

# Expectation Maximization

Data is generated from a mixture of  $k$  Gaussians, one per cluster

$$P(x_j) = \sum_{i=1}^k \underbrace{f(x_j | \mu_i, \Sigma_i)}_{\text{Normal Distribution}} \cdot \underbrace{P(c_i)}_{\text{Prior prob mixture coefficient}}$$

$$\theta = \{ \mu_1, \Sigma_1, P(c_1), \mu_2, \Sigma_2, P(c_2), \dots, \}$$

Maximize the likelihood

$$P(D|\theta) \equiv L = \prod_{j=1}^n P(x_j | \theta) \quad : \text{figure the "best" } \theta$$

log likelihood

$$\ln L = \sum_{j=1}^n \ln(P(x_j | \theta))$$

$$\frac{\partial \ln L}{\partial \mu_i}$$

$$\frac{\partial \ln L}{\partial \Sigma_i}$$

$$\frac{\partial \ln L}{\partial P(c_i)}$$

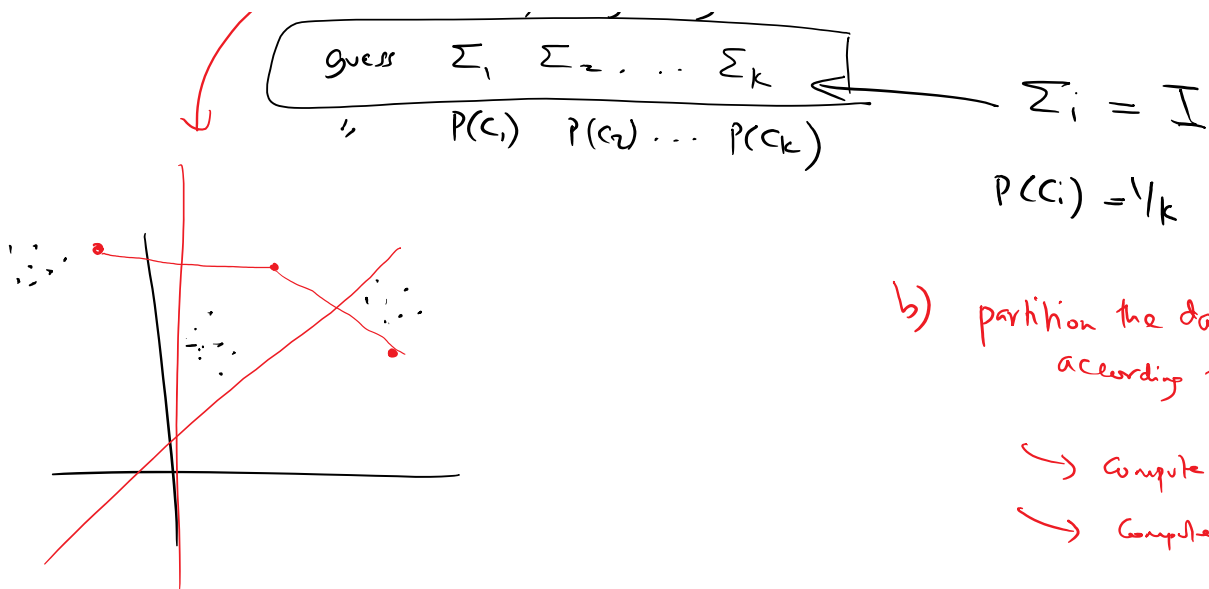
① Initialize

guess  $\vec{\mu}_1, \vec{\mu}_2, \dots, \vec{\mu}_k$

guess  $\Sigma_1, \Sigma_2, \dots, \Sigma_k$

Pick any set of  $k$  points in  $D$

$\Sigma_i = T$



b) partition the data according to closest mean

→ compute  $\Sigma_i$

→ compute  $P(c_i)$

② Expectation step:

given  $\vec{\mu}_i, \Sigma_i, P(c_i) \quad \forall i=1, \dots, k$

for all points  $\vec{x}_j \in D$

$$\boxed{\begin{array}{l} w_{ij} \\ \text{weight of} \\ \vec{x}_j \text{ for } c_i \end{array}} = P(c_i | x_j) = \frac{P(x_j | c_i) \cdot P(c_i)}{\sum_{a=1}^k P(x_j | c_a) \cdot P(c_a)}$$

$$N(x_j | \vec{\mu}_i, \Sigma_i)$$

$$\boxed{O(n \cdot k \cdot d)} \text{ operations}$$

$$O(nkd) + O(kd^3)$$

$$\frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} e^{-\frac{1}{2} (\vec{x}_j - \vec{\mu}_i)^T \Sigma_i^{-1} (\vec{x}_j - \vec{\mu}_i)}$$

Computing the inverse

$$\boxed{\sum_{i=1}^k = 1}$$

$$) \cdot k$$

Total prob of point  $x_j$

e.g.  $k=3$ .

$$\omega_{1j} = 0.5 \quad \omega_{2j} = 0.3 \quad \omega_{3j} = 0.2$$

② Maximization step

Given  $\omega_{ij} \forall c_i, \vec{x}_j$ ,

Update  $\vec{\mu}_i, \vec{\Sigma}_i, P(c_i)$

$$\vec{\mu}_i = \frac{\sum_{j=1}^n \omega_{ij} \cdot \vec{x}_j}{\sum \omega_{ij}}$$

weighted mean!

$$\vec{\Sigma}_i = \frac{\sum_{j=1}^n (\vec{x}_j - \vec{\mu}_i) (\vec{x}_j - \vec{\mu}_i)^T \cdot \omega_{ij}}{\sum_{j=1}^n \omega_{ij}}$$

weighted Covariance matrix

$$P(c_i) = \frac{\sum_{j=1}^n \omega_{ij}}{n}$$

Repeat E-M steps  
until convergence

(e.g. diff among the  $k$  means)

	$c_1$	$c_2$	$c_3$	
$\vec{x}_1$	0.5	0.5	0.2	1
$\vec{x}_2$	0.8	0.1	0.1	1
$\vec{x}_3$	0.1	0.7	0.2	1
	1.4	1.1	0.5	$\boxed{n}$
	2	2	2	

k-means : hard clustering

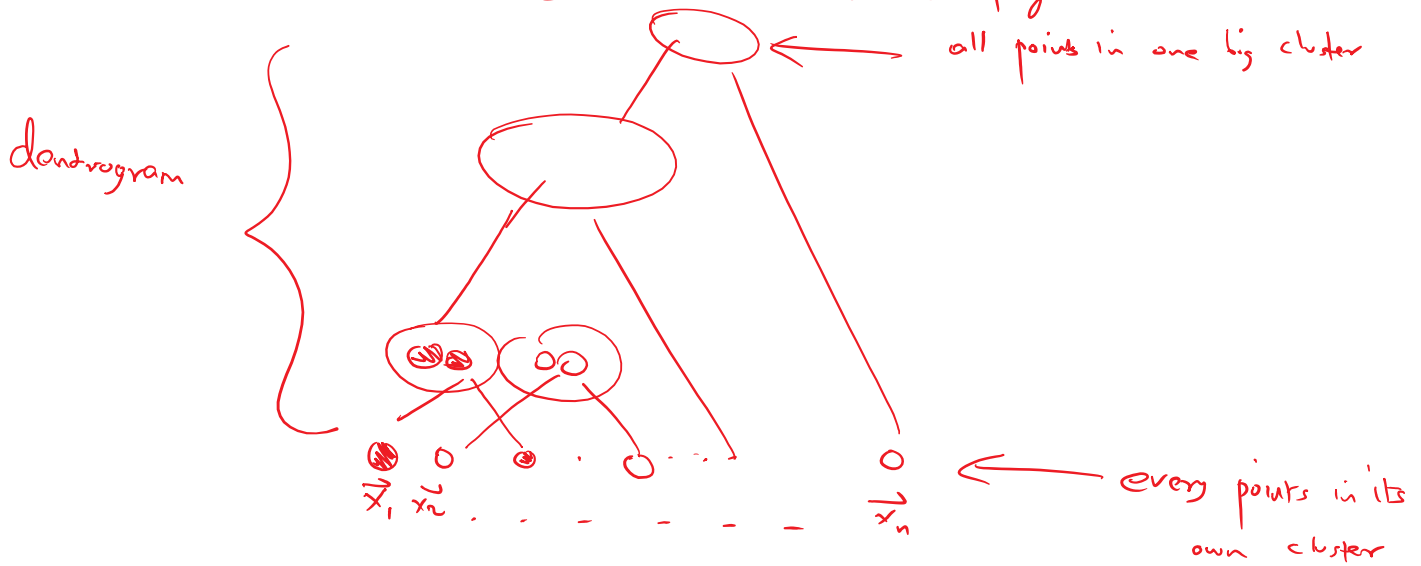
0/1 : prob of a point belonging to a cluster.

EM : soft clustering

$P(C_i | x_j)$  : prob of  $x_j$  belonging to  $C_i$

## Hierarchical clustering

→ to extract the "tree" of groupings



a) bottom up: agglomerative  
merge groups from bottom to top

b) top down: divisive  
recursive partitioning

graph clustering

# of clusters

① every point is in its own group

$$C_i = \{ \vec{x}_i \} \quad \forall i=1, \dots, n \quad | \quad n - \text{clusters}$$

to find

② merge the two closest clusters

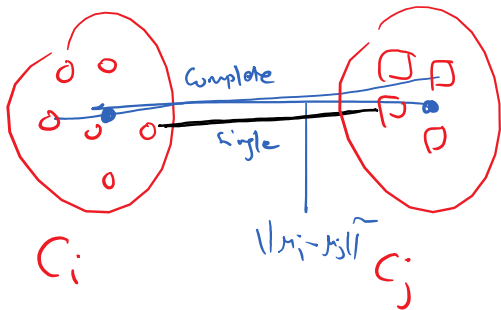
$$C_1 \cup C_2, \dots, C_{i-1} \cup C_{i+1}, \dots, C_n$$

(2) merge the two closest clusters

$S(C_i, C_j)$  : distance between clusters  
 $\forall$  all pairs  
 If  $C_i$  &  $C_j$  are the closest clusters  
 $C_{ij} = C_i \cup C_j$   
 remove  $C_i$  &  $C_j$   $O(n^2)$

(3) Update the distances from  $C_{ij}$  to every other cluster  
 repeat (2) until # of clusters  $\leq k$

How to measure distance between clusters



1) Single link

$$S(C_i, C_j) = \min \left\{ s(x_a, x_b) \mid \begin{array}{l} x_a \in C_i \\ x_b \in C_j \end{array} \right\}$$

2) Complete link

$$S(C_i, C_j) = \max \{ \quad \}$$

3) Avg link / Avg Distance

$$S(C_i, C_j) = \text{mean} \left\{ s(x_a, x_b) \mid x_a \in C_i, x_b \in C_j \right\}$$

Avg pair-wise distance

4)  $S(C_i, C_j) = \|\vec{\mu}_i - \vec{\mu}_j\|^2$   
 distance between the means.

5)  $S(C_i, C_j) = \Delta SSE_{ij}$

Change in the SSE values due to the merge

$$\text{sum of squared error} \quad \underline{\underline{SSE(C_i)}} = \sum_{x_j \in C_i} \underline{\underline{\| \vec{x}_j - \vec{\mu}_i \|^2}}$$

$$C_i \text{ \& } C_j \leftarrow \text{merge them into } \underline{C_{ij} = C_i \cup C_j}$$

$$\underline{SSE(C_{ij})}$$

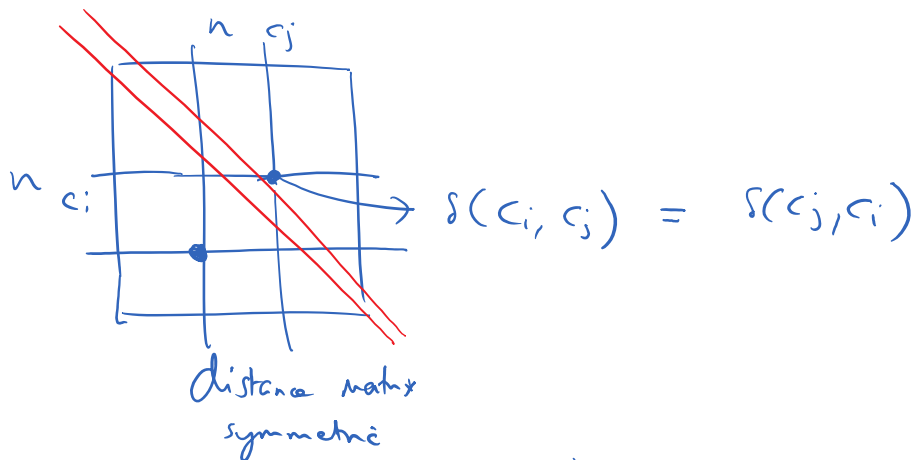
$$\underline{SSE(C_i) + SSE(C_j)}$$

$$\Delta SSE_{ij} = SSE(C_{ij}) - SSE(C_i) - SSE(C_j)$$

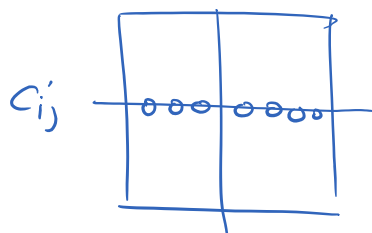
Smaller the better.

$$s(C_i, C_j) = \Delta SSE_{ij} = \underbrace{\left( \frac{n_i n_j}{n} \right)}_{\text{weighted}} \underbrace{\| \vec{\mu}_i - \vec{\mu}_j \|^2}$$

Updating the distance matrix



Remove  $C_i$  &  $C_j$   
create a new distance row for  $C_{ij}$



$$\underline{s(C_{ij}, C_k)} \quad \forall k \neq i, j$$



$$s(c_{ij}, c_r) \quad \forall k \neq i, j$$

updating depends on the distance function

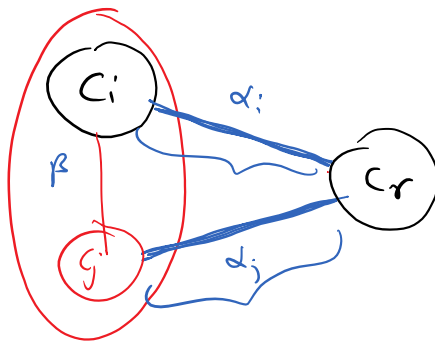
$$s(c_{ij}, c_r) = \alpha_i \cdot \underline{s(c_i, c_r)} + \alpha_j \cdot s(c_j, c_r) + \beta(c_i, c_j)$$

$r \neq i$   
 $r \neq j$

$$+ \gamma \cdot |s(c_i, c_r) - s(c_j, c_r)|$$

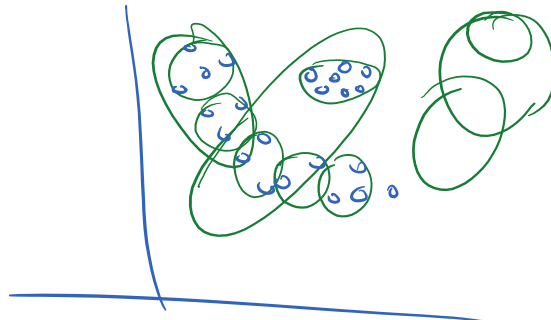
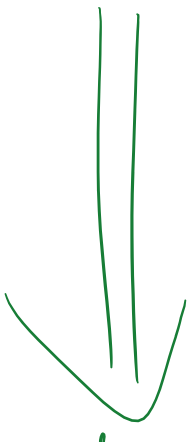
$$\alpha_i, \alpha_j, \beta, \gamma$$

$$\text{Single link: } \alpha_i = \alpha_j = 1/2 \quad \gamma = -1/2$$



difference between the distance  
 $\gamma$

## Density-based clustering

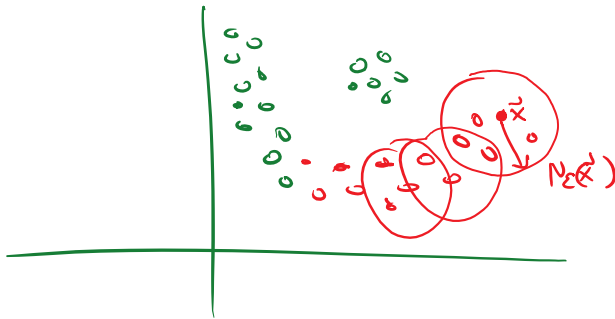


← hierarchical clustering has the ability to capture non-convex shape



design a method that is more suitable

non-convex  
shape



DBSCAN

$$N_{\epsilon}(\vec{x}) = \left\{ \vec{y} \mid \|\vec{x} - \vec{y}\| \leq \epsilon \right\}$$

radius

$\epsilon$ -Neighborhood

Core point  $\vec{x}$  :  $|N_{\epsilon}(\vec{x})| \geq \underline{\text{minpts}}$

border point :  $|N_{\epsilon}(\vec{x})| < \text{minpts}$

And  $\vec{x} \in N_{\epsilon}(\vec{y})$  for some core point  $\vec{y}$ .

Noise point: otherwise

① Compute  $N_{\epsilon}(\vec{x}_i) \forall \vec{x}_i \in D$

nearest neighbor search

② based on minpts  
label the core points

$O(n^2)$

③ for each unlabeled core point  $x$ :

label  $(x_i) \leftarrow \text{clusterid}$

also label all points in

$N_{\epsilon}(x_i)$

recursively jump to any core

point  $\rightarrow \dots$

Input parameters are

$\epsilon \leftarrow$  radius

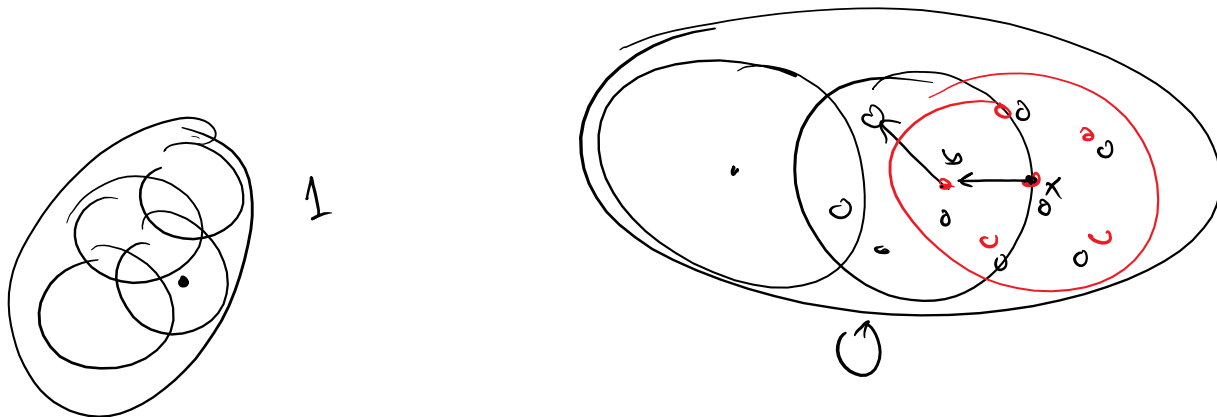
$\text{minpts} \leftarrow$  threshold for core.

Connected components  
over the core  
point areas



recursively jump to any core  
point  $\vec{y} \in N_{\epsilon}(x_i)$  & repeat

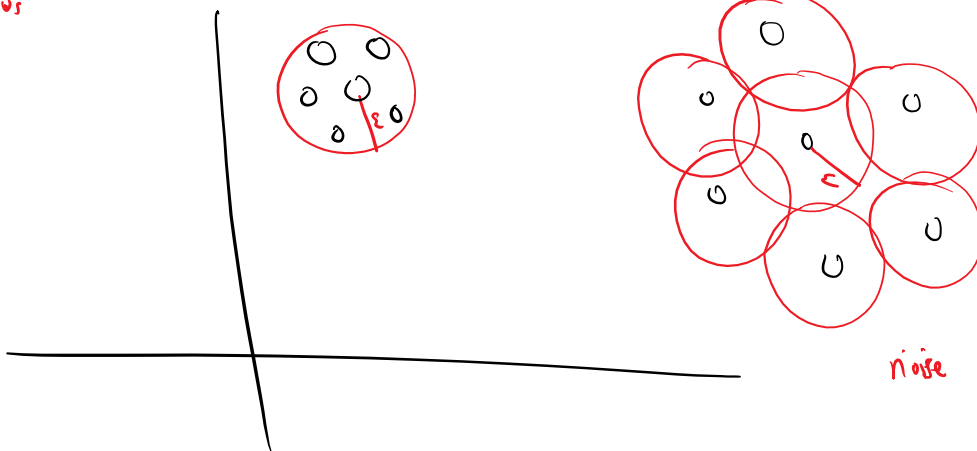
over the core  
point graph.



(4) any remaining unlabeled point is noise

$\epsilon \leftarrow$  radius  
is fixed

minpts



Kernel Density Estimation  $f(x)$



