

# Lecture 21 - Mar 20

Gaussian Mixture Clustering (by EM algorithm)  
Hierarchical Clustering

## References

*Data Mining and Machine Learning*

Ch 13 - Representation-based Clustering

Ch 14 - Hierarchical Clustering

*Elements of Statistical Learning*

14.3.7 Gaussian Mixtures as Soft K-Means Clustering

14.3.12 Hierarchical Clustering

Generalize k-means to soft assignment where each point has a probability of being in each cluster.

Assume each cluster  $C_i$  is characterized by a normal distribution

Assume pdf of  $x \in X$  is a Gaussian mixture over the clusters normal densities

$$f(x) = \sum_{i=1}^k f_i(x) P(C_i) = \sum_{i=1}^k f_i(x | \mu_i, \Sigma_i) P(C_i)$$

↑  
priors = "mixture parameters"

### Gaussian mixture model

Assume each cluster  $C_i$  is characterized by a univariate normal distribution

Multivariate normal distribution

$$f_i(x) = f(x | \mu_i, \Sigma_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left\{ -\frac{(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)}{2} \right\}$$

↑  
unknown parameters

Assume pdf of  $x \in X$  is a Gaussian mixture over the cluster normal densities

$$f(x) = \sum_{i=1}^k f_i(x) P(C_i) = \sum_{i=1}^k f_i(x | \mu_i, \Sigma_i) P(C_i)$$

↑  
priors = "mixture parameters"

Model parameters  $\Theta = \{ \mu_i, \Sigma_i, P(C_i) \mid i=1, 2, \dots, k \}$

## Gaussian Mixture Clustering in 1D

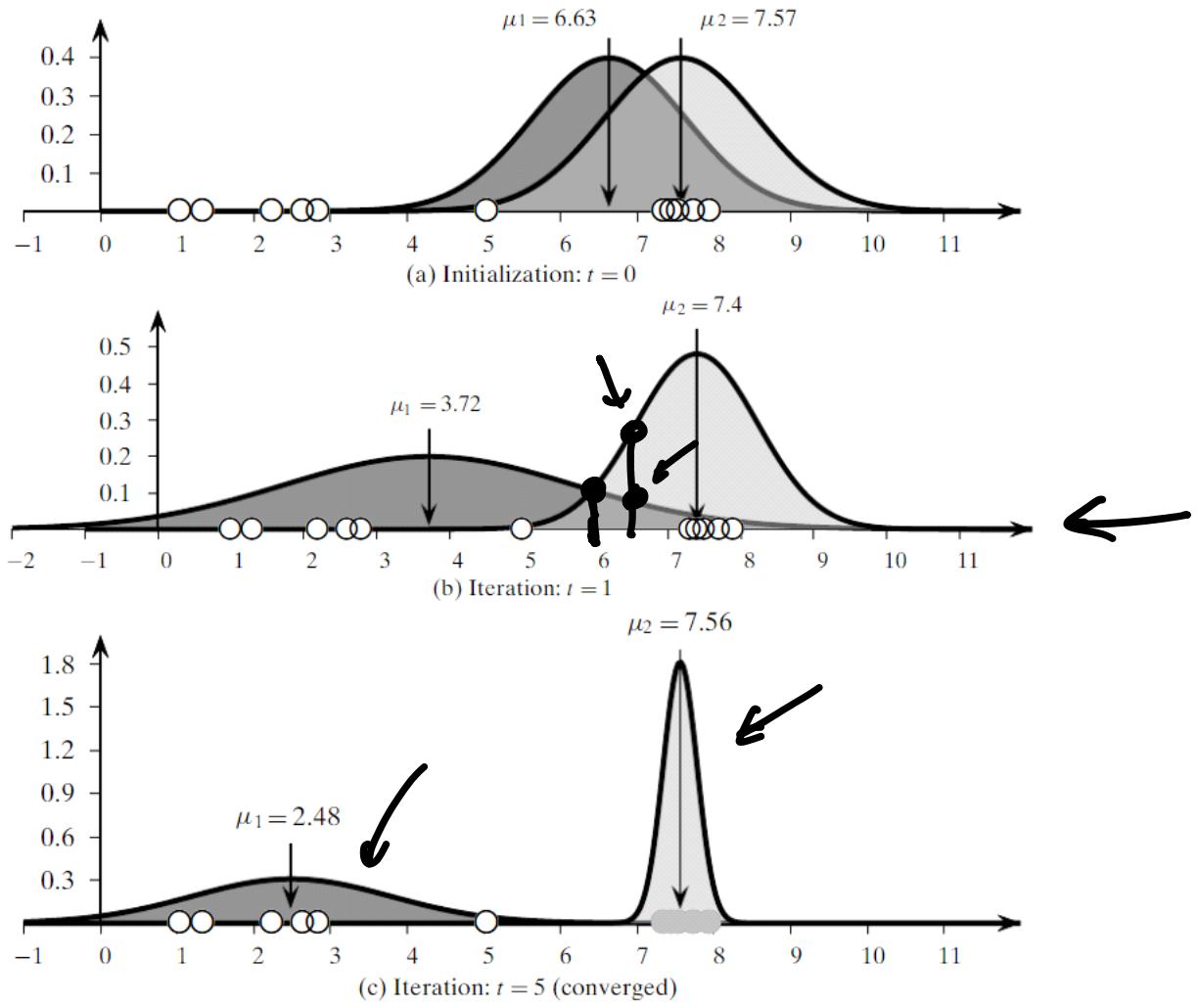
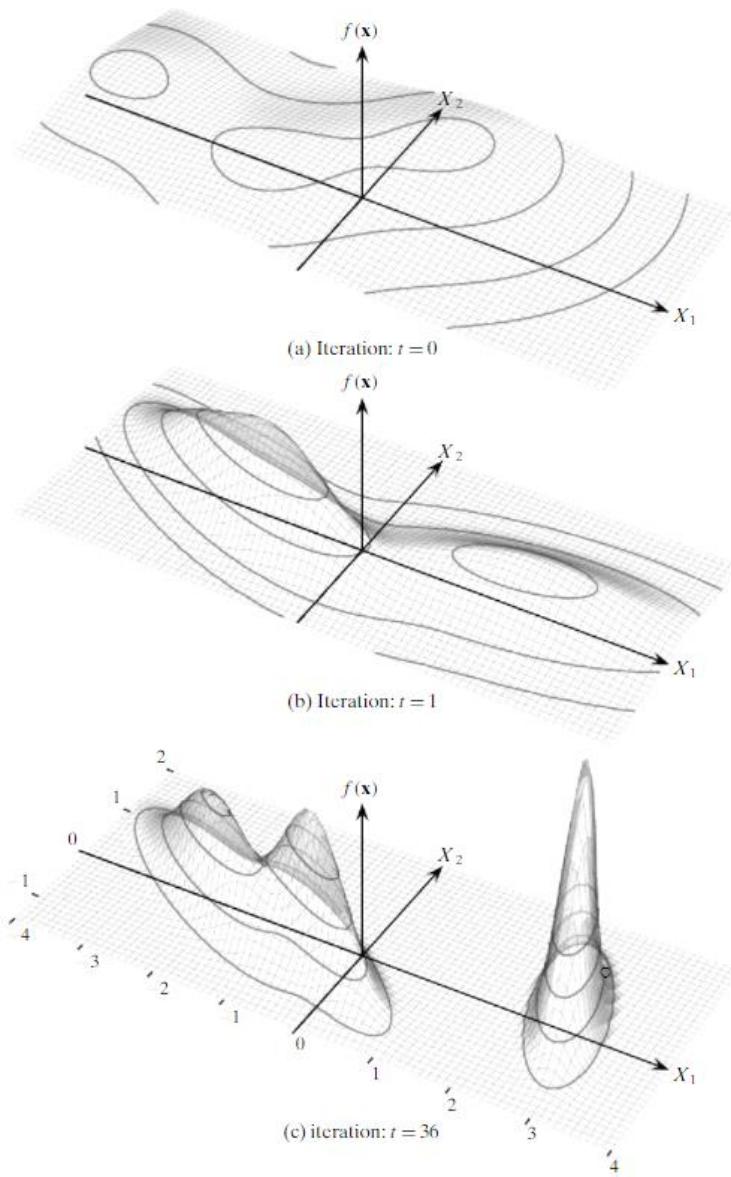


Figure 13.4. EM in one dimension.

## Gaussian Mixture Clustering in 2D



Ch 12 Lab  
~~Ch 12 Lab~~

Figure 13.5. EM algorithm in two dimensions: mixture of  $k = 3$  Gaussians.

## Maximum Likelihood Estimation of $\Theta$

Assuming  $x_1, \dots, x_n$  form a sample, likelihood of  $\Theta$  is

$$P(X|\theta) = \prod_{j=1}^n f(x_j)$$

$$\text{Goal: } \theta^* = \underset{\theta}{\operatorname{argmax}} \left\{ P(X|\theta) \right\} = \underset{\theta}{\operatorname{argmax}} \left\{ \ln(P(X|\theta)) \right\}$$

$$\begin{aligned} \ln(P(X|\theta)) &= \sum_{j=1}^n \ln(f(x_j)) \\ &= \sum_{j=1}^n \ln \left( \sum_{i=1}^k f(x_j | \mu_i, \Sigma_i) P(C_i) \right) \end{aligned}$$

Maximizing by hand is typically infeasible, so EM algorithm can iteratively seek  $\theta$  with two steps:

0. Make an initial estimate for  $\theta$

1. Given estimated  $\theta$ , compute posterior probabilities with

1. Given estimated  $\Theta$ , compute  $P(c_i | x_j)$   
 Bayes' Theorem

"Expectation  
step"

$$P(c_i | x_j) = \frac{P(x_j | c_i) P(c_i)}{\sum_{a=1}^k P(x_j | c_a) P(c_a)}$$

using  $P(x_j | c_i) \approx Z \sum f_i(x)$

$$\Rightarrow P(c_i | x_j) \approx \frac{f_i(x_j) P(c_i)}{\sum_{a=1}^k f_a(x_j) P(c_a)}$$

2. Given weights  $P(c_i | x_j)$ , estimate  $M_i, \Sigma_i, P(c_i)$  for each  $C_i$ .

"maximization  
step"

Initialization

initialize  $M_i$  uniformly randomly in feature space

$$\Sigma_i = I_d$$

$$P(c_i) = \frac{1}{k}$$

Expectation Step

## Expectation Step

Compute posterior probability  $\underbrace{P(C_i|x_j)}_{w_{ij}}$  for all  $i, j$ .

$$as \quad w_{ij} = \frac{f_i(x_j) P(C_i)}{\sum_{a \in I}^k f_a(x_j) P(C_a)}$$

## Maximization Step

Use posterior probabilities  $w_{ij}$  to re-compute  $M_i, \Sigma_i, P(C_i)$

$$M_i = \frac{\sum_{j=1}^n w_{ij} x_j}{\sum_{j=1}^n w_{ij}}$$

$$\Sigma_i = \frac{\sum_{j=1}^n w_{ij} (x_j - M_i)(x_j - M_i)^T}{\sum_{j=1}^n w_{ij}}$$

Considering the pairwise attribute view, the covariance between dimensions  $X_a$  and  $X_b$  is estimated as

$$\sigma_{ab}^i = \frac{\sum_{j=1}^n w_{ij} (x_{ja} - \mu_{ia})(x_{jb} - \mu_{ib})}{\sum_{j=1}^n w_{ij}}$$

$$P(C_i) = \frac{\sum_{j=1}^n w_{ij}}{n}$$

## K-means as a Special Case of EM

F-meas

Instead of multivariate normal cluster distribution, let

$$P(x_j | C_i) = \begin{cases} 1, & \text{if } \arg\max_{C_a} \{ \|x_j - \mu_a\|^2 \} \\ 0, & \text{otherwise} \end{cases}$$

posterior probabilities are binary

$$P(C_i | x_j) = \frac{P(x_j | C_i) P(C_i)}{\sum_{a=1}^k P(x_j | C_a) P(C_a)}$$

$$= \begin{cases} 0, & \text{if } i \neq \arg\max_i \{ \|x_j - \mu_i\|^2 \} \\ 1, & \text{if } i = \arg\max_i \{ \|x_j - \mu_i\|^2 \} \end{cases}$$

0 if 0  
1 if 1

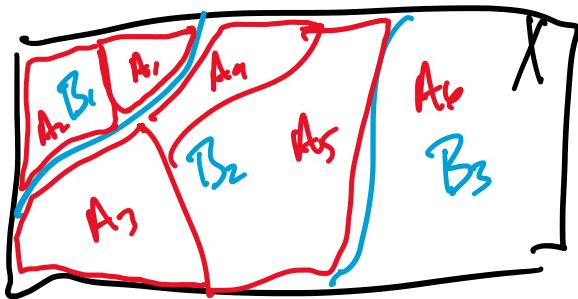
→ only  $\mu_i + P(C_i)$  are parameters

⇒ Simplifies to K-mean

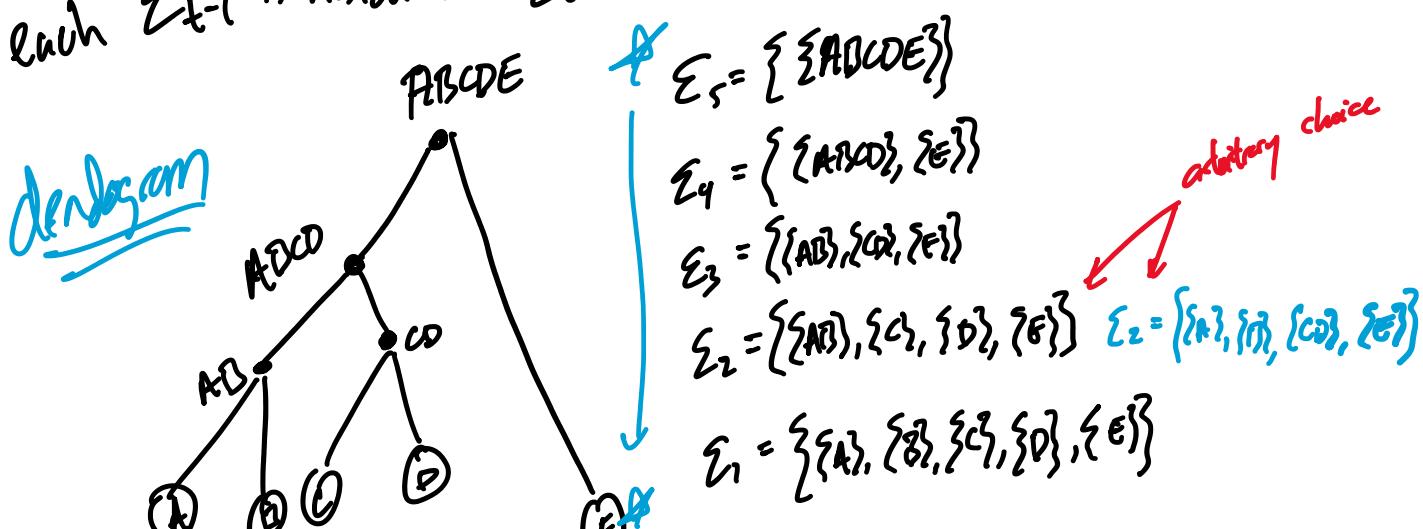
## Hierarchical Clustering

$X \in \mathbb{R}^{n \times d}$ ,  $x_i \in \mathbb{R}^d$ ,  $\Sigma = \{E_1, \dots, E_k\}$  set of clusters  
 (  $E_1, \dots, E_k$  partition  $X$  )  
 $E_i \cap E_j = \emptyset$  for  $i \neq j$   
 $E_1 \cup \dots \cup E_k = X$

Consider clusterings  $\mathcal{A} = \{A_1, \dots, A_r\}$  and  $\mathcal{B} = \{B_1, \dots, B_s\}$   
 $\mathcal{A}$  is nested in  $\mathcal{B}$  if for  $r > s$ , each  $A_i \subseteq B_j$  for some  $i, j$



Hierarchical clustering produces a hierarchy of clusterings  
 from  $\Sigma_1 = \{\{x_1\}, \{x_2\}, \dots, \{x_n\}\}$  up to  $\Sigma_n = \{\{x_1, \dots, x_n\}\}$  where  
 each  $\Sigma_{t-1}$  is nested in  $\Sigma_t$





The cluster dendrogram captures the nesting structure.  
If  $E_i \subset E_j$ , we draw an edge between them.

Number of Hierarchical Clusterings :  $(2n-3)!! \rightarrow$  enumerative techniques fail

$$(2n-3)(2n-5)(2n-7) \dots 1$$

## Agglomerative Hierarchical Clustering

Agglomerative techniques - "bottom-up" methods start with each pt. in its own cluster and merges them based on some criteria

Repeatedly merge the two closest clusters. If  $\mathcal{E} = \{E_1, \dots, E_m\}$ , then if  $E_i, E_j$  are the closest,  $\Sigma \rightarrow \Sigma \setminus \{E_i, E_j\} \cup \underbrace{\{E_i \cup E_j\}}_{E_{ij}}$

Repeat until  $\Sigma$  is one cluster of all points or until  $|\Sigma| = k$  (prespecified)

---

### Algorithm 14.1: Agglomerative Hierarchical Clustering Algorithm

---

#### AGGLOMERATIVECLUSTERING( $\mathbf{D}, k$ ):

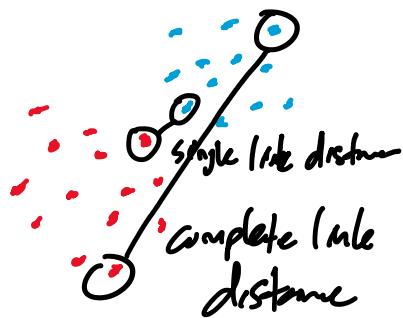
- 1  $\mathcal{C} \leftarrow \{C_i = \{\mathbf{x}_i\} \mid \mathbf{x}_i \in \mathbf{D}\}$  // Each point in separate cluster
- 2  $\Delta \leftarrow \{\|\mathbf{x}_i - \mathbf{x}_j\| : \mathbf{x}_i, \mathbf{x}_j \in \mathbf{D}\}$  // Compute distance matrix
- 3 **repeat**
- 4     Find the closest pair of clusters  $C_i, C_j \in \mathcal{C}$
- 5      $C_{ij} \leftarrow C_i \cup C_j$  // Merge the clusters
- 6      $\mathcal{C} \leftarrow (\mathcal{C} \setminus \{C_i, C_j\}) \cup \{C_{ij}\}$  // Update the clustering
- 7     Update distance matrix  $\Delta$  to reflect new clustering
- 8 **until**  $|\mathcal{C}| = k$

## Distance between Clusters

Several distance measures exist. - based on distance between points, e.g.  $L^2$  distance.

Single Link  $\delta(E_i, E_j) = \min \{ \|x-y\| : x \in E_i, y \in E_j \}$

min dist. between a point in  $E_i$  + a point in  $E_j$



Complete Link  $\delta(E_i, E_j) = \max \{ \|x-y\| : x \in E_i, y \in E_j \}$

max dist. between a pt in  $E_i$  + a pt in  $E_j$ . If we connected points from clusters with distance at most  $\delta(E_i, E_j)$ , we would connect all pairs.

Group Average  $\delta(E_i, E_j) = \frac{\sum_{x \in E_i} \sum_{y \in E_j} \|x-y\|}{n_i \cdot n_j}$

Average pairwise distance ( $|C_i| = n_i$ )

Mean Distance  $\delta(\bar{E}_i, E_j) = \|M_i - M_j\|$

Mean dist between centroids

Min Variance: Ward's Method

$$SSE_i = \sum_{x \in E_i} \|x - M_i\|^2 = \sum_{x \in E_i} \|x\|^2 - n_i \|M_i\|^2$$

$$\dots \rightarrow \leftarrow : \frac{m}{2} SSE$$

$$SSE \text{ of } \Sigma : \sum_{i=1}^m SSE_i$$

$$\Rightarrow \delta(E_i, E_j) = \Delta SSE_{ij} = SSE_{ij} - SSE_i - SSE_j$$

Net change in SSE if  $E_i, E_j$  are merged

$$\delta(E_i, E_j) = n_i \|M_i\|^2 + n_j \|M_j\|^2 - (n_i + n_j) \|M_{ij}\|^2$$

$$\text{where } M_{ij} = \frac{n_i M_i + n_j M_j}{n_i + n_j} \quad (\text{weighted average})$$

$$\delta(E_i, E_j) = \frac{n_i n_j}{n_i + n_j} \|M_i - M_j\|^2$$

When  $E_i, E_j$  merge to form  $E_{ij}$ , we need to update the distance from  $E_{ij}$  to each  $E_r, r \neq i, j$

Lance-Williams formula allows recomputing distances for all the cluster proximity measures

$$\delta(E_{ij}, E_r) = \alpha_i \delta(E_i, E_r) + \alpha_j \delta(E_j, E_r) + \beta \delta(E_i, E_j) + \gamma |\delta(E_i, E_r) - \delta(E_j, E_r)|$$

A

Measure	$\alpha_i$	$\alpha_j$	$\beta$	$\gamma$
Single Link	$\frac{1}{2}$	$\frac{1}{2}$	0	$\frac{-1}{2}$
Complete Link	$\frac{1}{2}$	$\frac{1}{2}$	0	$\frac{1}{2}$
Group Average	$\frac{n_i}{n_i + n_j}$	$\frac{n_j}{n_i + n_j}$	0	0
Mean Distance	$\frac{n_i}{n_i + n_j}$	$\frac{n_j}{n_i + n_j}$	$\frac{-n_i n_j}{(n_i + n_j)^2}$	0

Mean Distance	$\frac{n_i}{n_i+n_j}$	$\frac{n_j}{n_i+n_j}$	$\frac{-n_i n_j}{(n_i+n_j)^2}$	0
Ward's Measure	$\frac{n_i + n_j}{n_i + n_j + n_r}$	$\frac{n_j + n_r}{n_i + n_j + n_r}$	$\frac{-n_r}{n_i + n_j + n_r}$	0

Computational Complexity:  $O(n^2 \log n)$

