

Cluster Analysis

Berlin Chen 2004

References:

1. Foundations of Statistical Natural Language Processing, Chapter 14
2. Data Mining: Concepts, Models, Methods and Algorithms, Chapter 6
3. Modern Information Retrieval, chapters 5, 7

Clustering

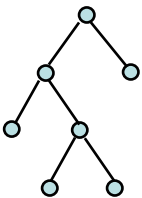
- Place similar objects in the same group and assign dissimilar objects to different groups
 - **Word clustering**
 - Neighbor overlap: words occur with the similar left and right neighbors (such as *in* and *on*)
 - **Document clustering**
 - Documents with the similar topics or concepts are put together
- But clustering cannot give a comprehensive description of the object
 - How to label objects shown on the visual display
- Clustering is a way of learning

Clustering vs. Classification

- Classification is **supervised** and requires a set of labeled training instances for each group (class)
- Clustering is **unsupervised** and learns without a teacher to provide the labeling information of the training data set
 - Also called automatic or unsupervised classification

Types of Clustering Algorithms

- Two types of structures produced by clustering algorithms
 - Flat or non-hierarchical clustering
 - Hierarchical clustering
- **Flat clustering**
 - Simply consisting of a certain number of clusters and the relation between clusters is often undetermined
- **Hierarchical clustering**
 - A hierarchy with usual interpretation that each node stands for a subclass of its mother's node
 - The leaves of the tree are the single objects
 - Each node represents the cluster that contains all the objects of its descendants



Hard Assignment vs. Soft Assignment

- Another important distinction between clustering algorithms is whether they perform soft or hard assignment
- **Hard Assignment**
 - Each object is assigned to one and only one cluster
- **Soft Assignment**
 - Each object may be assigned to multiple clusters
 - An object x_i has a probability distribution $P(\cdot|x_i)$ over clusters c_j where $P(x_i|c_j)$ is the probability that x_i is a member of c_j
 - Is somewhat more appropriate in many tasks such as NLP, IR, ...

Hard Assignment vs. Soft Assignment

- Hierarchical clustering usually adopts hard assignment
- While in flat clustering both types of clustering are common

Summarized Attributes of Clustering Algorithms

- Hierarchical Clustering
 - Preferable for detailed data analysis
 - Provide more information than flat clustering
 - No single best algorithm (each of the algorithms only optimal for some applications)
 - Less efficient than flat clustering (minimally have to compute $n \times n$ matrix of similarity coefficients)

Summarized Attributes of Clustering Algorithms

- Flat Clustering
 - Preferable if efficiency is a consideration or data sets are very large
 - K-means is the conceptually method and should probably be used on a new data because its results are often sufficient
 - K-means assumes a simple Euclidean representation space, and so cannot be used for many data sets, e.g., nominal data like colors
 - The EM algorithm is the most choice. It can accommodate definition of clusters and allocation of objects based on complex probabilistic models

Hierarchical Clustering

Hierarchical Clustering

- Can be in either bottom-up or top-down manners
 - **Bottom-up** (*agglomerative*) 凝集的
 - Start with individual objects and grouping the most similar ones
 - E.g., with the minimum distance apart

$$\text{sim}(x, y) = \frac{1}{1 + d(x, y)}$$

← distance measures will be discussed later on

- The procedure terminates when one cluster containing all objects has been formed
- **Top-down** (*divisive*) 分裂的
 - Start with all objects in a group and divide them into groups so as to maximize *within-group* similarity

Hierarchical Agglomerative Clustering (HAC)

- A bottom-up approach
- Assume a similarity measure for determining the similarity of two objects
- Start with all objects in a separate cluster and then repeatedly joins the two clusters that have the most similarity until there is one only cluster survived
- The history of merging/clustering forms a binary tree or hierarchy

Hierarchical Agglomerative Clustering (HAC)

- Algorithm

```
1 Given: a set  $\mathcal{X} = \{x_1, \dots, x_n\}$  of objects
2       a function  $\text{sim}: \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{X}) \rightarrow \mathbb{R}$ 
3 for  $i := 1$  to  $n$  do   Initialization (for tree leaves):
4    $c_i := \{x_i\}$  end   Each object is a cluster
5  $C := \{c_1, \dots, c_n\}$ 
6  $j := n + 1$ 
7 while  $|C| > 1$    cluster number
8    $(c_{n_1}, c_{n_2}) := \arg \max_{(c_u, c_v) \in C \times C} \text{sim}(c_u, c_v)$ 
9    $c_j = c_{n_1} \cup c_{n_2}$    merged as a new cluster
10   $C := C \setminus \{c_{n_1}, c_{n_2}\} \cup \{c_j\}$    The original two clusters
11   $j := j + 1$    are removed
```

Figure 14.2 Bottom-up hierarchical clustering.

Distance Metrics

- Euclidian Distance (L_2 norm)

$$L_2(\vec{x}, \vec{y}) = \sum_{i=1}^m (x_i - y_i)^2$$

- L_1 Norm

$$L_1(\vec{x}, \vec{y}) = \sum_{i=1}^m |x_i - y_i|$$

- Cosine Similarity (transform to a distance by subtracting from 1)

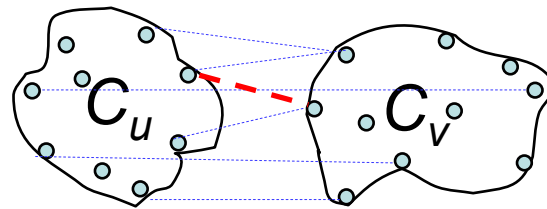
$$1 - \frac{\vec{x} \cdot \vec{y}}{|\vec{x}| \cdot |\vec{y}|}$$

ranged between 0 and 1

Measures of Cluster Similarity

- Especially for the bottom-up approaches
- **Single-link** clustering
 - The similarity between two clusters is the similarity of the two closest objects in the clusters
 - Search over all pairs of objects that are from the two different clusters and select the pair with the greatest similarity

$$\text{sim}(c_i, c_j) = \max_{\vec{x} \in c_i, \vec{y} \in c_j} \text{sim}(\vec{x}, \vec{y})$$

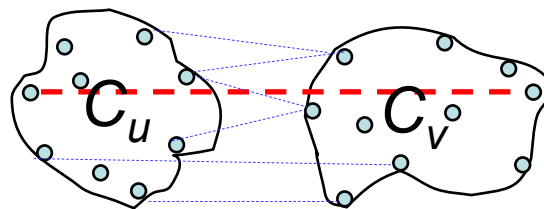


greatest similarity

Measures of Cluster Similarity

- **Complete-link** clustering
 - The similarity between two clusters is the similarity of their two most dissimilar members
 - Sphere-shaped clusters are achieved
 - Preferable for most IR and NLP applications

$$sim(c_i, c_j) = \min_{\vec{x} \in c_i, \vec{y} \in c_j} sim(\vec{x}, \vec{y})$$



least similarity

Measures of Cluster Similarity

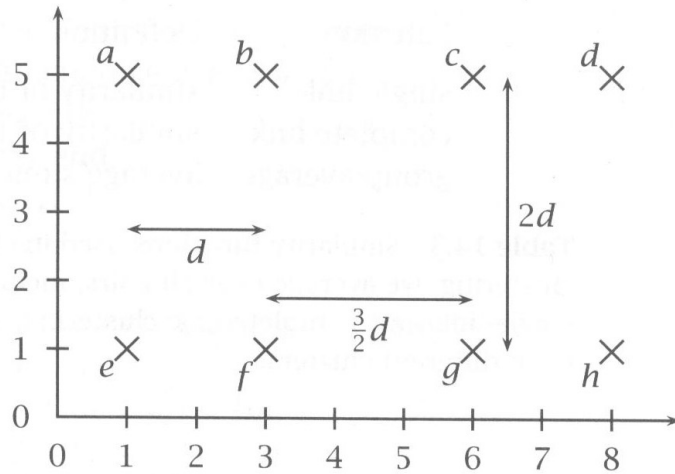


Figure 14.4 A cloud of points in a plane.

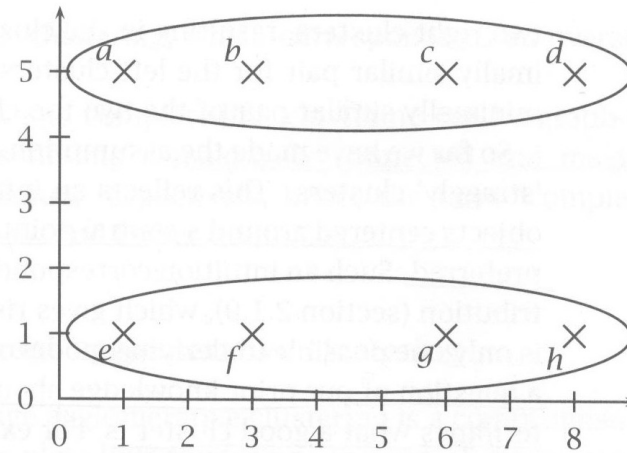


Figure 14.6 Single-link clustering of the points in figure 14.4.

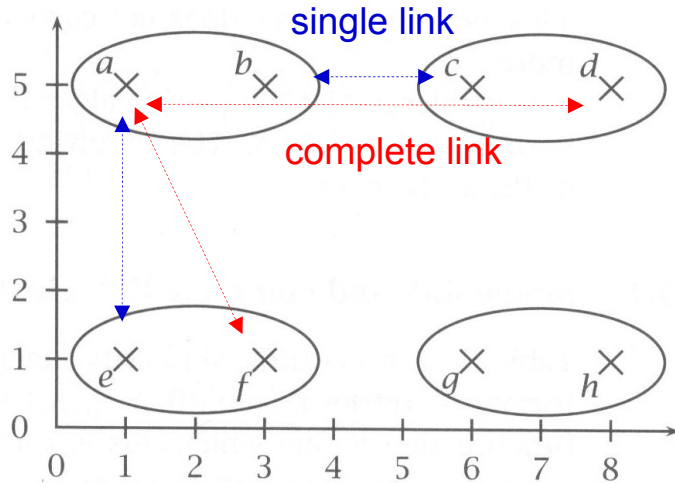


Figure 14.5 Intermediate clustering of the points in figure 14.4.

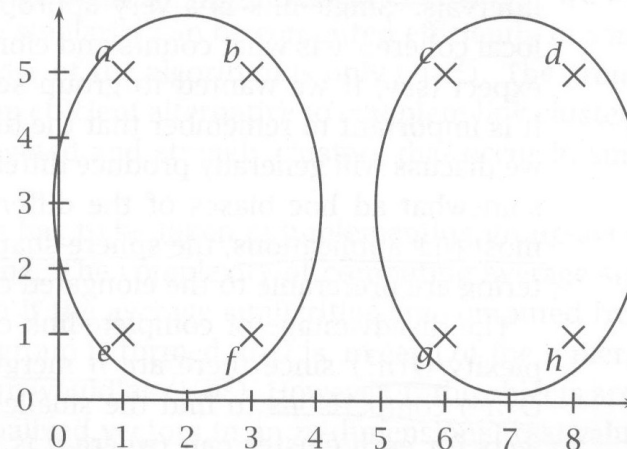


Figure 14.7 Complete-link clustering of the points in figure 14.4.

Measures of Cluster Similarity

- **Group-average** agglomerative clustering
 - A compromise between single-link and complete-link clustering
 - The similarity between two clusters is the average similarity between members
 - If the objects are represented as length-normalized vectors and the similarity measure is the cosine
 - There exists an fast algorithm for computing the average similarity

$$\text{sim}(\vec{x}, \vec{y}) = \cos(\vec{x}, \vec{y}) = \frac{\vec{x} \cdot \vec{y}}{|\vec{x}| |\vec{y}|} = \vec{x} \cdot \vec{y}$$

length-normalized vectors

Measures of Cluster Similarity

- **Group-average** agglomerative clustering (cont.)

- The average similarity *SIM* between vectors in a cluster c_j is defined as

$$SIM(c_j) = \frac{1}{|c_j|(|c_j| - 1)} \sum_{\substack{\vec{x} \in c_j \\ \vec{y} \in c_j \\ \vec{y} \neq \vec{x}}} sim(\vec{x}, \vec{y}) = \frac{1}{|c_j|(|c_j| - 1)} \sum_{\vec{x} \in c_j} \sum_{\substack{\vec{y} \in c_j \\ \vec{y} \neq \vec{x}}} \vec{x} \cdot \vec{y}$$

- The sum of members in a cluster c_j : $\vec{s}(c_j) = \sum_{\vec{x} \in c_j} \vec{x}$

- Express $SIM(c_j)$ in terms of $\vec{s}(c_j)$

$$\vec{s}(c_j) \cdot \vec{s}(c_j) = \sum_{\vec{x} \in c_j} \vec{x} \cdot \vec{s}(c_j) = \sum_{\vec{x} \in c_j} \sum_{\vec{y} \in c_j} \vec{x} \cdot \vec{y}$$

length-normalized vector

$$= |c_j|(|c_j| - 1)SIM(c_j) + \sum_{\vec{x} \in c_j} \vec{x} \cdot \vec{x} = 1$$

$$= |c_j|(|c_j| - 1)SIM(c_j) + |c_j|$$

$$\therefore SIM(c_j) = \frac{\vec{s}(c_j) \cdot \vec{s}(c_j) - |c_j|}{|c_j|(|c_j| - 1)}$$

Measures of Cluster Similarity

- **Group-average** agglomerative clustering (cont.)

-As merging two clusters c_i and c_j , the cluster sum vectors $\vec{s}(c_i)$ and $\vec{s}(c_j)$ are known in advance

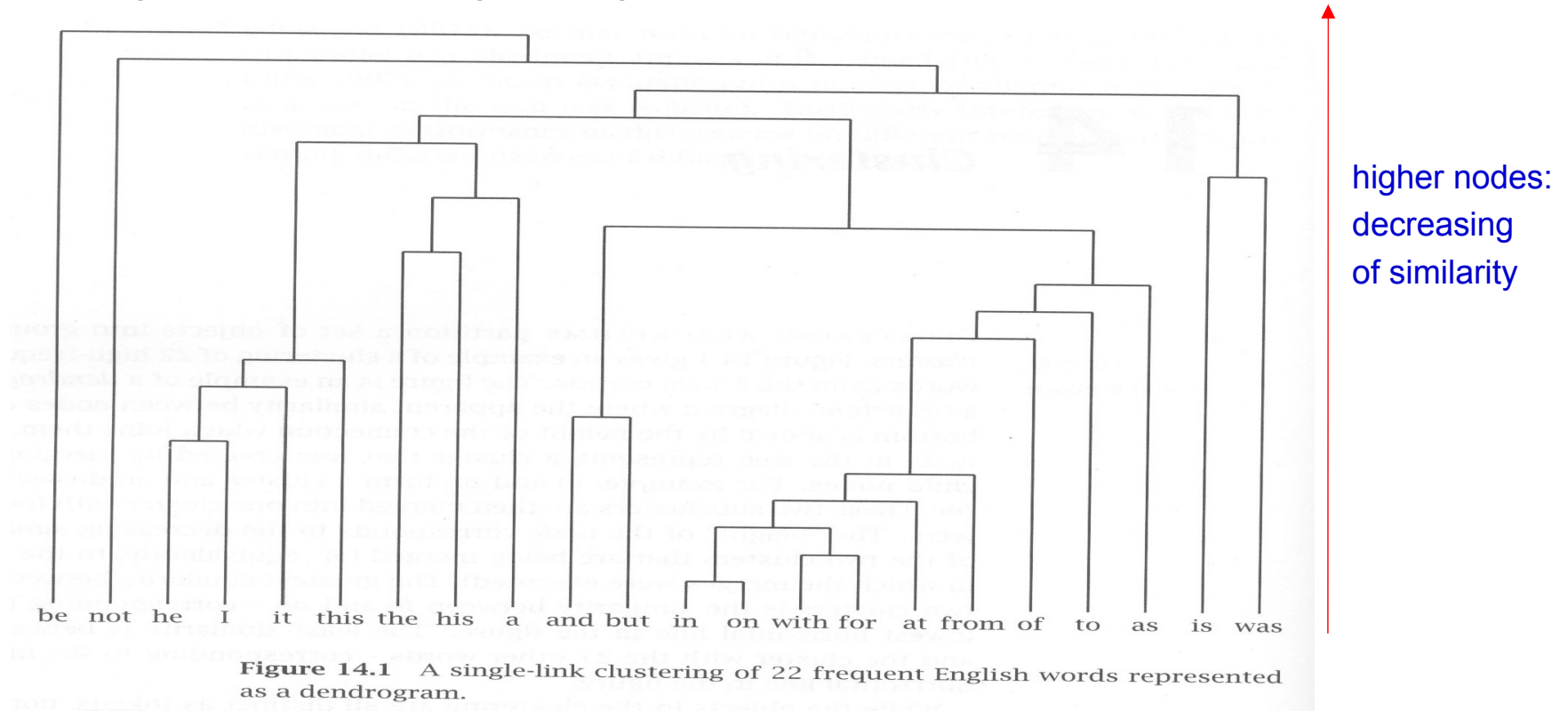
$$\Rightarrow \vec{s}(c_{New}) = \vec{s}(c_i) + \vec{s}(c_j), \quad |c_{New}| = |c_i| + |c_j|$$

– The average similarity for their union will be

$$SIM(c_i \cup c_j) = \frac{(\vec{s}(c_i) + \vec{s}(c_j)) \cdot (\vec{s}(c_i) + \vec{s}(c_j)) - (|c_i| + |c_j|)}{(|c_i| + |c_j|)(|c_i| + |c_j| - 1)}$$

Example: Word Clustering

- Words (objects) are described and clustered using a set of features and values
 - E.g., the left and right neighbors of tokens of words



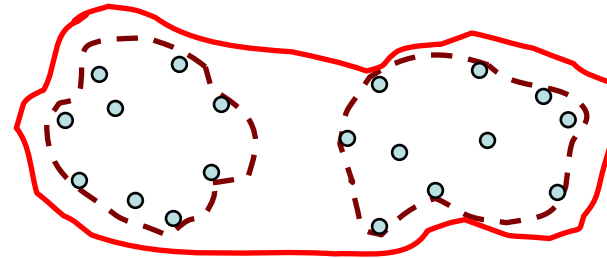
"be" has least similarity with the other 21 words !

Divisive Clustering

- A top-down approach
- Start with all objects in a single cluster
- At each iteration, select the least **coherent** cluster and **split** it
- Continue the iterations until a predefined criterion (e.g., the cluster number) is achieved
- The history of clustering forms a binary tree or hierarchy

Divisive Clustering

- To select the least **coherent** cluster, the measures used in bottom-up clustering can be used again here
 - Single link measure
 - Complete-link measure
 - Group-average measure



- How to **split** a cluster
 - Also is a clustering task (finding two sub-clusters)
 - Any clustering algorithm can be used for the splitting operation, e.g.,
 - Bottom-up (agglomerative) algorithms
 - Non-hierarchical clustering algorithms (e.g., *K*-means)

Divisive Clustering

- Algorithm

```
1 Given: a set  $\mathcal{X} = \{x_1, \dots, x_n\}$  of objects
2       a function  $\text{coh}: \mathcal{P}(\mathcal{X}) \rightarrow \mathbb{R}$ 
3       a function  $\text{split}: \mathcal{P}(\mathcal{X}) \rightarrow \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{X})$ 
4  $C := \{\mathcal{X}\}$  ( $= \{c_1\}$ )
5  $j := 1$ 
6 while  $\exists c_i \in C$  s.t.  $|c_i| > 1$ 
7      $c_u := \arg \min_{c_v \in C} \text{coh}(c_v)$ 
8      $(c_{j+1}, c_{j+2}) := \text{split}(c_u)$ 
9      $C := C \setminus \{c_u\} \cup \{c_{j+1}, c_{j+2}\}$ 
10     $j := j + 2$ 
```

split the least coherent cluster

Generate two new clusters and remove the original one

Figure 14.3 Top-down hierarchical clustering.

Non-Hierarchical Clustering

Non-hierarchical Clustering

- Start out with a partition based on **randomly selected seeds** (one seed per cluster) and then refine the initial partition
 - In a multi-pass manner
- **Problems** associated non-hierarchical clustering
 - When to stop **MI, group average similarity, likelihood**
 - What is the right number of clusters
- **Algorithms** introduced here
 - The *K*-means algorithm
 - The EM algorithm

$k-1 \rightarrow k \rightarrow k+1$

Hierarchical clustering
also has to face this problem

The *K*-means Algorithm

- A **hard clustering** algorithm
- Define clusters by the **center of mass** of their members
- **Initialization**
 - A set of initial cluster centers is needed
- **Recursion**
 - Assign each object to the cluster whose center is closet
 - Then, re-compute the center of each cluster as the **centroid** or mean (average) of its members
 - Using the **medoid** as the cluster center ?
(a medoid is one of the objects in the cluster)



The *K*-means Algorithm

- Algorithm

```
1 Given: a set  $\mathcal{X} = \{\vec{x}_1, \dots, \vec{x}_n\} \subseteq \mathbb{R}^m$ 
2   a distance measure  $d : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ 
3   a function for computing the mean  $\mu : \mathcal{P}(\mathbb{R}) \rightarrow \mathbb{R}^m$ 
4 Select  $k$  initial centers  $\vec{f}_1, \dots, \vec{f}_k$ 
5 while stopping criterion is not true do
6   for all clusters  $c_j$  do
7      $c_j = \{\vec{x}_i \mid \forall \vec{f}_l d(\vec{x}_i, \vec{f}_j) \leq d(\vec{x}_i, \vec{f}_l)\}$ 
8   end
9   for all means  $\vec{f}_j$  do
10     $\vec{f}_j = \mu(c_j)$ 
11  end
12 end
```

cluster centroid

cluster assignment

calculation of new centroids

Figure 14.8 The *K*-means clustering algorithm.

The *K*-means Algorithm

- Example 1

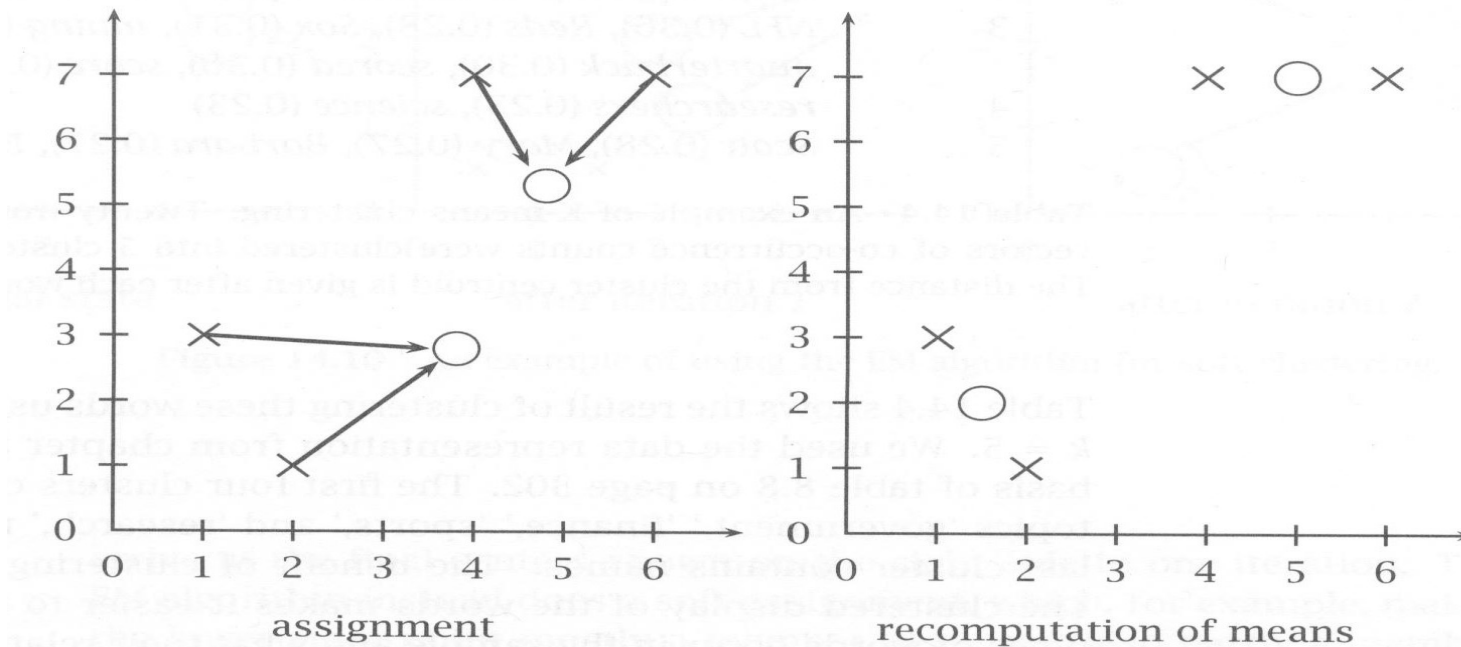


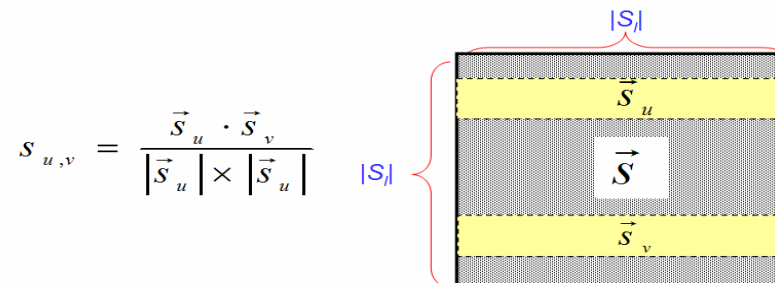
Figure 14.9 One iteration of the K-means algorithm. The first step assigns objects to the closest cluster mean. Cluster means are shown as circles. The second step recomputes cluster means as the center of mass of the set of objects that are members of the cluster.

The *K*-means Algorithm

- Example 2

	Cluster	Members	
$1 - \frac{\vec{x} \cdot \vec{y}}{ \vec{x} \cdot \vec{y} }$	1	<i>ballot</i> (0.28), <i>polls</i> (0.28), <i>Gov</i> (0.30), <i>seats</i> (0.32)	government
	2	<i>profit</i> (0.21), <i>finance</i> (0.21), <i>payments</i> (0.22)	finance
	3	<i>NFL</i> (0.36), <i>Reds</i> (0.28), <i>Sox</i> (0.31), <i>inning</i> (0.33), <i>quarterback</i> (0.30), <i>scored</i> (0.30), <i>score</i> (0.33)	sports
	4	<i>researchers</i> (0.23), <i>science</i> (0.23)	research
	5	<i>Scott</i> (0.28), <i>Mary</i> (0.27), <i>Barbara</i> (0.27), <i>Edward</i> (0.29)	name

Table 14.4 An example of K-means clustering. Twenty words represented as vectors of co-occurrence counts were clustered into 5 clusters using K-means. The distance from the cluster centroid is given after each word.

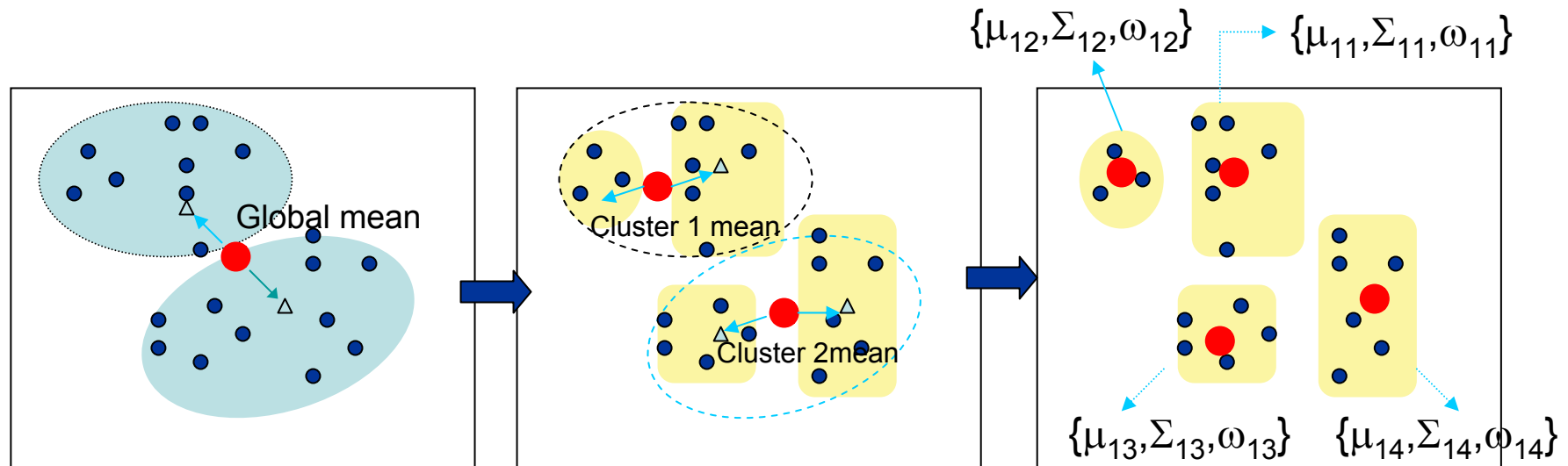


The *K*-means Algorithm

- Choice of initial cluster centers (seeds) is important
 - Pick at random
 - Or use another method such as hierarchical clustering algorithm on a subset of the objects
 - E.g., **buckshot algorithm** uses the group-average agglomerative clustering to randomly sample of the data that has size square root of the complete set
 - Poor seeds will result in **sub-optimal** clustering
- How to break ties in case there are several centers with the same distance from an object
 - Randomly assign the object to one of the candidate clusters
 - Or, perturb objects slightly

The K -means Algorithm

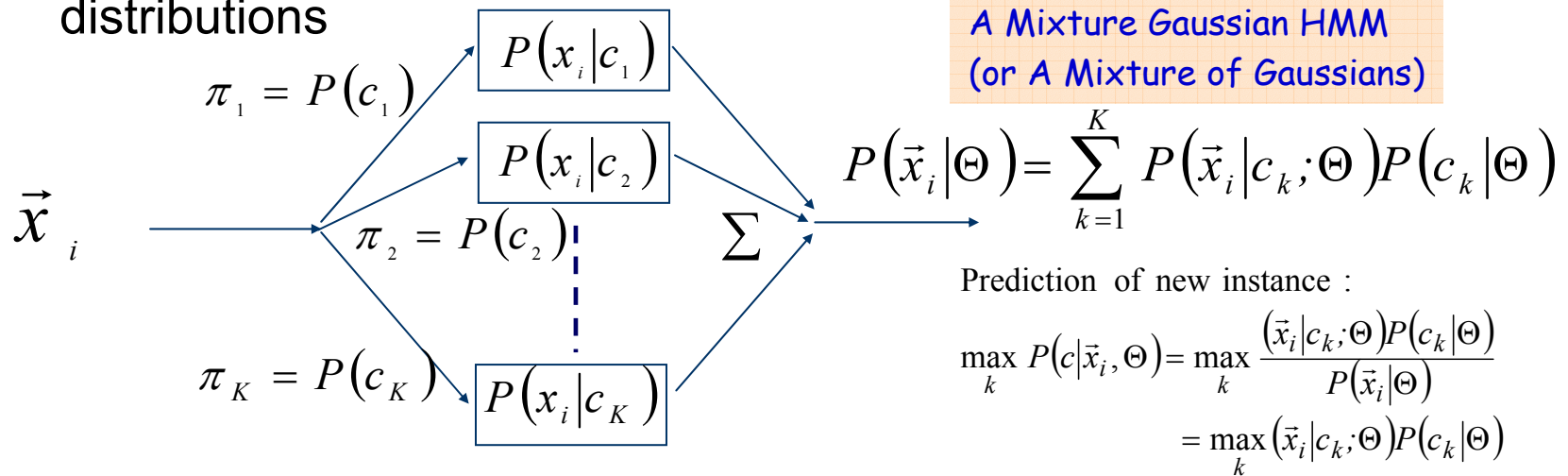
- E.g., the LBG algorithm
 - By Linde, Buzo, and Gray



$M \rightarrow 2M$ at each iteration

The EM Algorithm

- A **soft version** of the K -mean algorithm
 - Each object could be the member of multiple clusters
 - Clustering as estimating a mixture of (continuous) probability distributions



Likelihood function for data samples: $X = \vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$

Continuous case:

$$P(\vec{x}_i | c_k; \Theta) = \frac{1}{\sqrt{(2\pi)^m |\Sigma_k|}} \exp\left(-\frac{1}{2}(\vec{x}_i - \vec{\mu}_k)^T \Sigma_k^{-1} (\vec{x}_i - \vec{\mu}_k)\right)$$

\vec{x}_i 's are independent identically distributed (i.i.d.)

$$P(X | \Theta) = \prod_{i=1}^n P(\vec{x}_i | \Theta)$$

$$= \prod_{i=1}^n \sum_{k_i=1}^K P(\vec{x}_i | c_{k_i}; \Theta) P(c_{k_i} | \Theta)$$

The EM Algorithm

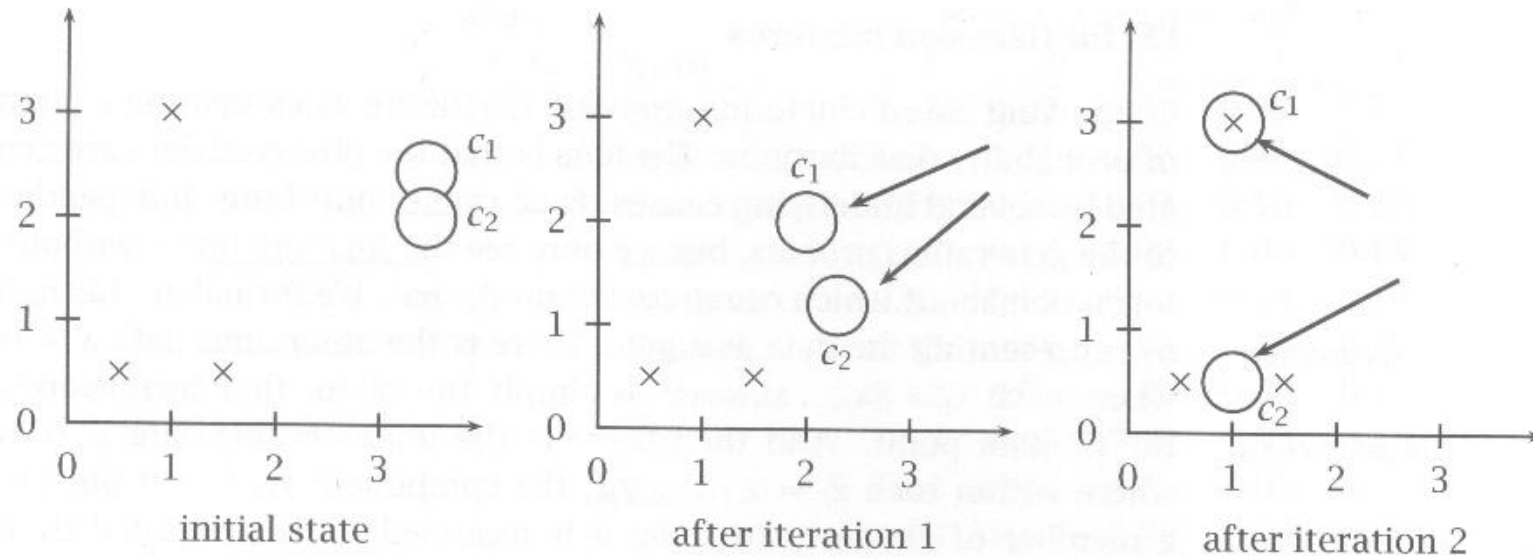


Figure 14.10 An example of using the EM algorithm for soft clustering.

The EM Algorithm

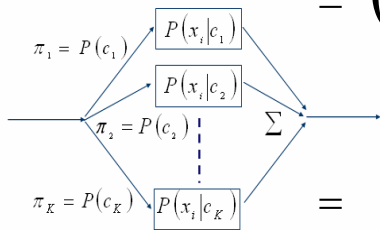
Note :

$$\prod_{i=1}^T \left(\sum_{k_i=1}^K a_{ik_i} \right) \\ = (a_{11} + a_{12} + \dots + a_{1K})(a_{21} + a_{22} + \dots + a_{2K}) \dots (a_{T1} + a_{T2} + \dots + a_{TK}) \\ = \sum_{k_1=1}^K \sum_{k_2=1}^K \dots \sum_{k_T=1}^K \prod_{i=1}^T a_{ik_i}$$

- E-step (Expectation)

- Derive the complete data likelihood function

likelihood function



$$P(\mathbf{X} | \hat{\Theta}) = \prod_{i=1}^n P(\vec{x}_i | \hat{\Theta}) = \prod_{i=1}^n \sum_{k_i=1}^K P(\vec{x}_i | c_{k_i}; \hat{\Theta}) P(c_{k_i} | \hat{\Theta})^{a_{ik_i}} \\ = \left(P(\vec{x}_1 | c_1; \hat{\Theta}) P(c_1 | \hat{\Theta}) + \dots + P(\vec{x}_1 | c_K; \hat{\Theta}) P(c_K | \hat{\Theta}) \right) \times \dots \\ \times \left(P(\vec{x}_n | c_1; \hat{\Theta}) P(c_1 | \hat{\Theta}) + \dots + P(\vec{x}_n | c_K; \hat{\Theta}) P(c_K | \hat{\Theta}) \right)$$

$$= \sum_{k_1=1}^K \sum_{k_2=1}^K \dots \sum_{k_n=1}^K \prod_{i=1}^n \left[P(\vec{x}_i | c_{k_i}; \hat{\Theta}) P(c_{k_i} | \hat{\Theta}) \right]$$

$$= \sum_{k_1=1}^K \sum_{k_2=1}^K \dots \sum_{k_n=1}^K \prod_{i=1}^n \left[P(\vec{x}_i, c_{k_i} | \hat{\Theta}) \right] \quad \text{?}$$

$$= \sum_{k_1=1}^K \sum_{k_2=1}^K \dots \sum_{k_n=1}^K \left[P(\vec{x}_1, c_{k_1}, \vec{x}_2, c_{k_2}, \dots, \vec{x}_n, c_{k_n} | \hat{\Theta}) \right]$$

$$\mathbf{X} = \vec{x}_1 \vec{x}_2 \dots \vec{x}_{n-1} \vec{x}_n \\ \mathbf{C} = c_{k_1} c_{k_2} \dots c_{k_{n-1}} c_{k_n}$$

How many kinds of \mathbf{C} ? (K^n kinds)

$$= \sum_{k_1=1}^K \sum_{k_2=1}^K \dots \sum_{k_n=1}^K \left[P(\vec{x}_1, c_{k_1}, \vec{x}_2, c_{k_2}, \dots, \vec{x}_n, c_{k_n} | \hat{\Theta}) \right]$$

$$= \sum_{\mathbf{C}} \left[P(\mathbf{X}, \mathbf{C} | \hat{\Theta}) \right]$$

the complete data likelihood function

The EM Algorithm

- E-step (Expectation)
 - Define the **auxiliary function** $\Phi(\Theta, \hat{\Theta})$ as the expectation of the **log complete likelihood function** L^{CM} with respect to the hidden/latent variable \mathbf{C} conditioned on known data (\mathbf{X}, Θ)

$$\begin{aligned}\Phi(\Theta, \hat{\Theta}) &= E \left[\log L^{CM} \right]_{\mathbf{C}|\mathbf{X}, \Theta} = E \left[\log P(\mathbf{X}, \mathbf{C} | \hat{\Theta}) \right]_{\mathbf{C}|\mathbf{X}, \Theta} \\ &= \sum_{\mathbf{c}} P(\mathbf{C} | \mathbf{X}, \Theta) \log P(\mathbf{X}, \mathbf{C} | \hat{\Theta}) \\ &= \sum_{\mathbf{c}} \frac{P(\mathbf{X}, \mathbf{C} | \Theta)}{P(\mathbf{X} | \Theta)} \log P(\mathbf{X}, \mathbf{C} | \hat{\Theta})\end{aligned}$$

- Maximize the log likelihood function $\log P(\mathbf{X} | \hat{\Theta})$ by maximizing the expectation of the log complete likelihood function $\Phi(\Theta, \hat{\Theta})$
 - We have shown this when deriving the HMM-based retrieval model

The EM Algorithm

- E-step (Expectation)

- The auxiliary function $\Phi(\Theta, \hat{\Theta})$

data will be aligned to a specific mixture

$$\Phi(\Theta, \hat{\Theta}) = \sum_{\mathbf{C}} \frac{P(\mathbf{X}, \mathbf{C} | \Theta)}{P(\mathbf{X} | \Theta)} \log P(\mathbf{X}, \mathbf{C} | \hat{\Theta})$$

$$= \sum_{\mathbf{C} = c_{k_1} c_{k_2} \dots c_{k_n}} \left[\prod_{j=1}^n \frac{P(\bar{x}_j, c_{k_j} | \Theta)}{P(\bar{x}_j | \Theta)} \right] \left[\log \prod_{i=1}^n P(\bar{x}_i, c_{k_i} | \hat{\Theta}) \right]$$

$$= \sum_{\mathbf{C} = c_{k_1} c_{k_2} \dots c_{k_n}} \left[\prod_{j=1}^n P(c_{k_j} | \bar{x}_j, \Theta) \right] \left[\sum_{i=1}^n \log P(\bar{x}_i, c_{k_i} | \hat{\Theta}) \right]$$

$$\delta_{k, k_i} = \begin{cases} 1 & \text{if } k_i = k \\ 0 & \text{otherwise} \end{cases}$$

$$= \sum_{\mathbf{C} = c_{k_1} c_{k_2} \dots c_{k_n}} \sum_{k=1}^K \left\{ \delta_{k, k_i} \left[\sum_{i=1}^n \log P(\bar{x}_i, c_k | \hat{\Theta}) \right] \left[\prod_{j=1}^n P(c_{k_j} | \bar{x}_j, \Theta) \right] \right\}$$

$$= \sum_{k=1}^K \sum_{i=1}^n \left\{ \left[\log P(\bar{x}_i, c_k | \hat{\Theta}) \right] \sum_{\mathbf{C} = c_{k_1} c_{k_2} \dots c_{k_n}} \delta_{k, k_i} \left[\prod_{j=1}^n P(c_{k_j} | \bar{x}_j, \Theta) \right] \right\}$$

$$= \sum_{k=1}^K \sum_{i=1}^n \left\{ \left[\log P(\bar{x}_i, c_k | \hat{\Theta}) \right] P(c_k | \bar{x}_i, \Theta) \right\}$$

$$= \sum_{k=1}^K \sum_{i=1}^n \left\{ P(c_{k_j} | \bar{x}_j, \Theta) \log P(\bar{x}_i, c_k | \hat{\Theta}) \right\}$$

each mixture collects the corresponding training data

$$= \sum_{k=1}^K \sum_{i=1}^n \left\{ P(c_{k_j} | \bar{x}_j, \Theta) \log \left[P(\bar{x}_i | c_k, \hat{\Theta}) P(c_k | \hat{\Theta}) \right] \right\}$$

$$= \sum_{k=1}^K \sum_{i=1}^n \left\{ P(c_{k_j} | \bar{x}_j, \Theta) \log P(c_k | \hat{\Theta}) \right\} + \sum_{k=1}^K \sum_{i=1}^n \left\{ P(c_{k_j} | \bar{x}_j, \Theta) \log P(\bar{x}_i | c_k, \hat{\Theta}) \right\}$$

The EM Algorithm

– Note that

$$\delta_{k,k_i} = \begin{cases} 1 & \text{if } k_i = k \\ 0 & \text{otherwise} \end{cases}$$

Note:

$$\prod_{t=1}^T \left(\sum_{k_i=1}^K a_{tk_i} \right)$$

$$= (a_{11} + a_{12} + \dots + a_{1K})(a_{21} + a_{22} + \dots + a_{2K}) \dots (a_{T1} + a_{T2} + \dots + a_{TK})$$

$$= \sum_{k_1=1}^K \sum_{k_2=1}^K \dots \sum_{k_T=1}^K \prod_{t=1}^T a_{tk_t}$$

$$\begin{aligned} \mathbf{C} &= c_{k_1} c_{k_2} \dots c_{k_n} \delta_{k,k_i} \left[\prod_{j=1}^n P(c_{k_j} | \vec{x}_j, \Theta) \right] \\ &= \sum_{c_{k_1}=1}^K \sum_{c_{k_2}=1}^K \dots \sum_{c_{k_n}=1}^K \prod_{j=1}^n \left[\delta_{k,k_i} P(c_{k_j} | \vec{x}_j, \Theta) \right] \\ &= \left[\prod_{j=1, j \neq i}^n \left[\sum_{k_j=1}^K P(c_{k_j} | \vec{x}_j, \Theta) \right] \right] \left[\sum_{c_{k_i}=1}^K \delta_{k,k_i} P(c_{k_i} | \vec{x}_i, \Theta) \right] \\ &= \left[\prod_{j=1, j \neq i}^n 1 \right] P(c_k | \vec{x}_i, \Theta) \\ &= P(c_k | \vec{x}_i, \Theta) \end{aligned}$$

\vec{x}_i can only be aligned to c_k

The EM Algorithm

- E-step (Expectation)
 - The auxiliary function can also be divided into two:

$$\Phi(\Theta, \hat{\Theta}) = \Phi_a(\Theta, \hat{\Theta}) + \Phi_b(\Theta, \hat{\Theta})$$

where

$$\begin{aligned} \Phi_a(\Theta, \hat{\Theta}) &= \sum_{i=1}^n \sum_{k=1}^K P(c_k | \vec{x}_i, \Theta) \log P(c_k | \hat{\Theta}) \\ &= \sum_{i=1}^n \sum_{k=1}^K \frac{P(\vec{x}_i | c_k, \Theta) P(c_k | \Theta)}{P(\vec{x}_i | \Theta)} \log P(c_k | \hat{\Theta}) \\ &= \sum_{i=1}^n \sum_{k=1}^K \frac{P(\vec{x}_i | c_k, \Theta) P(c_k | \Theta)}{\sum_{l=1}^K P(\vec{x}_i | c_l, \Theta) P(c_l | \Theta)} \log P(c_k | \hat{\Theta}) \end{aligned}$$

auxiliary function for
mixture weights

$$\begin{aligned} \Phi_b(\Theta, \hat{\Theta}) &= \sum_{i=1}^n \sum_{k=1}^K P(c_k | \vec{x}_i, \Theta) \log P(\vec{x}_i | c_k, \hat{\Theta}) \\ &= \sum_{i=1}^n \sum_{k=1}^K \frac{P(\vec{x}_i | c_k, \Theta) P(c_k | \Theta)}{\sum_{l=1}^K P(\vec{x}_i | c_l, \Theta) P(c_l | \Theta)} \log P(\vec{x}_i | c_k, \hat{\Theta}) \end{aligned}$$

auxiliary function for
cluster distributions

The EM Algorithm

- M-step (Maximization)
 - Remember that
 - Maximize a function F by applying Lagrange multiplier

By applying Lagrange Multiplier ℓ

Suppose that $F = \sum_{j=1}^N w_j \log y_j \Rightarrow \hat{F} = \sum_{j=1}^N w_j \log y_j + \ell \left(\sum_{j=1}^N y_j - 1 \right)$

$$\frac{\partial \hat{F}}{\partial y_j} = \frac{w_j}{y_j} + \ell = 0 \Rightarrow \ell = -\frac{w_j}{y_j} \quad \forall j$$

$$\ell \sum_{j=1}^N y_j = -\sum_{j=1}^N w_j \Rightarrow \ell = -\sum_{j=1}^N w_j$$

$$\therefore y_j = \frac{w_j}{\sum_{j=1}^N w_j}$$

Constraint

Note :

$$\frac{\partial \log y_j}{\partial y_j} = \frac{1}{y_j}$$

The EM Algorithm

- M-step (Maximization)

auxiliary function for mixture weights (or priors for Gaussians)

- Maximize $\Phi_a(\Theta, \hat{\Theta})$

$$\begin{aligned} \bar{\Phi}_a(\Theta, \hat{\Theta}) &= \Phi_a(\Theta, \hat{\Theta}) + l \left(\sum_{k=1}^K P(c_k | \hat{\Theta}) - 1 \right) \\ &= \sum_{k=1}^K \underbrace{\sum_{i=1}^n \frac{P(\vec{x}_i | c_k, \Theta) P(c_k | \Theta)}{\sum_{l=1}^K P(\vec{x}_i | c_l, \Theta) P(c_l | \Theta)}}_{w_k} \log \underbrace{P(c_k | \hat{\Theta})}_{y_k} + l \left(\sum_{k=1}^K P(c_k | \hat{\Theta}) - 1 \right) \end{aligned}$$

$$\Rightarrow \hat{\pi}_k = P(c_k | \hat{\Theta}) = \frac{w_k}{\sum_{k=1}^K w_k} = \frac{\sum_{i=1}^n \frac{P(\vec{x}_i | c_k, \Theta) P(c_k | \Theta)}{\sum_{l=1}^K P(\vec{x}_i | c_l, \Theta) P(c_l | \Theta)}}{\sum_{k=1}^K \sum_{i=1}^n \frac{P(\vec{x}_i | c_k, \Theta) P(c_k | \Theta)}{\sum_{l=1}^K P(\vec{x}_i | c_l, \Theta) P(c_l | \Theta)}} = \frac{\sum_{i=1}^n \frac{P(\vec{x}_i | c_k, \Theta) P(c_k | \Theta)}{\sum_{l=1}^K P(\vec{x}_i | c_l, \Theta) P(c_l | \Theta)}}{n}$$

The EM Algorithm

- M-step (Maximization)

- Maximize $\Phi_b(\Theta, \hat{\Theta})$

auxiliary function for
Gaussian Means and Variances

$$P(\vec{x}_i | c_k; \Theta) = \frac{1}{\sqrt{(2\pi)^m |\Sigma_k|}} \exp\left(-\frac{1}{2}(\vec{x}_i - \vec{\mu}_k)^T \Sigma_k^{-1} (\vec{x}_i - \vec{\mu}_k)\right)$$

$$\Phi_b(\Theta, \hat{\Theta}) = \sum_{i=1}^n \sum_{k=1}^K \frac{P(\vec{x}_i | c_k, \Theta) P(c_k | \Theta)}{\sum_{l=1}^K P(\vec{x}_i | c_l, \Theta) P(c_l | \Theta)} \log P(\vec{x}_i | c_k, \hat{\Theta})$$

Let $w_{k,i} = \frac{P(\vec{x}_i | c_k, \Theta) P(c_k | \Theta)}{\sum_{l=1}^K P(\vec{x}_i | c_l, \Theta) P(c_l | \Theta)}$ and $\log P(\vec{x}_i | c_k; \Theta) =$

$$-m/2 \cdot \log(2\pi) - 1/2 \log |\Sigma_k| - \frac{1}{2} (\vec{x}_i - \vec{\mu}_k)^T \Sigma_k^{-1} (\vec{x}_i - \vec{\mu}_k)$$

$$\Rightarrow \Phi_b(\Theta, \hat{\Theta}) = - \sum_{i=1}^n \sum_{k=1}^K w_{k,i} \left[\frac{1}{2} \log |\Sigma_k| + \frac{1}{2} (\vec{x}_i - \hat{\mu}_k)^T \hat{\Sigma}_k^{-1} (\vec{x}_i - \hat{\mu}_k) \right] + D$$

constant

The EM Algorithm

- M-step (Maximization)

- Maximize $\Phi_b(\Theta, \hat{\Theta})$ with respect to $\vec{\mu}_k$

$$\Phi_b(\Theta, \hat{\Theta}) = - \sum_{i=1}^n \sum_{k=1}^K w_{k,i} \left[\frac{1}{2} \log |\hat{\Sigma}_k| + \frac{1}{2} (\vec{x}_i - \hat{\mu}_k)^T \hat{\Sigma}_k^{-1} (\vec{x}_i - \hat{\mu}_k) \right] + D$$

$$\frac{\partial \Phi_b(\Theta, \hat{\Theta})}{\partial \hat{\mu}_k} = - \sum_{i=1}^n w_{k,i} \cdot \frac{1}{2} \cdot (2) \cdot \hat{\Sigma}_k^{-1} (\vec{x}_i - \hat{\mu}_k) (-1) = 0$$

$$\Rightarrow \hat{\mu}_k = \frac{\sum_{i=1}^n w_{k,i} \cdot \vec{x}_i}{\sum_{i=1}^n w_{k,i}} = \frac{\sum_{i=1}^n \frac{P(\vec{x}_i | c_k, \Theta) P(c_k | \Theta) \cdot \vec{x}_i}{\sum_{l=1}^K P(\vec{x}_i | c_l, \Theta) P(c_l | \Theta)}}{\sum_{i=1}^n \frac{P(\vec{x}_i | c_k, \Theta) P(c_k | \Theta)}{\sum_{l=1}^K P(\vec{x}_i | c_l, \Theta) P(c_l | \Theta)}}$$

$$\frac{d(\mathbf{x}^T \mathbf{C} \mathbf{x})}{d\mathbf{x}} = (\mathbf{C} + \mathbf{C}^T) \mathbf{x}$$

and Σ_k^{-1} is symmetric here

The EM Algorithm

- M-step (Maximization)

- Maximize $\Phi_b(\Theta, \hat{\Theta})$ with respect to Σ_k

$$\Phi_b(\Theta, \hat{\Theta}) = - \sum_{i=1}^n \sum_{k=1}^K w_{k,i} \left[\frac{1}{2} \log |\hat{\Sigma}_k| + \frac{1}{2} (\bar{x}_i - \hat{\mu}_k)^T \hat{\Sigma}_k^{-1} (\bar{x}_i - \hat{\mu}_k) \right] + D$$

$$\frac{\partial \Phi_b(\Theta, \hat{\Theta})}{\partial \hat{\Sigma}_k} = - \sum_{i=1}^n \frac{1}{2} \cdot w_{k,i} \left[\cancel{|\hat{\Sigma}_k|^{-1}} \cdot \cancel{|\hat{\Sigma}_k|} \cdot \hat{\Sigma}_k^{-1} - \Sigma_k^{-1} (\bar{x}_i - \hat{\mu}_k) (\bar{x}_i - \hat{\mu}_k)^T \Sigma_k^{-1} \right] = 0$$

$$\Rightarrow \sum_{i=1}^n w_{k,i} \cdot \hat{\Sigma}_k^{-1} = \sum_{i=1}^n w_{k,i} \cdot \hat{\Sigma}_k^{-1} (\bar{x}_i - \hat{\mu}_k) (\bar{x}_i - \hat{\mu}_k)^T \hat{\Sigma}_k^{-1}$$

$$\Rightarrow \sum_{i=1}^n w_{k,i} \cdot \hat{\Sigma}_k \hat{\Sigma}_k^{-1} \hat{\Sigma}_k = \sum_{i=1}^n w_{k,i} \cdot \hat{\Sigma}_k \hat{\Sigma}_k^{-1} (\bar{x}_i - \hat{\mu}_k) (\bar{x}_i - \hat{\mu}_k)^T \hat{\Sigma}_k^{-1} \hat{\Sigma}_k$$

$$\Rightarrow \sum_{i=1}^n w_{k,i} \cdot \hat{\Sigma}_k = \sum_{i=1}^n w_{k,i} \cdot (\bar{x}_i - \hat{\mu}_k) (\bar{x}_i - \hat{\mu}_k)^T$$

$$\Rightarrow \hat{\Sigma}_k = \frac{\sum_{i=1}^n w_{k,i} \cdot (\bar{x}_i - \hat{\mu}_k) (\bar{x}_i - \hat{\mu}_k)^T}{\sum_{i=1}^n w_{k,i}} = \frac{\sum_{i=1}^n \frac{P(\bar{x}_i | c_k, \Theta) P(c_k | \Theta)}{\sum_{l=1}^K P(\bar{x}_i | c_l, \Theta) P(c_l | \Theta)} \cdot (\bar{x}_i - \hat{\mu}_k) (\bar{x}_i - \hat{\mu}_k)^T}{\sum_{i=1}^n \frac{P(\bar{x}_i | c_k, \Theta) P(c_k | \Theta)}{\sum_{l=1}^K P(\bar{x}_i | c_l, \Theta) P(c_l | \Theta)}}$$

$$\frac{d[\det(\mathbf{X})]}{d\mathbf{X}} = \det(\mathbf{X}) \cdot \mathbf{X}^{-T}$$

and

$$\frac{d(\mathbf{a}^T \mathbf{X}^{-1} \mathbf{b})}{d\mathbf{X}} = -\mathbf{X}^{-1} \mathbf{a} \mathbf{b}^T \mathbf{X}^{-1}$$

The EM Algorithm

- The initial cluster distributions can be estimated using the K -means algorithm
- The procedure terminates when the likelihood function $P(X | \Theta)$ is converged or maximum number of iterations is reached