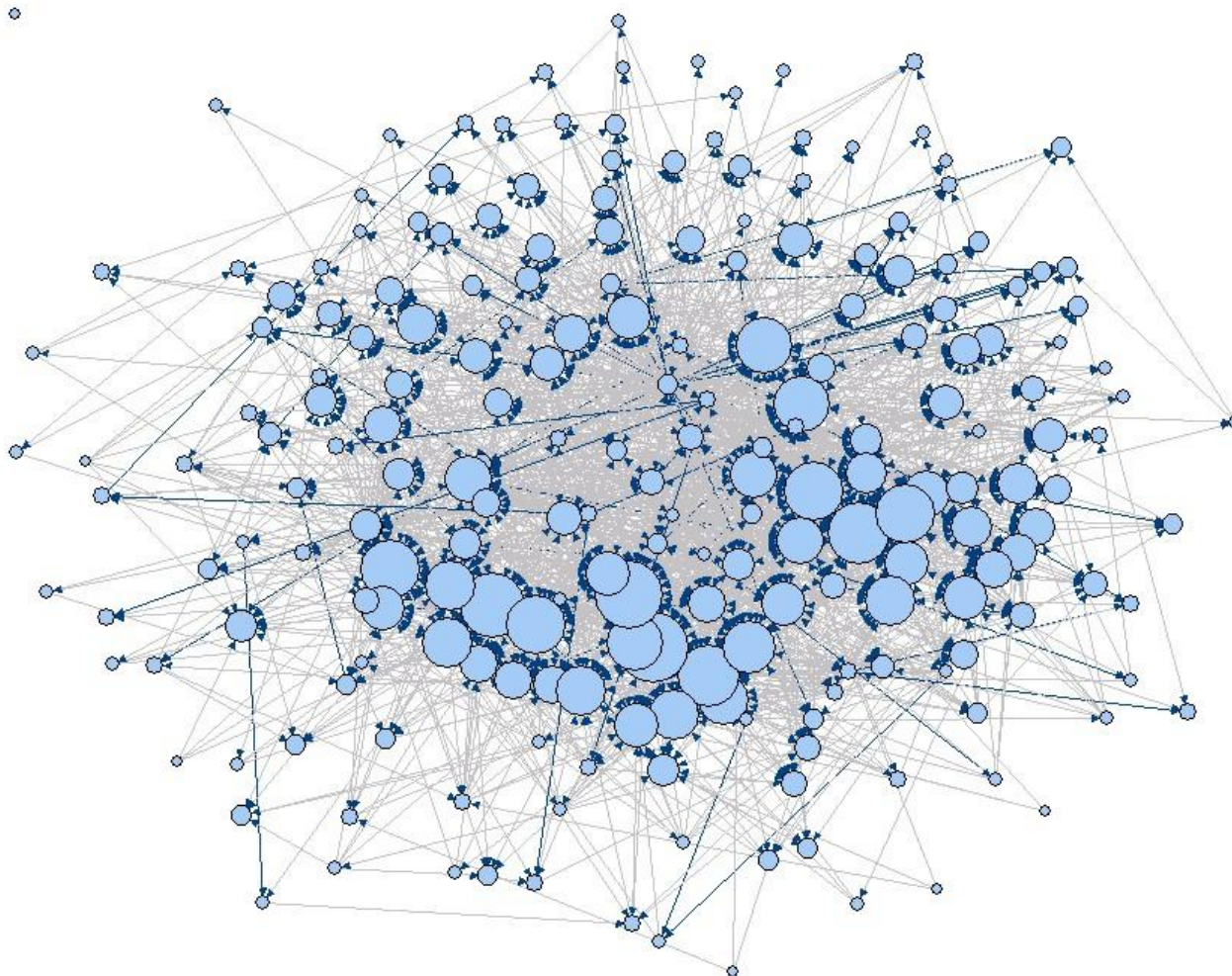
- 如何确定应该取多少个cluster?
 - 2000个样本，假设每次合并两个cluster
 - 每次合并得到的两个cluster的距离





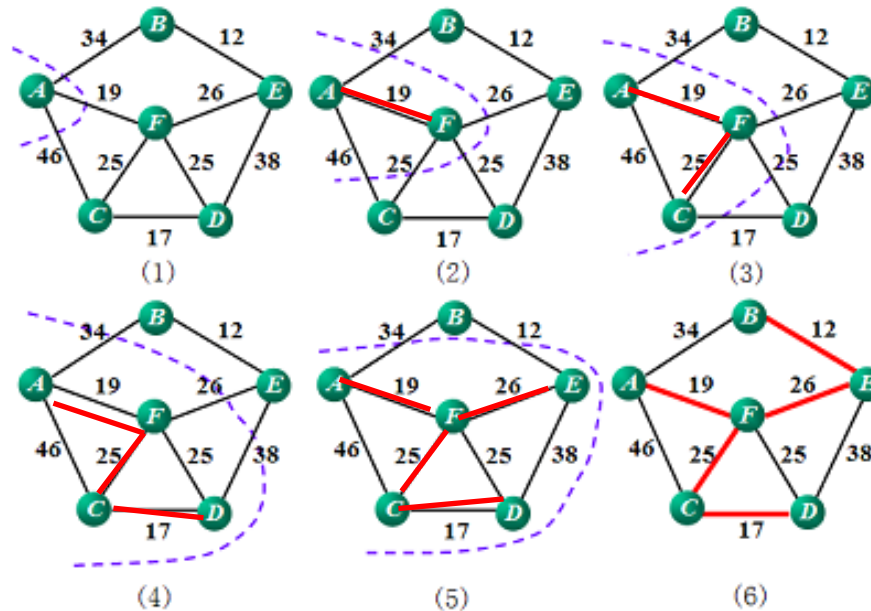
Divisive Hierarchical Clustering

Social Network Graphs



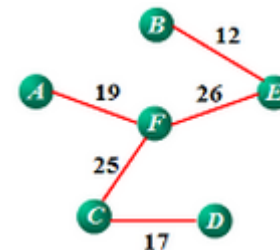
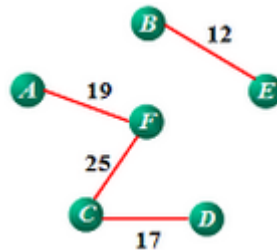
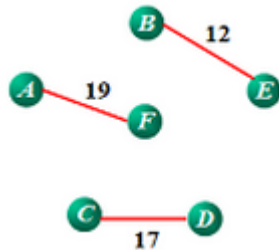
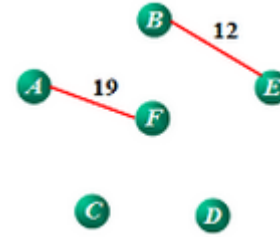
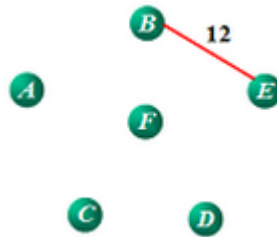
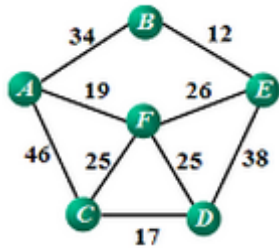
Divisive Hierarchical Clustering

- 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

聚类模型

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

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*

高斯分布

- 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：

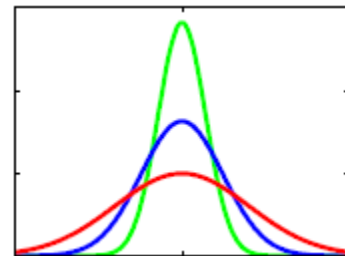
$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2}(x - \mu)^2 \right\}$$

- μ : 模型均值, σ^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, μ : 模型均值, Σ 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：

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

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:

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

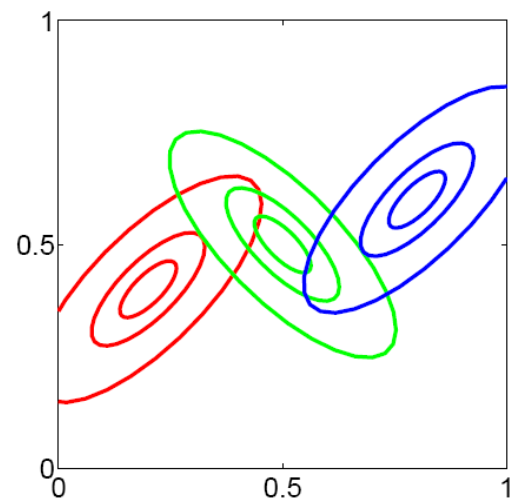
- 混合系数: $\sum_{k=1}^K \pi_k = 1 \quad 0 \leq \pi_k \leq 1$

- 可以看做一种先验概率:

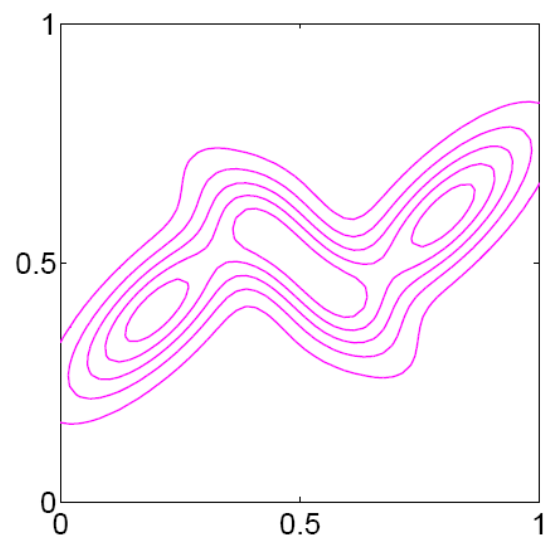
$$p(\mathbf{x}) = \sum_{k=1}^K p(k) p(\mathbf{x} | k)$$

举例

3个Gaussian的混合：

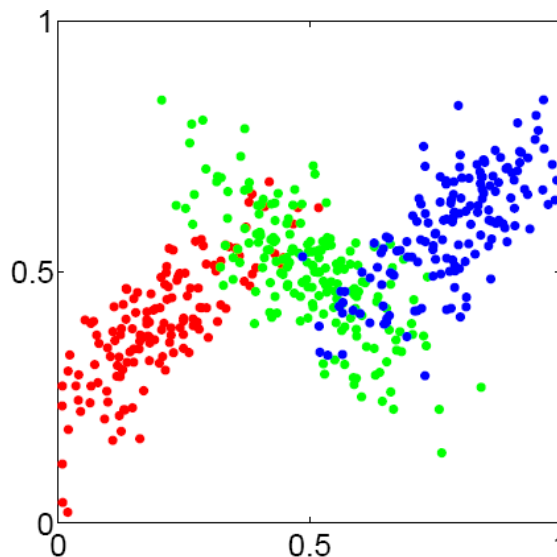


概率分布的等高线：



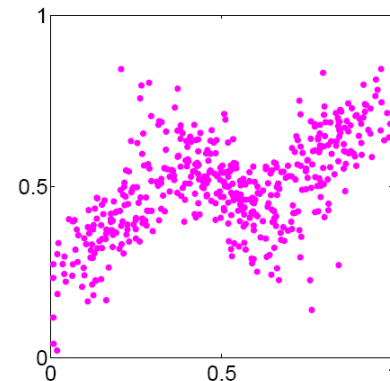
从Gaussian混合分布中抽样得到数据

- 样本点 x_n 的生成过程：
 - 首先，以概率 π_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

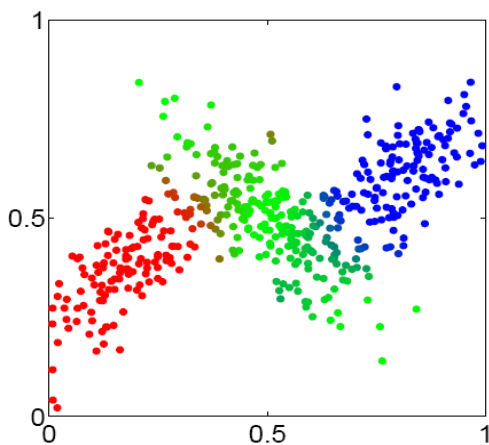


- 如果知道每个样本点由哪个成分抽样得到，则通过极大似然方法可以得到每个类对应的Gaussian模型的参数
- 问题：数据集缺少类别标注
- 因此类别labels可以看做是隐变量 (latent/hidden variable)

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 ，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

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 μ_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

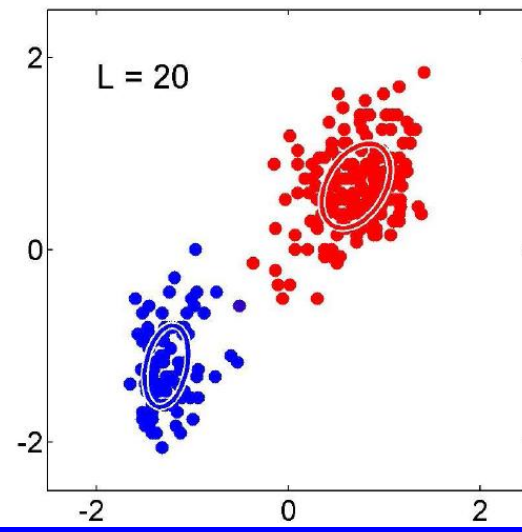
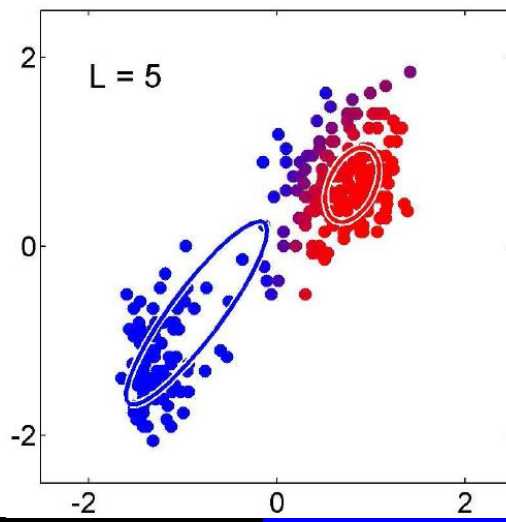
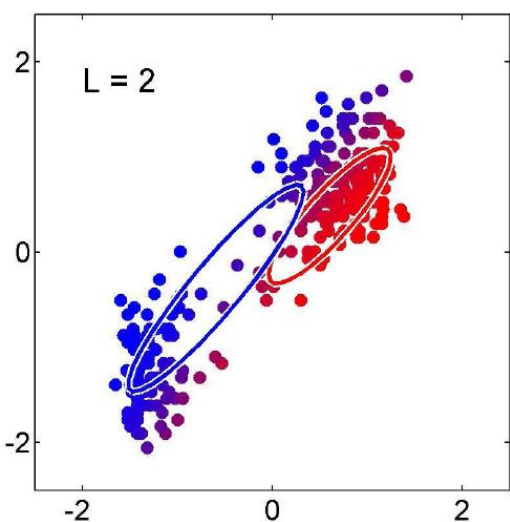
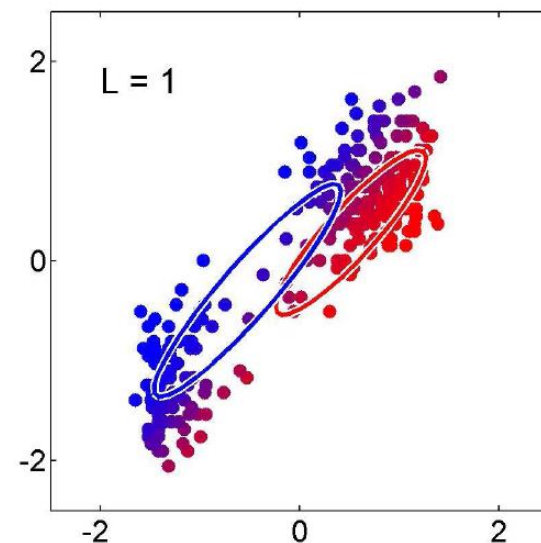
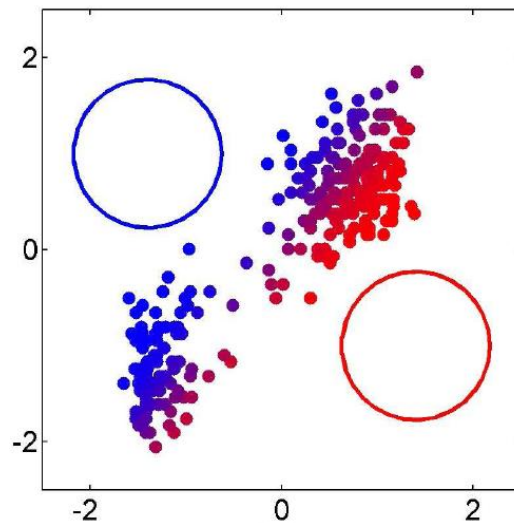
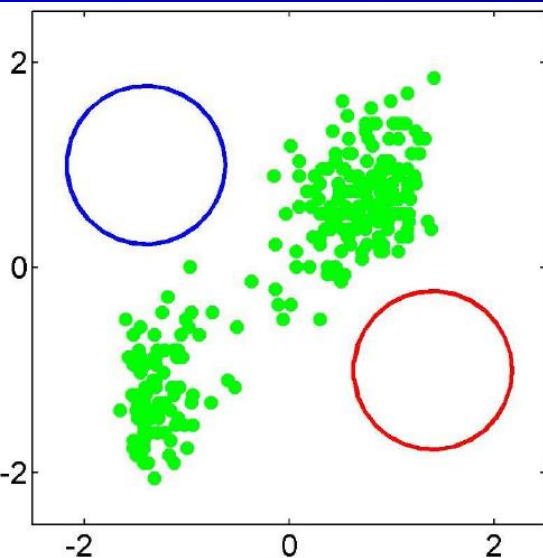
$$\gamma_k(x_i) = \frac{\pi_k N(x_i | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_i | \mu_j, \Sigma_j)}$$

- 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 μ 和 Σ 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 ϵ
- 令极限 $\epsilon \rightarrow 0$
- Responsibilities取两个值:

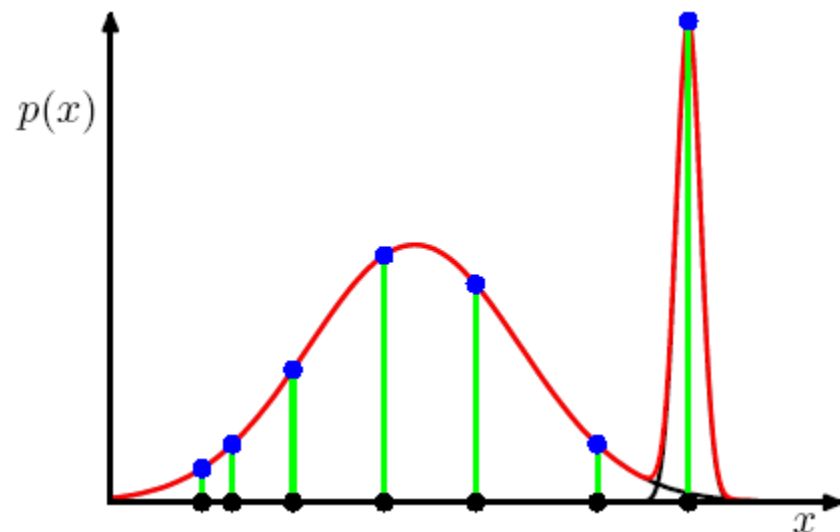
$$\gamma_i(\mathbf{x}_n) = \frac{\pi_i \exp \left\{ -\|\mathbf{x}_n - \boldsymbol{\mu}_i\|^2 / 2\epsilon \right\}}{\sum_j \pi_j \exp \left\{ -\|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2 / 2\epsilon \right\}} \rightarrow r_{ni} \in \{0, 1\}$$

- 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

Problems and Solutions

- 如何避免似然函数中的奇点
 - Bayesian方法
- 如何选取混合成分的个数 K
 - 呃。。。也采用Bayesian方法

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

Next lecture

- Latent Semantic Analysis