

# Clustering

**Exemplifying the application of  
hierarchical agglomerative clustering  
(single-, complete- and average-linkage)**

CMU, 2012 fall, Tom Mitchell, Ziv Bar-Joseph, HW4, pr. 2.a  
extended by Liviu Ciortuz

The table below is a distance matrix for 6 objects.

	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>
<i>A</i>	0					
<i>B</i>	0.12	0				
<i>C</i>	0.51	0.25	0			
<i>D</i>	0.84	0.16	0.14	0		
<i>E</i>	0.28	0.77	0.70	0.45	0	
<i>F</i>	0.34	0.61	0.93	0.20	0.67	0

Show the final result of hierarchical clustering with single-, complete- and average-linkage by drawing the corresponding dendrograms.

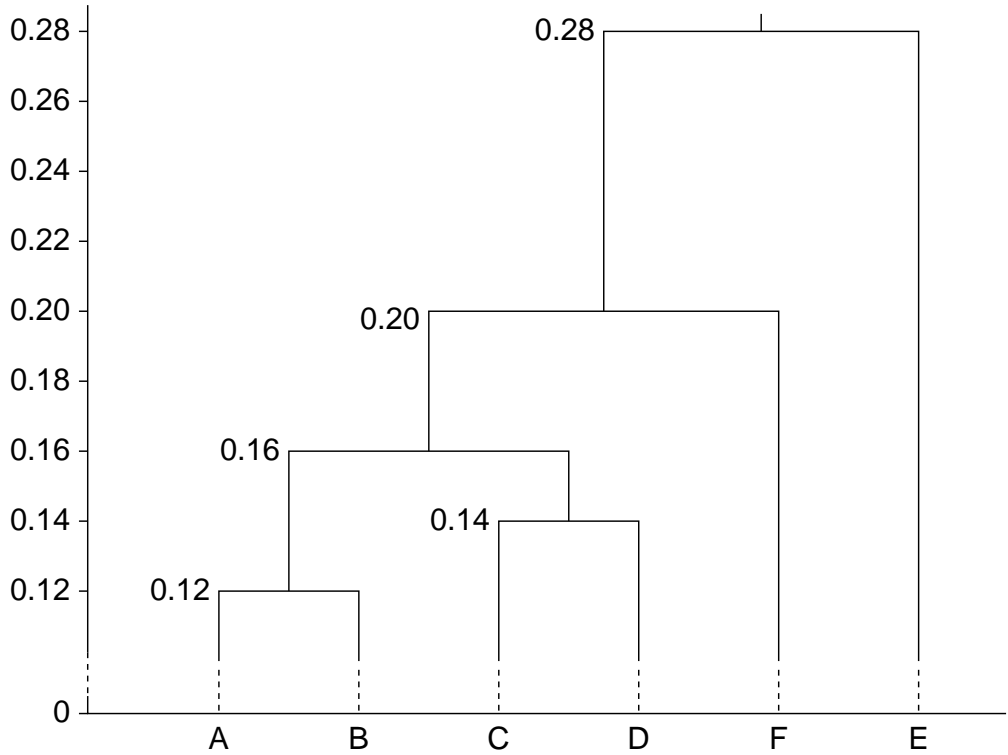
Solution:  
Single-linkage:

	<i>AB</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>
<i>AB</i>	0				
<i>C</i>	0.25	0			
<i>D</i>	0.16	<b>0.14</b>	0		
<i>E</i>	0.28	0.70	0.45	0	
<i>F</i>	0.34	0.93	0.20	0.67	0

	<i>AB</i>	<i>CD</i>	<i>E</i>	<i>F</i>
<i>AB</i>	0			
<i>CD</i>	<b>0.16</b>	0		
<i>E</i>	0.28	0.45	0	
<i>F</i>	0.34	0.20	0.67	0

	<i>ABCD</i>	<i>E</i>	<i>F</i>
<i>ABCD</i>	0		
<i>E</i>	0.28	0	
<i>F</i>	<b>0.20</b>	0.67	0

	<i>ABCDF</i>	<i>E</i>
<i>ABCDF</i>	0	
<i>E</i>	<b>0.28</b>	0



## Complete-linkage:

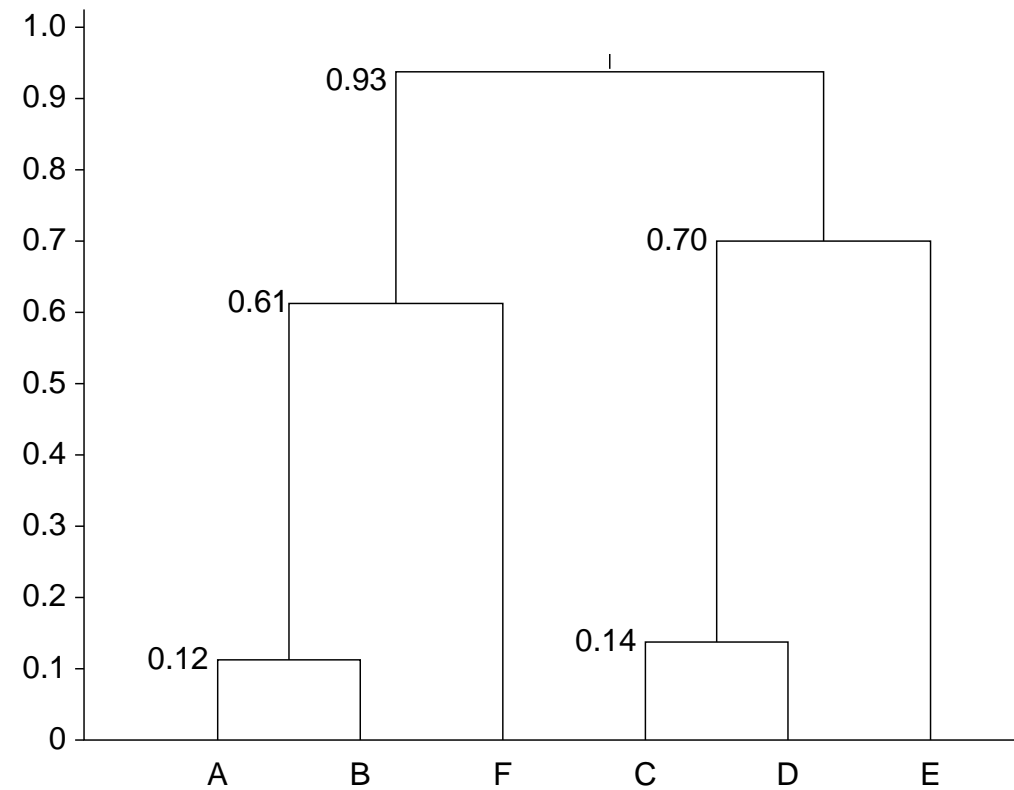
4.

	<i>AB</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>
<i>AB</i>	0				
<i>C</i>	0.51	0			
<i>D</i>	0.84	<b>0.14</b>	0		
<i>E</i>	0.77	0.70	0.45	0	
<i>F</i>	0.61	0.93	0.20	0.67	0

	<i>AB</i>	<i>CD</i>	<i>E</i>	<i>F</i>
<i>AB</i>	0			
<i>CD</i>	0.84	0		
<i>E</i>	0.77	0.70	0	
<i>F</i>	<b>0.61</b>	0.93	0.67	0

	<i>ABF</i>	<i>CD</i>	<i>E</i>
<i>ABF</i>	0		
<i>CD</i>	0.93	0	
<i>E</i>	0.77	<b>0.70</b>	0

	<i>ABF</i>	<i>CDE</i>
<i>ABF</i>	0	
<i>CDE</i>	<b>0.93</b>	0



## Average-linkage:

	<i>AB</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>
<i>AB</i>	0				
<i>C</i>	0.38	0			
<i>D</i>	0.50	<b>0.14</b>	0		
<i>E</i>	0.525	0.70	0.45	0	
<i>F</i>	0.475	0.93	0.20	0.67	0

	<i>AB</i>	<i>CD</i>	<i>E</i>	<i>F</i>
<i>AB</i>	0			
<i>CD</i>	<b>0.44</b>	0		
<i>E</i>	0.525	0.575	0	
<i>F</i>	0.475	0.565	0.67	0

	<i>ABCD</i>	<i>E</i>	<i>F</i>
<i>ABCD</i>	0		
<i>E</i>	0.55	0	
<i>F</i>	<b>0.52</b>	0.67	0

$$d(AB, C) = \frac{d(A, C) + d(B, C)}{2} = \frac{0.51 + 0.25}{2} = 0.38$$

$$d(AB, D) = \frac{d(A, D) + d(B, D)}{2} = \frac{0.84 + 0.16}{2} = 0.5$$

$$d(AB, E) = \frac{d(A, E) + d(B, E)}{2} = \frac{0.28 + 0.77}{2} = 0.525$$

$$d(AB, F) = \frac{d(A, F) + d(B, F)}{2} = \frac{0.34 + 0.61}{2} = 0.475$$

$$d(AB, CD) = \frac{2d(AB, C) + 2d(AB, D)}{4} \stackrel{(*)}{=} \frac{0.38 + 0.50}{2} = 0.44$$

$$d(CD, E) = \frac{d(C, E) + d(D, E)}{2} = \frac{0.76 + 0.45}{2} = 0.575$$

$$d(CD, F) = \frac{d(C, F) + d(D, F)}{2} = \frac{0.93 + 0.20}{2} = 0.565$$

$$d(ABCD, E) \stackrel{(*)}{=} \frac{2d(AB, E) + 2d(CD, E)}{4} = \frac{0.525 + 0.575}{2} = 0.55$$

$$d(ABCD, F) \stackrel{(*)}{=} \frac{2d(AB, F) + 2d(CD, F)}{4} = \frac{0.475 + 0.565}{2} = 0.52$$

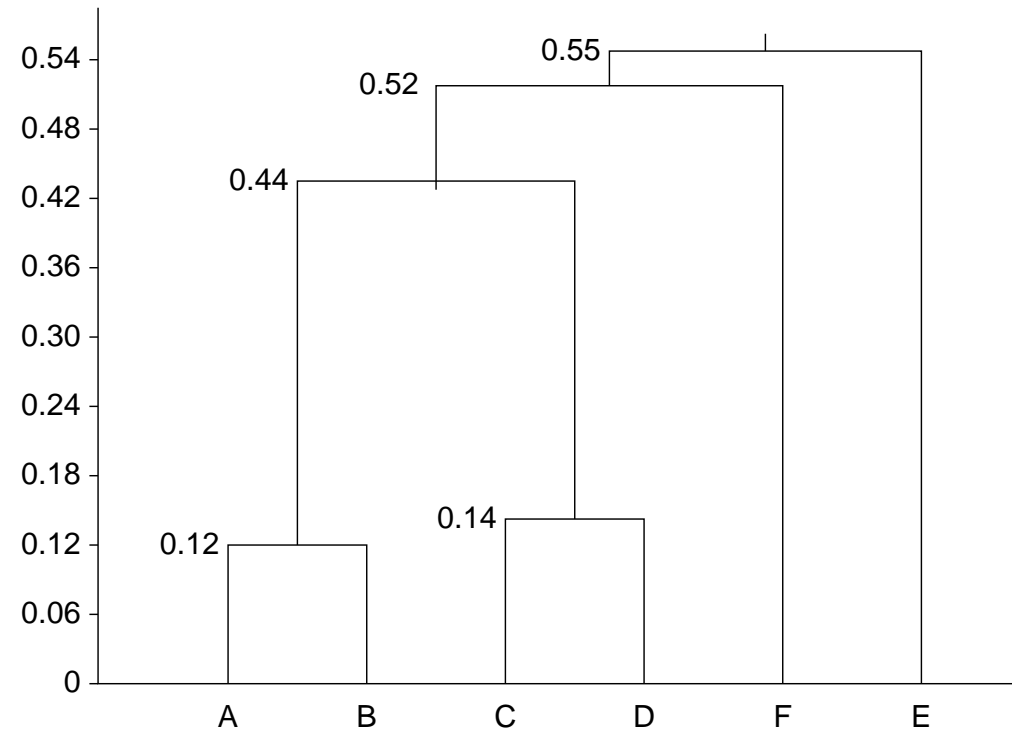
## Average-linkage (cont'd):

	<i>ABCD</i>	<i>F</i>
<i>ABCD</i>	0	
<i>F</i>	<b>0.55</b>	0

$$d(ABCD, F) \stackrel{(*)}{=} \frac{4d(ABCD, E) + d(F, E)}{5} = \frac{2.87}{5} = 0.574$$

**Note:** For the proof of the (\*) relation, see problem CMU, 2010 fall, Aarti Singh, HW3, pr. 4.2:

$$d(X \cup Y, Z) = \frac{|X| d(X, Z) + |Y| d(Y, Z)}{|X| + |Y|}$$



**Exemplifying**  
the application of hierarchical divisive clustering  
and the relationship between single-linkage hierarchies and  
**Minimum Spanning Trees (MSTs)**

CMU, 2009 spring, Ziv Bar-Joseph, final exam, pr. 9.3



Hierarchical clustering may be bottom-up or top-down.

In this problem we will see whether a top-down clustering algorithm can be exactly analogous to a bottom-up clustering algorithm.

Consider the following *top-down clustering algorithm*:

1. Calculate the pairwise distance  $d(P_i, P_j)$  between every two objects  $P_i$  and  $P_j$  in the set of objects to be clustered, and build a complete graph on the set of objects with edge weights being the corresponding distances.
2. Generate the Minimum Spanning Tree of the graph, i.e. choose the subset of edges  $E'$  with minimum sum of weights such that  $G' = (P, E')$  is a single connected tree.
3. Throw out the edge with the heaviest weight to generate two disconnected trees corresponding to top level clusters.
4. Repeat the previous step recursively on the lower level clusters to generate a top-down clustering on the set of  $n$  objects.

- a. Apply this algorithm on the dataset given in the nearby table, using the Euclidian distance.
- b. Does this top-down algorithm perform analogously to any bottom-up algorithm that you have encountered in class? Why?

Point	$x$	$y$
$P_1$	1	2
$P_2$	2	2
$P_3$	3	6
$P_4$	6	4
$P_5$	6	6
$P_6$	12	12

Solution:

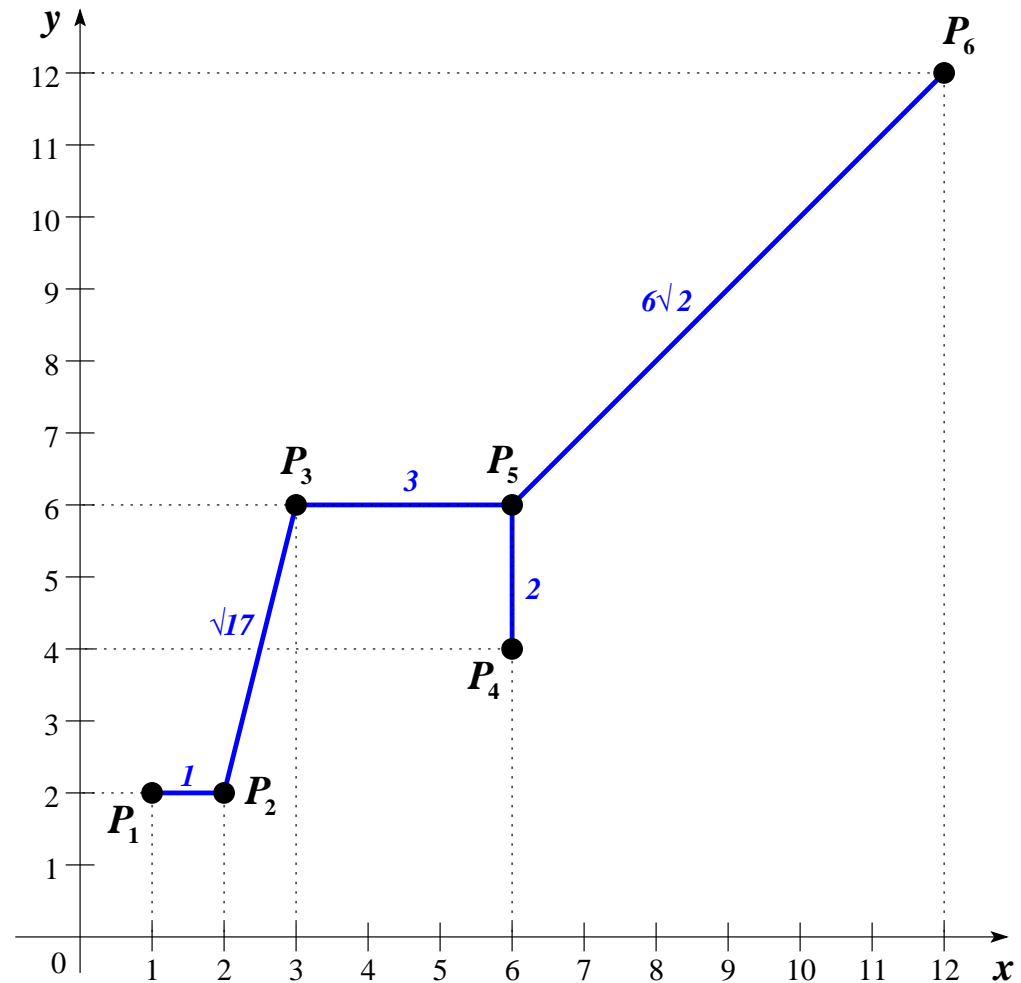
a.

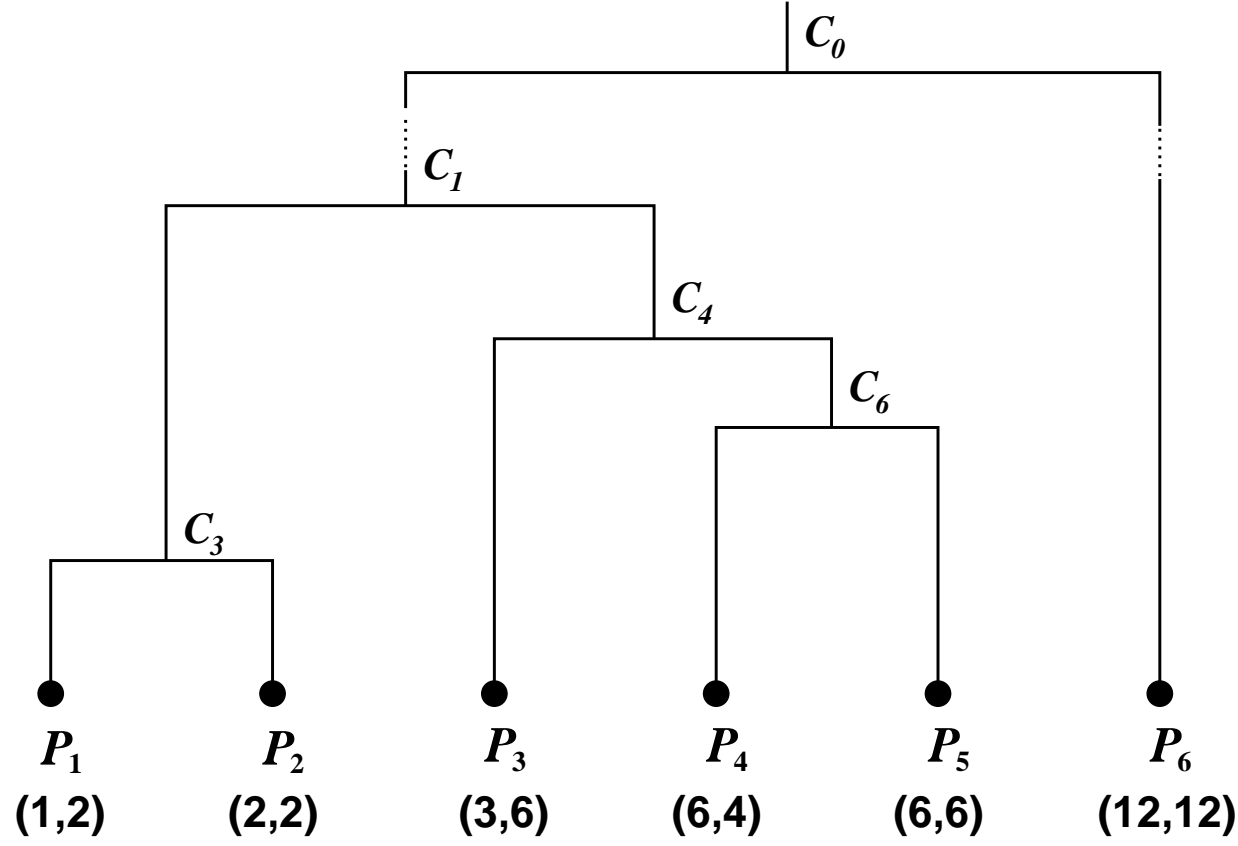
Kruskal algorithm:

1.  $(P_1, P_2)$ , cost 1,
2.  $(P_4, P_5)$ , cost 2,
3.  $(P_3, P_5)$ , cost 3,
4.  $(P_2, P_3)$ , cost  $\sqrt{17}$
5.  $(P_5, P_6)$ , cost  $6\sqrt{2}$ .

Prim algorithm:

1.  $(P_1, P_2)$ , cost 1,
2.  $(P_2, P_3)$ , cost  $\sqrt{17}$
3.  $(P_3, P_5)$ , cost 3,
4.  $(P_5, P_4)$ , cost 2,
5.  $(P_5, P_6)$ , cost  $6\sqrt{2}$ .





*Note:* If there is only one MST for the given dataset, then both Kruskal's and Prim's algorithm will find it. Otherwise, the two algorithms can produce different results.

One can see (both on this dataset and also in general) that Kruskal's algorithm is exactly analogous to the single-linkage bottom-up clustering algorithm.

Therefore, there is indeed a bottom-up equivalent to the top-down clustering algorithm presented in this exercise.

Hierarchical [top-down] clustering

Ward's metric

CMU, 2010 fall, Aarti Singh, HW3, pr. 4.1

In this problem you will analyze an alternative approach to quantify the distance between two disjoint clusters, proposed by Joe H. Ward in 1963. We will call it Ward's metric.

Ward's metric simply says that the distance between two disjoint clusters,  $X$  and  $Y$ , is how much the **sum of squares** will increase when we merge them. More formally,

$$\Delta(X, Y) = \sum_{i \in X \cup Y} \|x_i - \mu_{X \cup Y}\|^2 - \sum_{i \in X} \|x_i - \mu_X\|^2 - \sum_{i \in Y} \|x_i - \mu_Y\|^2 \quad (1)$$

where  $\mu_i$  is the centroid of cluster  $i$  and  $x_i$  is a data point in a given cluster.

Here,  $\Delta(X, Y)$  can be thought as the *merging cost* of combining clusters  $X$  and  $Y$  into one cluster. That is, in agglomerative clustering those two clusters with the lowest *merging cost* is merged using the Ward's metric as a *closeness* measure.

a. Can you reduce the formula in Equation (1) for  $\Delta(X, Y)$  to a simpler form? Give the simplified formula.

*Hint:* Your formula should be in terms of the cluster sizes (let's denote them as  $n_X$  and  $n_Y$ ) and the distance  $\|\mu_X - \mu_Y\|^2$  between cluster centroids  $\mu_X$  and  $\mu_Y$  *only*.



## Solution

$$\begin{aligned}
\Delta(X, Y) &= \sum_{i \in X \cup Y} \|x_i - \mu_{X \cup Y}\|^2 - \sum_{i \in X} \|x_i - \mu_X\|^2 - \sum_{i \in Y} \|x_i - \mu_Y\|^2 \\
&= \sum_{i \in X \cup Y} \left( x_i - \frac{n_X \mu_X + n_Y \mu_Y}{n_X + n_Y} \right)^2 - \sum_{i \in X} (x_i - \mu_X)^2 - \sum_{i \in Y} (x_i - \mu_Y)^2 \\
&= \sum_{i \in X \cup Y} x_i^2 - \left( \frac{2n_X \mu_X + 2n_Y \mu_Y}{n_X + n_Y} \right) \sum_{i \in X \cup Y} x_i + \sum_{i \in X \cup Y} \left( \frac{n_X \mu_X + n_Y \mu_Y}{n_X + n_Y} \right)^2 \\
&\quad - \sum_{i \in X} x_i^2 + 2\mu_X \sum_{i \in X} x_i - \sum_{i \in X} \mu_X^2 - \sum_{i \in Y} x_i^2 + 2\mu_Y \sum_{i \in Y} x_i - \sum_{i \in Y} \mu_Y^2 \\
&= -\frac{2(n_X \mu_X + n_Y \mu_Y)^2}{n_X + n_Y} + \frac{(n_X \mu_X + n_Y \mu_Y)^2}{n_X + n_Y} + 2\mu_X n_X \mu_X + 2\mu_Y n_Y \mu_Y - n_X \mu_X^2 - n_Y \mu_Y^2 \\
&= -\frac{(n_X \mu_X + n_Y \mu_Y)^2}{n_X + n_Y} + n_X \mu_X^2 + n_Y \mu_Y^2 = \frac{n_X n_Y}{n_X + n_Y} (\mu_X^2 - 2\mu_X \mu_Y + \mu_Y^2) \\
&= \frac{n_X n_Y}{n_X + n_Y} \|\mu_X - \mu_Y\|^2
\end{aligned}$$

b. Give an interpretation for Ward's metric. What do you think it is trying to achieve?

*Hint:* The simplified formula from above will be helpful to answer this part.

Solution:

With hierarchical clustering, the sum of squares starts out at zero (because every point is in its own cluster) and then grows as we merge clusters. Ward's metric keeps this growth in sum of squares as small as possible. This is nice if we believe that the sum of squares (as a measure of cluster coherence) should be small.

Notice that the number of points also shows up in  $\Delta$ , as well as their geometric separation (Harmonic mean). The *intuition* is that given two pairs of clusters whose centers are equally far apart, Ward's method will prefer to merge the smaller ones.

c. Assume that you are given two *pairs* of clusters  $P_1$  and  $P_2$ . The centers of the two clusters in  $P_1$  is farther apart than the centers of the two clusters in  $P_2$ . Using Ward's metric, does agglomerative clustering always choose to merge the two clusters in  $P_2$  (those with less 'distance' between their centers)? Why (not)? Justify your answer with a simple example.

Solution:

No, not always. Which pair will be merged also depends on the size of the clusters. One simple counter example is where the size of the clusters in  $P_1$  are 1 and 99, respectively; and similarly 50 and 50 for  $P_2$ .

Then the harmonic mean  $\frac{2n_X n_Y}{n_X + n_Y}$  of cluster sizes is  $2 \cdot 0.99$  for  $P_1$  and  $2 \cdot 25$  for  $P_2$ .

If the distance between the centers of the clusters in  $P_1$  is less than  $25/0.99 = 25.25$  times the distance between the centers of the clusters in  $P_2$ , then  $\Delta(P_1)$  will be still smaller and so clusters in  $P_1$  are merged.

d. In clustering it is usually not trivial to decide what is the right number of clusters the data falls into. Using Ward's metric for agglomerative clustering, can you come up with a simple heuristic to pick the number of clusters  $k$ ?

Solution:

Ward's algorithm can give us a hint to pick a reasonable  $k$  through the merging cost. If the cost of merging increases a lot, it is probably going too far, and losing a lot of structure.

So one possible *heuristic* is to keep reducing  $k$  until the cost jumps, and then use the  $k$  right before the jump. In other words, pick the  $k$  just before the merging cost takes off.

Exemplifying non-hierarchical clustering  
using the *K*-means algorithm

T.U. Dresden, 2006 summer, Steffen Höldobler, Axel Grossmann, HW3

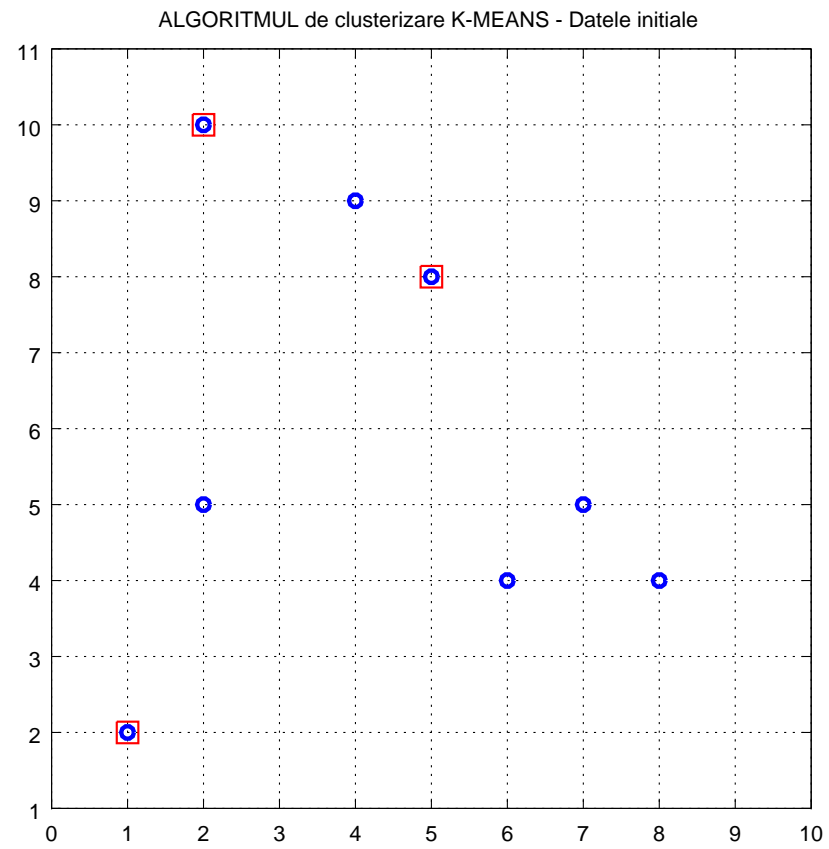
Folosiți algoritmul  $K$ -means și distanța euclidiană pentru a grupa următoarele 8 instanțe din  $\mathbb{R}^2$  în 3 clustere:

$$A(2, 10), B(2, 5), C(8, 4), D(5, 8), E(7, 5), F(6, 4), G(1, 2), H(4, 9).$$

Se vor lua drept centroizi inițiali punctele  $A$ ,  $D$  și  $G$ .

a. Rulați prima iterație a algoritmului  $K$ -means. Pe un grid de valori  $10 \times 10$  veți marca instanțele date, pozițiile centroizilor la începutul primei iterații și componența fiecărui cluster la finalul acestei iterații. (Trasați mediatoarele segmentelor determinate de centroizi, ca separatori ai clusterelor.)

b. Câte iterații sunt necesare pentru ca algoritmul  $K$ -means să convergă? Desenați pe câte un grid rezultatul rulării fiecărei iterații.



**Solution:**

**Iteration 0:**

$$\left. \begin{array}{l} \mu_1^0 = (2, 10) \\ \mu_2^0 = (5, 8) \\ \mu_3^0 = (1, 2) \end{array} \right\} \Rightarrow$$

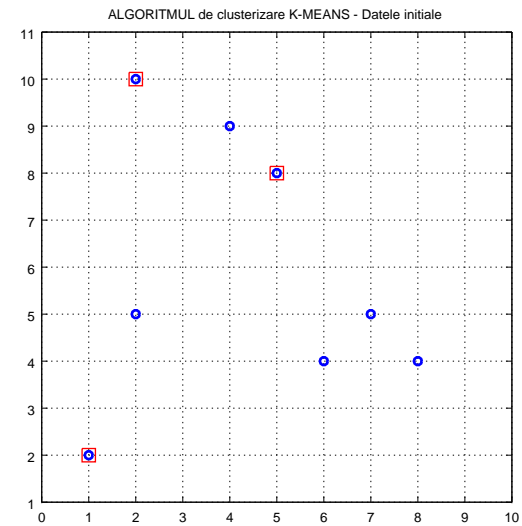
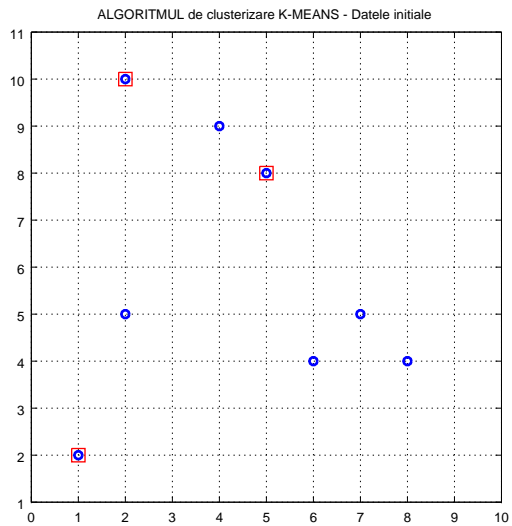
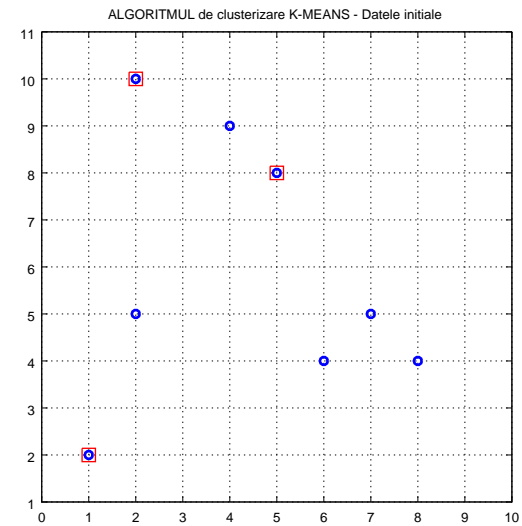
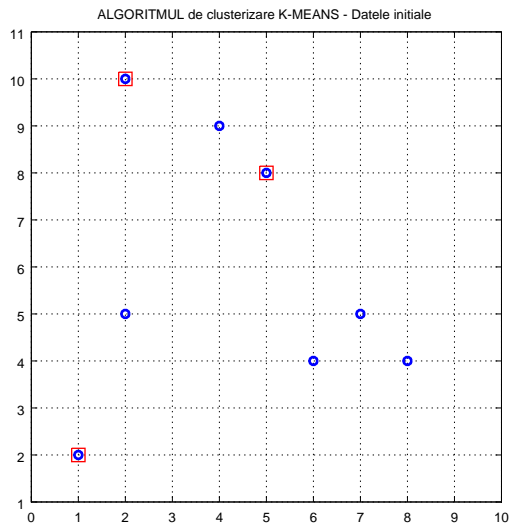
$P_i$	$d(\mu_1^0, P_i)$	$d(\mu_2^0, P_i)$	$d(\mu_3^0, P_i)$
$A(2, 10)$	0	...	...
$B(2, 5)$	5	$3\sqrt{2}$	$\sqrt{10}$
$C(8, 4)$			
$D(5, 8)$	...	0	...
$E(7, 5)$			
$F(6, 4)$			
$G(1, 2)$	...	...	0
$H(4, 9)$			

$$\Rightarrow \left\{ \begin{array}{l} C_1^0 = \{A\} \\ C_2^0 = \{C, D, E, F, H\} \\ C_3^0 = \{B, G\} \end{array} \right.$$

**Iteration 1:**

$$\left. \begin{array}{l} \mu_1^1 = \mu_1^0 = (2, 10) \\ \mu_2^1 = \left( \frac{4+5+6+7+8}{5}, \frac{4+4+5+8+9}{5} \right) = (6, 6) \\ \mu_3^1 = \left( \frac{2+1}{2}, \frac{5+2}{2} \right) = (1.5, 3.5) \end{array} \right\} \Rightarrow \dots \Rightarrow \left\{ \begin{array}{l} C_1^0 = \{A, H\} \\ C_2^0 = \{C, D, E, F\} \\ C_3^0 = \{B, G\} \end{array} \right.$$



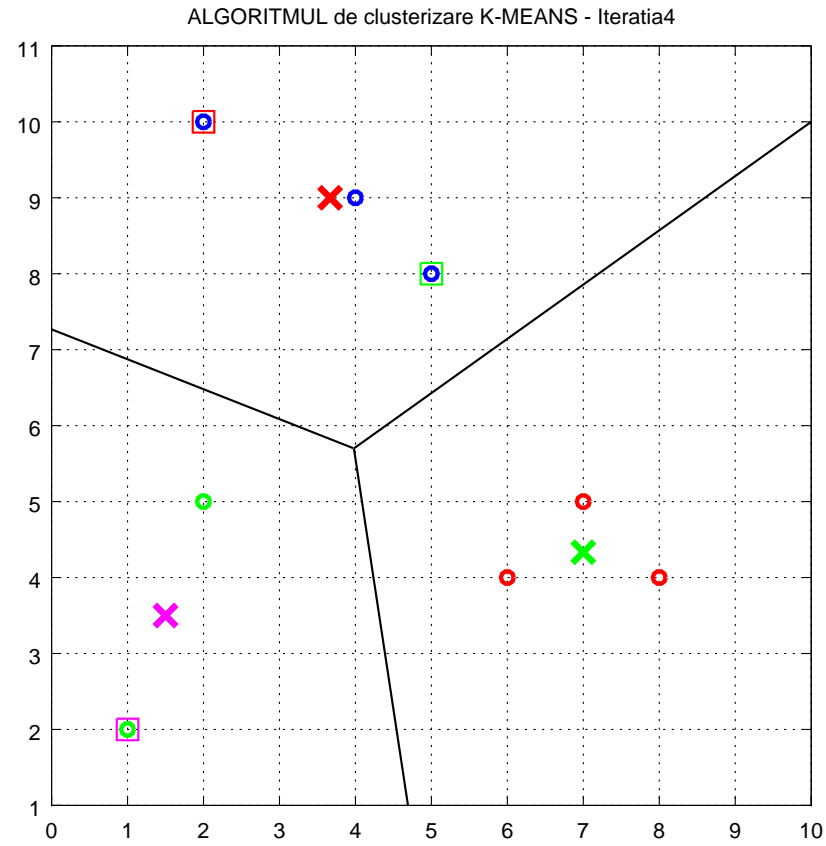


**Iteration 2:**

$$\left. \begin{array}{l} \mu_1^2 = (3, 9.5) \\ \mu_2^2 = \left(\frac{26}{4}, \frac{21}{4}\right) = (6.5, 5.25) \\ \mu_3^2 = \mu_3^1 = (1.5, 3.5) \end{array} \right\} \Rightarrow \dots \Rightarrow \left\{ \begin{array}{l} C_1^2 = \{A, D, H\} \\ C_2^2 = \{C, E, F\} \\ C_3^2 = \{B, G\} \end{array} \right.$$

**Iteration 3:**

$$\left. \begin{array}{l} \mu_1^3 = \left(\frac{2+4+5}{3}, \frac{8+9+10}{3}\right) = (11/3, 9) \\ \mu_2^3 = (7, 13/3) \\ \mu_3^3 = \mu_3^2 = (1.5, 3.5) \end{array} \right\} \Rightarrow \dots \Rightarrow \left\{ \begin{array}{l} C_1^3 = \{A, D, H\} = C_1^2 \\ C_2^3 = \{C, E, F\} = C_2^2 \\ C_3^3 = \{B, G\} = C_3^2 \end{array} \right\} \Rightarrow \textit{Stop}$$



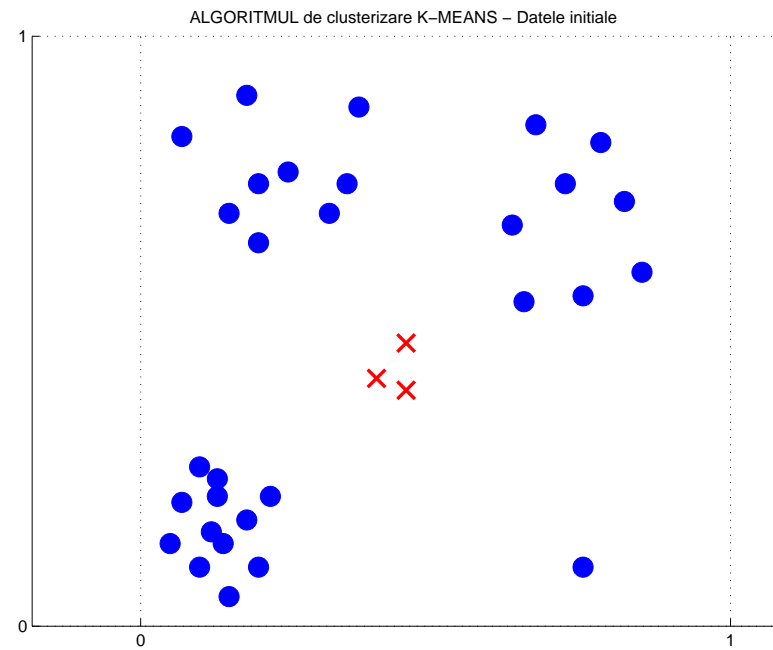
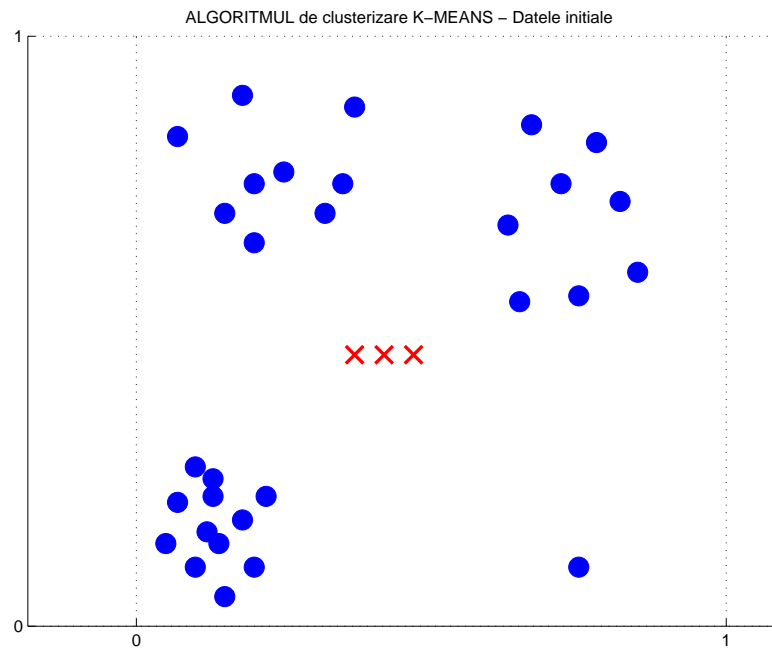
Exemplifying one property of  $K$ -means:

**The clusterisation result depends on initialisation**

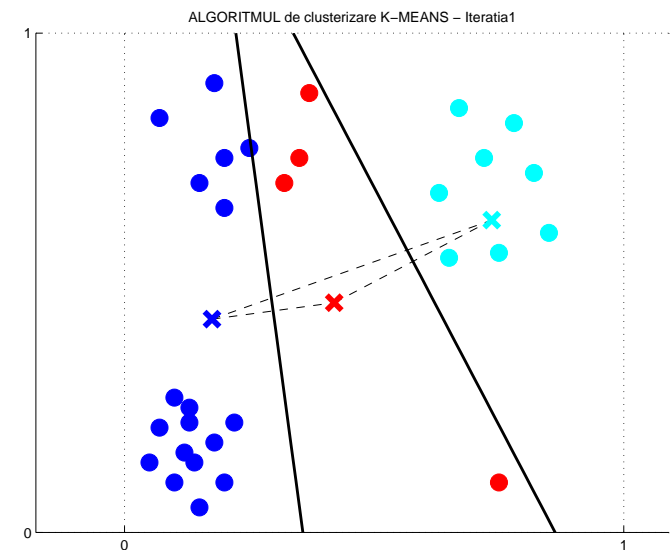
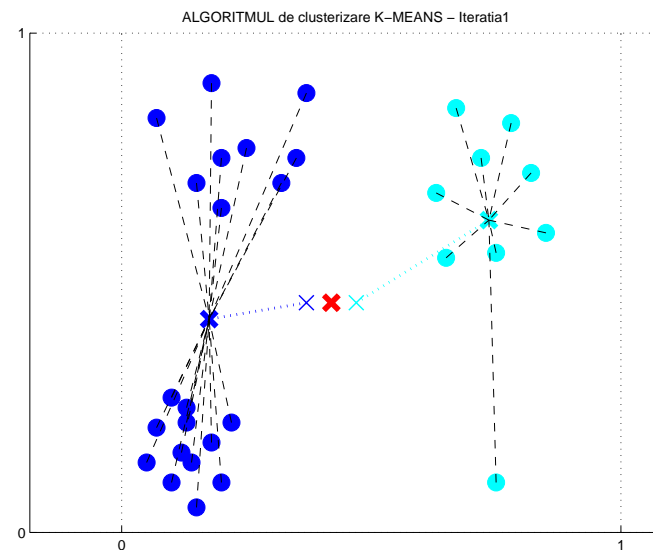
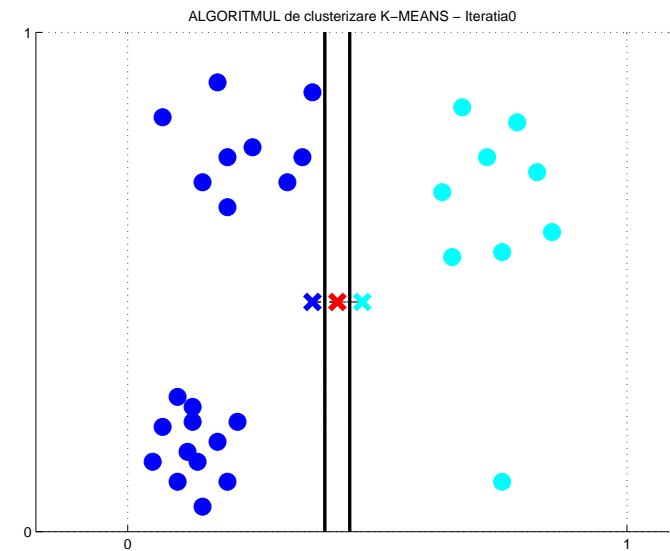
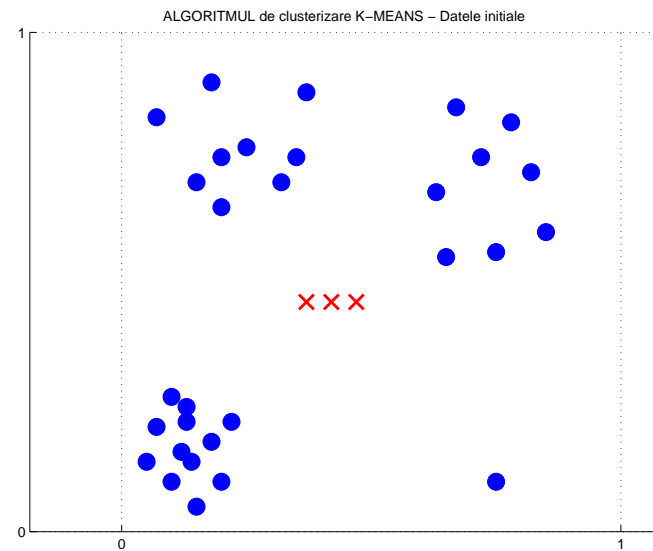
CMU, 2006 spring, Carlos Guestrin, HW5, pr. 1

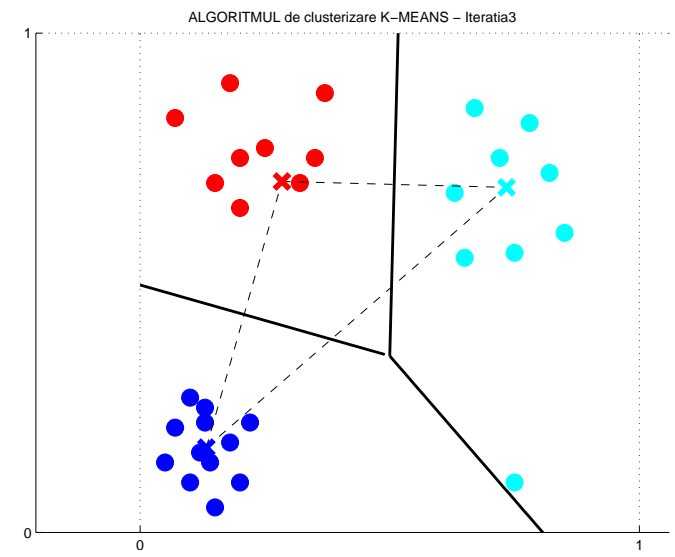
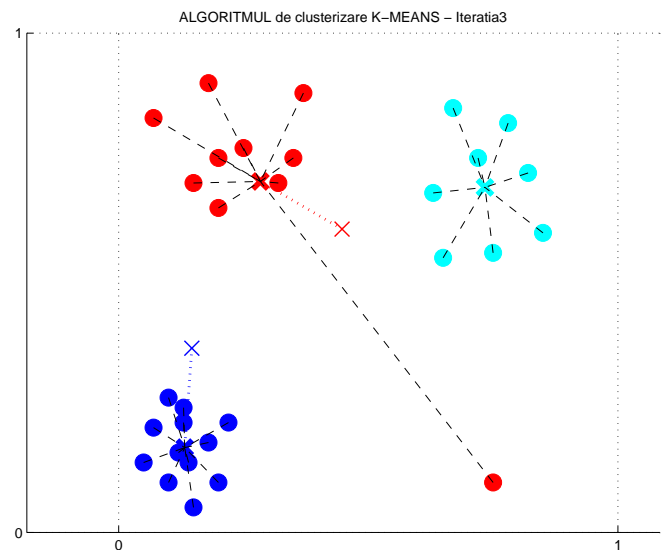
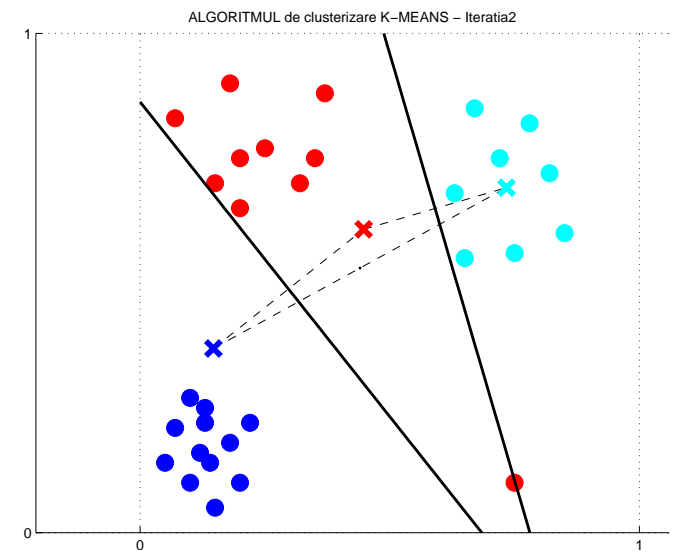
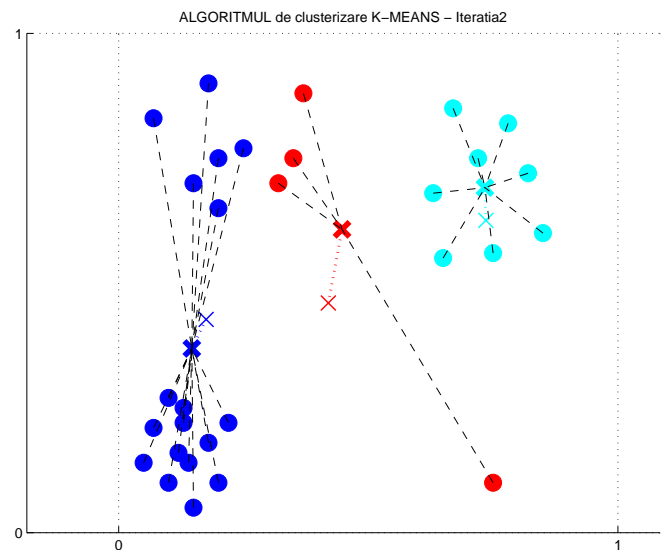
a-b. Consider the data set in the following figures. The  $\bullet$  symbols indicate data points, while the crosses ( $\times$ ) indicate the current cluster centers. For each one of these figures, show the progress of the  $K$ -means algorithm by showing how the class centers move with each iteration until convergence. For each iteration, indicate which data points will be associated with each of the clusters, as well as the updated class centers. If during the cluster update step, a cluster center has no points associated with it, it will not move. Use/produce as many figures you need until convergence of the algorithm.

c. What does this imply about the behavior of the K-Means algorithm?

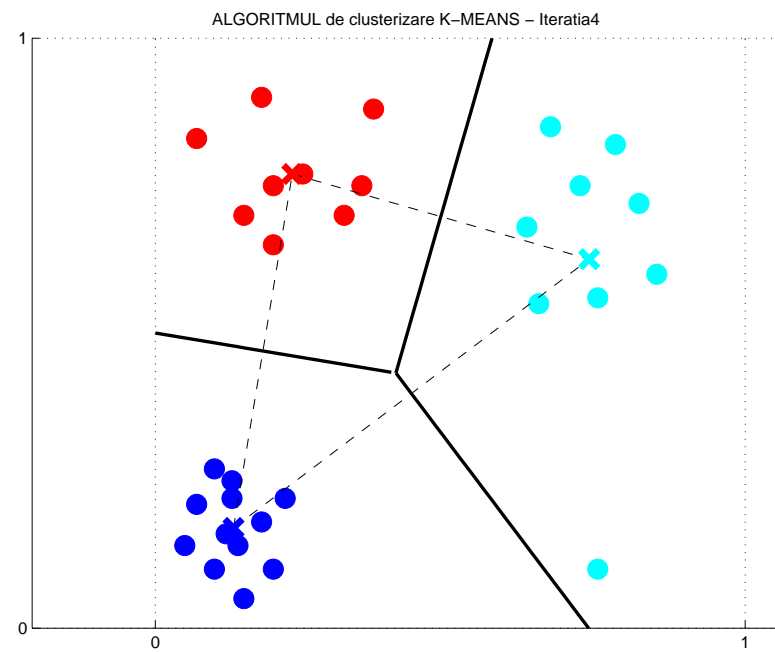
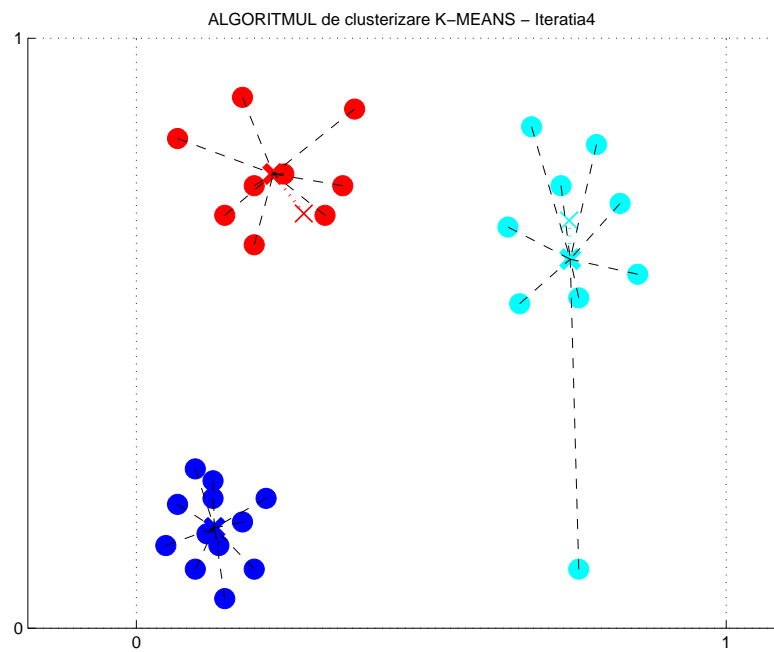


Solution:  
a.

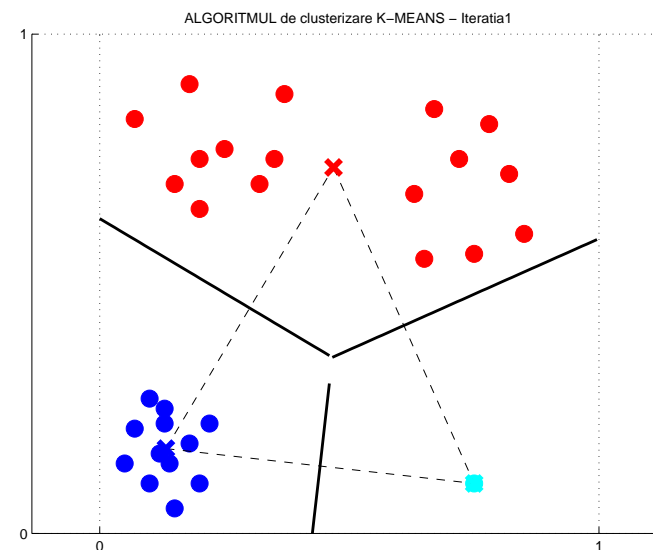
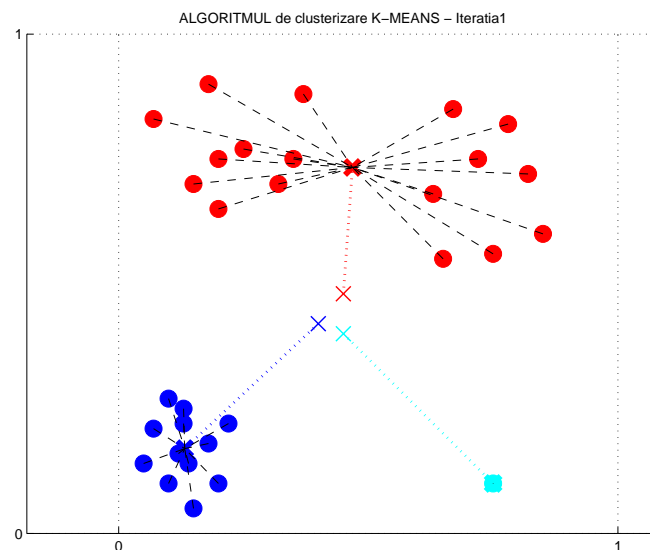
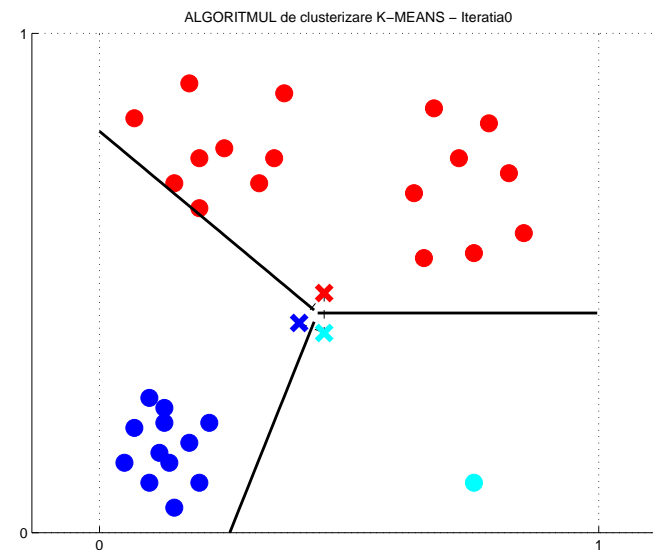
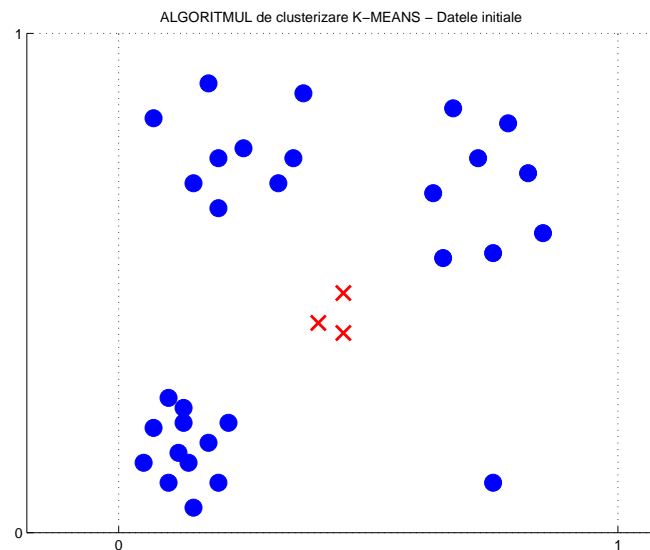








Solution:  
b.



### Solution: c.

Este evident din acest exercițiu că rezultatul algoritmului  $K$ -means depinde de poziționarea inițială a centroizilor. În cazul inițializării de la punctul  $a$  au fost necesare 4 iterații până a se ajunge la convergență, pe când la punctul  $b$  algoritmul a converș după doar o iterație.

Mai este încă ceva important *de remarcat*: faptul că la prima variantă de inițializare, punctul din dreapta jos, care este un *outlier* (rom., excepție, caz particular, aberație) este până la urmă asociat clusterului format de punctele din partea dreaptă (sus), în vreme ce la cea de-a doua variantă de inițializare el constituie un cluster aparte/“singleton”, obligând în mod indirect grupările de puncte din stânga-sus și dreapta-sus să formeze împreună un singur cluster.

# Some proofs

*K*-means as an optimisation algorithm:

The monotonicity of the  $J_K$  criterion

[CMU, 2009 spring, Ziv Bar-Joseph, HW5, pr. 2.1]

## Algoritmul $K$ -means (S. P. Lloyd, 1957)

**Input:**  $x_1, \dots, x_n \in \mathbb{R}^d$ , cu  $n \geq K$ .

**Output:** o anumită  $K$ -partiție pentru  $\{x_1, \dots, x_n\}$ .

**Procedură:**

[*Inițializare/Iterația 0:*]  $t \leftarrow 0$ ;

se fixează în mod arbitrar  $\mu_1^0, \dots, \mu_K^0$ , centroizii inițiali ai clusterelor, și se asignează fiecare instanță  $x_i$  la centroidul cel mai apropiat, formând astfel clusterelor  $C_1^0, \dots, C_K^0$ .

[*Recursivitate:*] Se execută iterația  $++ t$ :

**Pasul 1:** se calculează noile poziții ale centroizilor:

$$\mu_j^t = \frac{1}{|C_j^{t-1}|} \sum_{x_i \in C_j^{t-1}} x_i \text{ pentru } j = \overline{1, K};$$

**Pasul 2:**

se reasignează fiecare  $x_i$  la [clusterul cu] centroidul cel mai apropiat, adică se stabilește noua componență a clusterelor la iterația  $t$ :  $C_1^t, \dots, C_K^t$ ;

[*Terminare:*] până când o anumită condiție este îndeplinită

(de exemplu: până când pozițiile centroizilor — sau: componența clusterelor — nu se mai modifică de la o iterație la alta).

a. Demonstrați că, de la o iterație la alta, algoritmul  $K$ -means mărește *coeziunea de ansamblu* a clusterelor. I.e., considerând funcția

$$J(C^t, \mu^t) \stackrel{\text{def.}}{=} \sum_{i=1}^n \|x_i - \mu_{C^t(x_i)}^t\|^2 \stackrel{\text{def.}}{=} \sum_{i=1}^n (x_i - \mu_{C^t(x_i)}^t) \cdot (x_i - \mu_{C^t(x_i)}^t),$$

unde:

$C^t = (C_1^t, C_2^t, \dots, C_K^t)$  este colecția de clusterare (i.e.,  $K$ -partiția) la momentul  $t$ ,

$\mu^t = (\mu_1^t, \mu_2^t, \dots, \mu_K^t)$  este colecția de centroizi ai clusterelor ( $K$ -configurația)

la momentul  $t$ ,

$C^t(x_i)$  desemnează clusterul la care este asignat elementul  $x_i$  la iterația  $t$ ,

operatorul  $\cdot$  desemnează produsul scalar al vectorilor din  $\mathbb{R}^d$ ,

arătați că  $J(C^t, \mu^t) \geq J(C^{t+1}, \mu^{t+1})$  pentru orice  $t$ .

## Ideea demonstrației

Inegalitatea de mai sus rezultă din două inegalități (care corespund pașilor 1 și 2 de la iterația  $t$ ):

$$J(C^t, \mu^t) \stackrel{(1)}{\geq} J(C^t, \mu^{t+1}) \stackrel{(2)}{\geq} J(C^{t+1}, \mu^{t+1})$$

La prima inegalitate (cea corespunzătoare pasului 1) se poate considera că parametrul  $C^t$  este fixat iar  $\mu$  este variabil, în vreme ce la a doua inegalitate (cea corespunzătoare pasului 2) se consideră  $\mu^t$  fixat și  $C$  variabil.

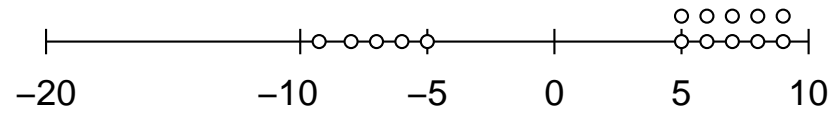
Prima inegalitate se poate obține însumând o serie de inegalități, și anume câte una pentru fiecare cluster  $C_j^t$ . A doua inegalitate se demonstrează imediat.

## Ilustrarea acestei idei, pe un exemplu particular:

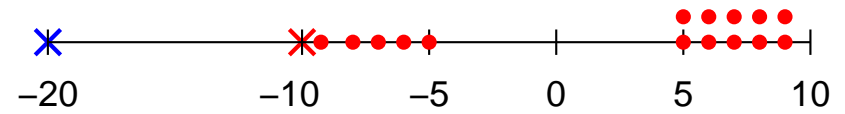
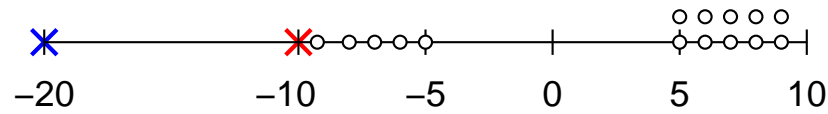
Vezi următoarele 3 slide-uri

[Edinburgh, 2009 fall, C. Williams, V. Lavrenko, HW4, pr. 3]

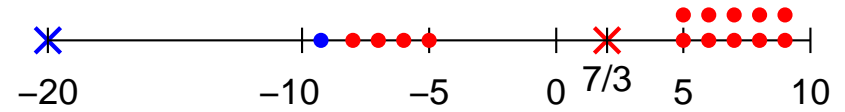
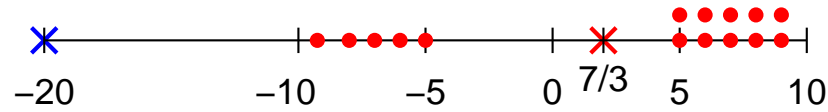




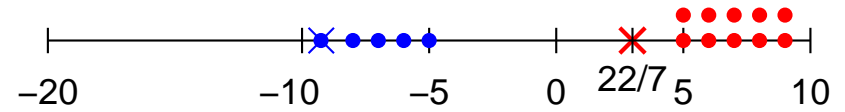
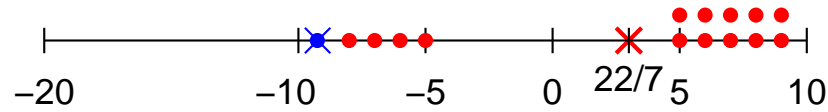
init. / iter. 0



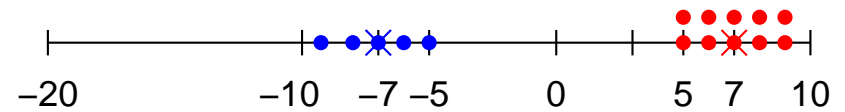
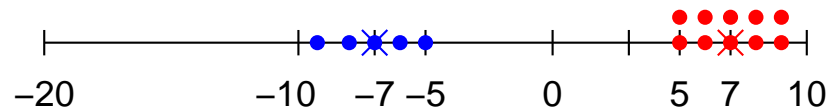
iter. 1



iter. 2



iter. 3



Pentru acest exemplu de aplicare a algoritmului  $K$ -means, scriem expresiile numerice pentru valoarea criteriului  $J_2(C^t, \mu^t)$  pentru fiecare iterație ( $t = 0, 1, 2, 3$ ).

iter.	$J_2(C^t, \mu^t)$
0.	$0 + \{(-9 - (-10))^2 + \dots + (-5 - (-10))^2 + 2[(5 - (-10))^2 + \dots + (9 - (-10))^2]\} \geq$
1.	$(-9 - (-20))^2 + \{(-8 - 7/3))^2 + \dots + (-5 - 7/3))^2 + 2[(5 - 7/3)^2 + \dots + (9 - 7/3)^2]\} \geq$
2.	$(-9 - (-9))^2 + \dots + (-5 - (-9))^2 + 2[(5 - 22/7)^2 + \dots + (9 - 22/7)^2] \geq$
3.	$(-9 - (-7))^2 + \dots + (-5 - (-7))^2 + 2[(5 - 7)^2 + \dots + (9 - 7)^2]$

**Observație:** La prima vedere, este greu să dovedim aceste inegalități ( $J_2(C^{t-1}, \mu^{t-1}) \geq J_2(C^t, \mu^t)$ , pentru  $t = 1, 2, 3$ ) ...altfel decât calculând efectiv valoarea expresiilor care se compară. Însă, introducând niște termeni intermediari, inegalitățile acestea se vor demonstra într-un mod foarte elegant...

iter.	$J_2(C^{t-1}, \mu^t)$	$J_2(C^t, \mu^t)$
0.		$0 + \{(-9 - (-10))^2 + \dots + (-5 - (-10))^2 + 2[(5 - (-10))^2 + \dots + (9 - (-10))^2]\} \geq$
1.	$0 + \{(-9 - 7/3)^2 + (-8 - 7/3)^2 + \dots + (-5 - 7/3)^2 + 2[(5 - 7/3)^2 + \dots + (9 - 7/3)^2]\} \geq$	$(-9 - (-20))^2 + \{(-8 - 7/3)^2 + \dots + (-5 - 7/3)^2 + 2[(5 - 7/3)^2 + \dots + (9 - 7/3)^2]\} \geq$
2.	$(-9 - (-9))^2 + \{(-8 - 22/7)^2 + \dots + (-5 - 22/7)^2 + 2[(5 - 22/7)^2 + \dots + (9 - 22/7)^2]\} \geq$	$(-9 - (-9))^2 + (-8 - (-9))^2 + \dots + (-5 - (-9))^2 + 2[(5 - 22/7)^2 + \dots + (9 - 22/7)^2] \geq$
3.	$(-9 - (-7))^2 + \dots + (-5 - (-7))^2 + 2[(5 - 7)^2 + \dots + (9 - 7)^2] =$	$(-9 - (-7))^2 + \dots + (-5 - (-7))^2 + 2[(5 - 7)^2 + \dots + (9 - 7)^2]$

### Explicații:

1. Inegalitățile pe orizontală ( $J_2(C^{t-1}, \mu^t) \geq J_2(C^t, \mu^t)$ , pentru  $t = 1, 2, 3$ ) sunt ușor de demonstrat, pe baza corespondenței termen cu termen. (Ele corespund eventualelor micșorări ale distanțelor atunci când se face reassignarea instanțelor la centroizi.)

2. Restul inegalităților ( $J_2(C^t, \mu^t) \geq J_2(C^t, \mu^{t+1})$ , pentru  $t = 1, 2, 3$ ) se rezolvă printr-o metodă de optimizare simplă. De exemplu, pentru  $t = 1$  este imediat că funcția  $(-9 - x)^2 + (-8 - x)^2 + \dots + (-5 - x)^2 + 2[(5 - x)^2 + \dots + (9 - x)^2]$  își atinge minimumul pentru  $x = 7/3$ , deci  $J_2(C^1, \mu^2) \geq J_2(C^1, \mu^1)$ .

## Demonstrație, pentru cazul general

**Observație:** Pentru conveniență, ne vom limita la cazul  $d = 1$ . Extinderea demonstrației la cazul  $d > 1$  nu comportă dificultăți.

**Demonstrarea inegalității (1):**  $J(C^t, \mu^t) \geq J(C^t, \mu^{t+1})$

(Vezi pasul 1 al iterației  $t$ .)

Fixăm  $j \in \{1, \dots, K\}$ . Dacă notăm cu  $C_j^t = \{x_{i_1}, x_{i_2}, \dots, x_{i_l}\}$ , unde  $l \stackrel{\text{not.}}{=} |C_j^t|$ , atunci

$$J(C_j^t, \mu_j^t) = \sum_{p=1}^l (x_{i_p} - \mu_j^t)^2, \text{ deci } J(C^t, \mu^t) = \sum_{j=1}^K J(C_j^t, \mu_j^t).$$

Dacă se consideră  $C_j^t$  fixat, iar  $\mu_j^t$  variabil, atunci putem minimiza imediat funcția

$$f(\mu) \stackrel{\text{def.}}{=} J(C_j^t, \mu) = l\mu^2 - 2\mu \sum_{p=1}^l x_{i_p} + \sum_{p=1}^l x_{i_p}^2 \Rightarrow \arg \min_{\mu} J(C_j^t, \mu) = \frac{1}{l} \sum_{p=1}^l x_{i_p} \stackrel{\text{def.}}{=} \mu_j^{t+1}.$$

Așadar,  $J(C_j^t, \mu) \geq J(C_j^t, \mu_j^{t+1})$ , pentru  $\forall \mu$ . În particular, pentru  $\mu = \mu_j^t$  vom avea:  $J(C_j^t, \mu_j^t) \geq J(C_j^t, \mu_j^{t+1})$ . Inegalitatea aceasta este valabilă pentru toate clusterelor  $j = 1, \dots, K$ . Dacă sumăm toate aceste inegalități, rezultă:  $J(C^t, \mu^t) \geq J(C^t, \mu^{t+1})$ .

## Demonstrarea inegalității (2): $J(C^t, \mu^{t+1}) \geq J(C^{t+1}, \mu^{t+1})$

(Vezi pasul 2 al iterației  $t$ .)

La acest pas, o instanță oarecare  $x_i$ , unde  $i \in \{1, \dots, n\}$ , este reassignată de la clusterul cu centroidul  $\mu_j^{t+1}$ , la un alt centroid  $\mu_q^{t+1}$ , dacă

$$\|x_i - \mu_{j'}^{t+1}\|^2 \geq \|x_i - \mu_q^{t+1}\|^2 \Leftrightarrow (x_i - \mu_{j'}^{t+1})^2 \geq (x_i - \mu_q^{t+1})^2, \text{ pentru orice } j' = 1, \dots, K.$$

În contextul iterației  $t$ , acest lucru implică

$$\left(x_i - \mu_{C^t(x_i)}^{t+1}\right)^2 \geq \left(x_i - \mu_{C^{t+1}(x_i)}^{t+1}\right)^2.$$

Sumând membru cu membru inegalitățile de acest tip obținute pentru  $i = \overline{1, n}$ , rezultă:  $J(C^t, \mu^{t+1}) \geq J(C^{t+1}, \mu^{t+1})$ , ceea ce era de demonstrat.

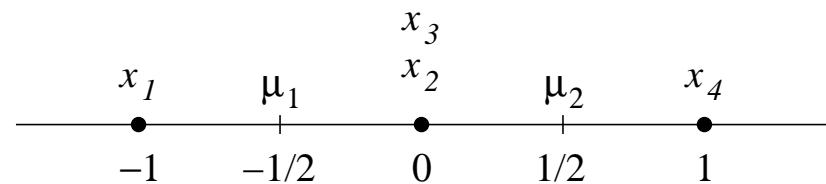
b. Ce puteți spune despre oprirea algoritmului  $K$ -means? (Termină oare acest algoritm într-un număr finit de pași, sau este posibil ca el să reviziteze o  $K$ -configurație anterioară  $\mu = (\mu_1, \dots, \mu_K)$ ?)

Răspuns:

Dacă algoritmul revizitează o  $K$ -partiție, atunci rezultă că pentru un anumit  $t$  avem  $J(C^{t-1}, \mu^t) = J(C^t, \mu^{t+1})$ . Este posibil ca acest fapt să se întâmple, și anume atunci când:

- există instanțe multiple (i.e.,  $x_i = x_j$ , deși  $i \neq j$ ),
- criteriul de oprire al algoritmului  $K$ -means este de forma “până când componența clusterelor nu se mai modifică”,
- se presupune că, în cazul în care o instanță  $x_i$  este situată la egală distanță față de doi sau mai mulți centroizi, ea poate fi asignată în mod aleatoriu la oricare dintre ei.

Așa se întâmplă în *exemplul* din figura alăturată dacă se consideră că la o iterație  $t$  avem  $x_2 = 0 \in C_1^t$  și  $x_3 = 0 \in C_2^t$ , iar la iterația următoare alegem ca  $x_3 = 0 \in C_1^{t+1}$  și  $x_2 = 0 \in C_2^{t+1}$  și, din nou, invers la iterația  $t + 2$ .



## Observații

- Dacă se păstrează criteriul dat ca exemplu în enunțul problemei – adică se iterează până când centroizii “staționează” – algoritmul se poate opri fără ca la ultima iterație  $J(C, \mu)$  să fi atins minimul posibil. În cazul exemplului de mai sus, vom avea  $\frac{1}{4} + 2 \cdot \frac{1}{4} + \frac{1}{4} = 1 > \frac{2}{3}$ .
- Dacă nu există instanțe multiple care să fie situate la distanțe egale față de doi sau mai mulți centroizi la o iterație oarecare a algoritmului  $K$ -means (precum sunt  $x_2$  și  $x_3$  în *exemplul* de mai sus), sau dacă se impune *restricția* ca în astfel de situații instanțele identice să fie asignate la un singur cluster, este evident că algoritmul  $K$ -means se oprește într-un număr finit de pași.



## Concluzii

- Algoritmul  $K$ -means explorează — pornind de la o anumită inițializare a celor  $K$  centroizi —, doar un subset din totalul de  $K^n$   $K$ -partiții, asigurându-ne însă că are loc proprietatea  $J(C^0, \mu^1) \geq J(C^1, \mu^2) \geq \dots \geq J(C^{t-1}, \mu^t) \geq J(C^t, \mu^{t+1})$ , conform punctului  $a$  al acestei probleme.
- **Atingerea minimului global** al funcției  $J(C, \mu)$  — unde  $C$  este o variabilă care parcurge mulțimea tuturor  $K$ -partițiilor care se pot forma cu instanțele  $\{x_1, \dots, x_n\}$  — **nu este garantată pentru algoritmul  $K$ -means**. Valoarea funcției  $J$  care se obține la oprirea algoritmului  $K$ -means este dependentă de plasarea inițială a centroizilor  $\mu$  precum și de modul concret în care sunt alcătuite clusterelor în cazul în care o instanță oarecare se află la distanță egală de doi sau mai mulți centroizi, după cum am arătat în exemplul de mai sus.

*K*-means algorithm:

The “approximate” maximization of the “distance” between clusters

[CMU, 2010 fall, Aarti Singh, HW3, pr. 5.2]

**Note:** In this problem we will work with a version of the  $K$ -means algorithm which is slightly modified w.r.t. the one given in the problem CMU, 2009 spring, Ziv Bar-Joseph, HW5, pr. 2.1, where we have proved the monotonicity of the criterion  $J$ .

Let  $X := \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  be our sample points, and  $K$  denote the number of clusters to use. We represent the cluster assignments of the data points by an *indicator matrix*  $\gamma \in \{0, 1\}^{n \times K}$  such that  $\gamma_{ij} = 1$  means  $\mathbf{x}_i$  belongs to cluster  $j$ . We require that each point belongs to exactly one cluster, so  $\sum_{j=1}^K \gamma_{ij} = 1$ .

[We already know that] the  $K$ -means algorithm “estimates”  $\gamma$  by minimizing the following “cohesion criterion” (or, “measure of distortion”, or simply “sum of squares”):

$$J(\gamma, \mu_1, \mu_2, \dots, \mu_K) := \sum_{i=1}^n \sum_{j=1}^K \gamma_{ij} \|\mathbf{x}_i - \mu_j\|^2,$$

where  $\|\cdot\|$  denotes the vector 2-norm.

$K$ -means alternates between estimating  $\gamma$  and re-computing  $\mu_j$ 's.

## The $K$ -means algorithm (...yet another version!)

- Initialize  $\mu_1, \mu_2, \dots, \mu_K$ , and let  $C := \{1, \dots, K\}$ .
- While the value of  $J$  is still decreasing, repeat the following:

1. Determine  $\gamma$  by

$$\gamma_{ij} \leftarrow \begin{cases} 1, & \|\mathbf{x}_i - \mu_j\|^2 \leq \|\mathbf{x}_i - \mu_{j'}\|^2, \forall j' \in C, \\ 0, & \text{otherwise.} \end{cases}$$

Break ties arbitrarily.

2. Recompute  $\mu_j$  using the updated  $\gamma$ :

For each  $j \in C$ , if  $\sum_{i=1}^n \gamma_{ij} > 0$  set

$$\mu_j \leftarrow \frac{\sum_{i=1}^n \gamma_{ij} \mathbf{x}_i}{\sum_{i=1}^n \gamma_{ij}}.$$

Otherwise, don't change  $\mu_j$ .

Let  $\bar{\mathbf{x}}$  denote the sample mean.

Consider the following three quantities:

$$\text{Total variation:} \quad V(X) = \frac{\sum_{i=1}^n \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2}{n}.$$

$$\text{Within-cluster variation:} \quad V_j(X) = \frac{\sum_{i=1}^n \gamma_{ij} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2}{\sum_{i=1}^n \gamma_{ij}}.$$

$$\text{Between-cluster variation:} \quad \tilde{V}(X) = \sum_{j=1}^K \left( \frac{\sum_{i=1}^n \gamma_{ij}}{n} \right) \|\boldsymbol{\mu}_j - \bar{\mathbf{x}}\|^2.$$

What is the relation between these three quantities?

Based on this relation, show that  $K$ -means can be interpreted as **minimizing a weighted average of within-cluster variations while approximately(!) maximizing the between-cluster variation**. Note that the relation may contain an extra term that does not appear above.

## Solution

To simplify the notation, we define  $n_j = \sum_{i=1}^n \gamma_{ij}$ .

We then have:

$$\begin{aligned}
 V(X) &= \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K \gamma_{ij} \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 = \frac{1}{n} \sum_{j=1}^K \sum_{i=1}^n \gamma_{ij} \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 \\
 &= \frac{1}{n} \sum_{j=1}^K \sum_{i=1}^n \gamma_{ij} \|\mathbf{x}_i - \boldsymbol{\mu}_j + \boldsymbol{\mu}_j - \bar{\mathbf{x}}\|^2 \tag{2} \\
 &= \frac{1}{n} \sum_{j=1}^K \sum_{i=1}^n \gamma_{ij} \left( \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2 + \|\boldsymbol{\mu}_j - \bar{\mathbf{x}}\|^2 + 2(\mathbf{x}_i - \boldsymbol{\mu}_j) \cdot (\boldsymbol{\mu}_j - \bar{\mathbf{x}}) \right) \\
 &= \sum_{j=1}^K \frac{n_j}{n} \frac{\sum_{i=1}^n \gamma_{ij} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2}{n_j} + \sum_{j=1}^K \frac{n_j \|\boldsymbol{\mu}_j - \bar{\mathbf{x}}\|^2}{n} + \frac{2}{n} \sum_{j=1}^K (\boldsymbol{\mu}_j - \bar{\mathbf{x}}) \cdot \sum_{i=1}^n \gamma_{ij} (\mathbf{x}_i - \boldsymbol{\mu}_j) \\
 &= \sum_{j=1}^K \frac{n_j}{n} V_j(X) + \tilde{V}(X) - \frac{2}{n} \sum_{j=1}^K n_j (\boldsymbol{\mu}_j - \bar{\mathbf{x}}) \cdot (\boldsymbol{\mu}_j - \bar{\boldsymbol{\mu}}_j), \text{ where } \bar{\boldsymbol{\mu}}_j \stackrel{\text{not.}}{=} \frac{\sum_{i=1}^n \gamma_{ij} \mathbf{x}_i}{n_j}. \tag{3}
 \end{aligned}$$

## Notes

1. In (2), the quantities  $\mu_j$  can be taken arbitrarily;  
however, they will be thought of in the end — see relation (3) and the Conclusion on the last slide — as the centroids of the clusters  $1, \dots, K$ , as computed at Step 1 in the  $K$ -means algorithm.
2. The equality (3) on the previous slide holds because

$$\begin{aligned}
 \sum_{i=1}^n \gamma_{ij}(x_i - \mu_j) &= \left( \sum_{i=1}^n \gamma_{ij} x_i \right) - n_j \mu_j = n_j \frac{\sum_{i=1}^n \gamma_{ij} x_i}{n_j} - n_j \mu_j \\
 &= n_j \left( \frac{\sum_{i=1}^n \gamma_{ij} x_i}{n_j} - \mu_j \right) = n_j (\bar{\mu}_j - \mu_j)
 \end{aligned}$$

## Conclusion

We already know — see CMU, 2009 spring, Ziv Bar-Joseph, HW5, pr. 2.1 — that  $K$ -means aims to minimize  $J$ , and consequently  $\frac{1}{n}J$ , which coincides with the first term in the expression we obtained for  $V(X)$ , namely  $\sum_{j=1}^K \frac{n_j}{n} V_j(X)$ .

Since the total variation  $V(X)$  is constant, minimizing the first term is equivalent to maximizing the sum of the other two terms, which is expected to be dominated by the between-cluster variation  $\tilde{V}(X)$  since a good  $\mu_j$  should be close to  $\bar{\mu}_j$ , making the third term small in absolute value.



Exemplifying the application of a simple version of  
EM/GMM

on data from  $\mathbb{R}$

$(\sigma_1 = \sigma_2 = 1, \pi_1 = \pi_2 = 1/2)$

CMU, 2012 spring, Ziv Bar-Joseph, final exam, pr. 3.1  
enhanced by Liviu Ciortuz

Suppose a GMM has two components with known variance and an equal prior distribution

$$\frac{1}{2}N(\mu_1, 1) + \frac{1}{2}N(\mu_2, 1).$$

The observed data are  $x_1 = 0.5$  and  $x_2 = 2$ , and the current estimates of  $\mu_1$  and  $\mu_2$  are 1 and 2 respectively.

a. Execute the first iteration of the EM algorithm.

*Hint:* Normal densities for the standardized variable  $y_{(\mu=0, \sigma=1)}$  at 0, 0.5, 1, 1.5, 2 are 0.4, 0.35, 0.24, 0.13, 0.05 respectively.

b. Consider the log-likelihood function for the “observable” data,

$$\ell(\mu_1, \mu_2) \stackrel{\text{def.}}{=} \ln P(x_1, x_2 | \mu_1, \mu_2) \stackrel{\text{indep.}}{=} \sum_{i=1}^2 \ln P(x_i | \mu_1, \mu_2) = \sum_{i=1}^2 \ln \left( \sum_{z_{ij}} P(x_i, z_{ij} | \mu_1, \mu_2) \right),$$

where  $z_{ij} \in \{0, 1\}$  and  $\sum_{j=1}^2 z_{ij} = 1$  for all  $i \in \{1, 2\}$ .

Compute the values of  $\ell$  function at the beginning and also at the end of the first iteration of the EM algorithm.

What do you see?

## Solution (a.)

58.

### The E-step:

$$\begin{aligned}
 E[Z_{i1}] &= P(Z_{i1} = 1|x_i, \mu) \stackrel{B.Th.}{=} \frac{P(x_i|Z_{i1} = 1, \mu_1)P(Z_{i1} = 1)}{P(x_i|Z_{i1} = 1, \mu_1)P(Z_{i1} = 1) + P(x_i|Z_{i2} = 1, \mu_2)P(Z_{i2} = 1)} \\
 &= \frac{P(x_i|Z_{i1} = 1, \mu_1) \cdot \frac{1}{2}}{P(x_i|Z_{i1} = 1, \mu_1) \cdot \frac{1}{2} + P(x_i|Z_{i2} = 1, \mu_2) \cdot \frac{1}{2}} \\
 &= \frac{P(x_i|Z_{i1} = 1, \mu_1)}{P(x_i|Z_{i1} = 1, \mu_1) + P(x_i|Z_{i2} = 1, \mu_2)} \text{ for } i \in \{1, 2\}.
 \end{aligned}$$

Therefore,

$$\begin{aligned}
 P(Z_{11} = 1|x_1, \mu) &= \frac{N(0.5; 1, 1)}{N(0.5; 1, 1) + N(0.5; 2, 1)} = \frac{N(0.5; 0, 1)}{N(0.5; 0, 1) + N(1.5; 0, 1)} = \frac{0.35}{0.35 + 0.13} = \frac{35}{48} \\
 P(Z_{21} = 1|x_2, \mu) &= \frac{N(2; 1, 1)}{N(2; 1, 1) + N(2; 2, 1)} = \frac{N(1; 0, 1)}{N(1; 0, 1) + N(0; 0, 1)} = \frac{0.24}{0.24 + 0.4} = \frac{0.24}{0.64} = \frac{3}{8}
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 P(Z_{12} = 1|x_1, \mu) &= P(Z_{11} = 0|x_1, \mu) = 1 - P(Z_{11} = 1|x_1, \mu) = \frac{13}{48} \\
 P(Z_{22} = 1|x_2, \mu) &= P(Z_{21} = 0|x_2, \mu) = 1 - P(Z_{21} = 1|x_2, \mu) = \frac{5}{8}
 \end{aligned}$$

The M-step:

$$\mu_j^{(t+1)} = \frac{\sum_{i=1}^2 E[Z_{ij}] x_i}{\sum_{i=1}^2 E[Z_{ij}]} = \frac{\sum_{i=1}^2 P(Z_{ij} = 1|x_i, \mu^{(t)}) x_i}{\sum_{i=1}^2 P(Z_{ij} = 1|x_i, \mu^{(t)})}$$

Therefore,

$$\mu_1^{(1)} = \frac{\frac{35}{48} \cdot 0.5 + \frac{3}{8} \cdot 2}{\frac{35}{48} + \frac{3}{8}} = \frac{107}{106} \approx 1.009 \quad \text{and} \quad \mu_2^{(1)} = \frac{\frac{13}{48} \cdot 0.5 + \frac{5}{8} \cdot 2}{\frac{13}{48} + \frac{5}{8}} = \frac{133}{86} \approx 1.54$$

### Solution (b.)

$$\begin{aligned}
 \ell(\mu_1, \mu_2) &\stackrel{def.}{=} \sum_{i=1}^2 \ln P(x_i | \mu_1, \mu_2) = \sum_{i=1}^2 \ln \left( \sum_{z_{ij}} P(x_i, z_{ij} | \mu_1, \mu_2) \right) \\
 &= \sum_{i=1}^2 \ln \left( \sum_{z_{ij}} P(x_i | z_{ij}, \mu_1, \mu_2) \cdot \underbrace{P(z_{ij} | \mu_1, \mu_2)}_{1/2} \right) \\
 &= \sum_{i=1}^2 \ln \left( \frac{1}{2} \sum_{z_{ij}} P(x_i | z_{ij}, \mu_1, \mu_2) \right) = \sum_{i=1}^2 \left[ -\ln 2 + \ln \left( \sum_{j=1}^2 P(x_i | z_{ij} = 1, \mu_j) \right) \right]
 \end{aligned}$$

$z_{ij}$	$P(x_i   z_{ij}, \mu_1, \mu_2)$
$z_{11} = 1, z_{12} = 0$	$\frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}(x_1 - \mu_1)^2)$
$z_{11} = 0, z_{21} = 1$	$\frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}(x_1 - \mu_2)^2)$
$z_{21} = 1, z_{22} = 0$	$\frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}(x_2 - \mu_1)^2)$
$z_{21} = 0, z_{22} = 1$	$\frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}(x_2 - \mu_2)^2)$

Therefore,

$$\begin{aligned}
 \ell(\mu_1, \mu_2) &= -2 \ln 2 + \ln \left( \frac{1}{\sqrt{2\pi}} \left( \exp\left(-\frac{1}{2}(x_1 - \mu_1)^2\right) + \exp\left(-\frac{1}{2}(x_1 - \mu_2)^2\right) \right) \right) \\
 &\quad + \ln \left( \frac{1}{\sqrt{2\pi}} \left( \exp\left(-\frac{1}{2}(x_2 - \mu_1)^2\right) + \exp\left(-\frac{1}{2}(x_2 - \mu_2)^2\right) \right) \right) \\
 &= -2 \ln 2 - \ln(2\pi) + \ln \left( \exp\left(-\frac{1}{2}(x_1 - \mu_1)^2\right) + \exp\left(-\frac{1}{2}(x_1 - \mu_2)^2\right) \right) \\
 &\quad + \ln \left( \exp\left(-\frac{1}{2}(x_2 - \mu_1)^2\right) + \exp\left(-\frac{1}{2}(x_2 - \mu_2)^2\right) \right)
 \end{aligned}$$

and

$$\ell(\mu_1^{(0)}, \mu_2^{(0)}) = \ell(1, 2) = -2.561833 \leq \ell(\mu_1^{(1)}, \mu_2^{(1)}) = \ell\left(\frac{107}{106}, \frac{133}{86}\right) = -2.462877,$$

meaning that the value of the log-likelihood function increases, which is in line with the theoretical result concerning the correctness of the EM algorithm.

Derivation of the EM algorithm for  
a mixture of  $K$  uni-variate Gaussians:  
the general case (i.e., when all parameters  $\pi, \mu, \sigma^2$  are free)

following Dahua Lin,  
*An Introduction to Expectation-Maximization*  
(MIT, ML 6768 course, 2012 fall)

## Note

We will first consider  $K = 2$ .

Generalization to  $K > 2$  will be shown afterwards.



### Estimation (E) Step:

$$\begin{aligned}
 p_{ij} &\stackrel{not.}{=} P(Z_{ij} = 1 \mid X_i, \mu, \sigma, \pi) \stackrel{calcul}{=} E[Z_{ij} \mid X_i, \mu, \sigma, \pi] \\
 &= \frac{P(X_i = x_i \mid Z_{ij} = 1, \mu, \sigma, \pi) \cdot P(Z_{ij} = 1 \mid \mu, \sigma, \pi)}{\sum_{j'=1}^2 P(X_i = x_i \mid Z_{ij'} = 1, \mu, \sigma, \pi) \cdot P(Z_{ij'} = 1 \mid \mu, \sigma, \pi)} \\
 &= \frac{\frac{1}{\sqrt{2\pi}\sigma_j} \cdot \exp\left(-\frac{(x_i - \mu_j)^2}{2\sigma_j^2}\right) \cdot \pi_j}{\frac{1}{\sqrt{2\pi}\sigma_1} \cdot \exp\left(-\frac{(x_i - \mu_1)^2}{2\sigma_1^2}\right) \cdot \pi_1 + \frac{1}{\sqrt{2\pi}\sigma_2} \cdot \exp\left(-\frac{(x_i - \mu_2)^2}{2\sigma_2^2}\right) \cdot \pi_2}
 \end{aligned}$$

Therefore, for  $t > 0$  we will have:

$$p_{ij}^{(t)} = \frac{\frac{\pi_j^{(t-1)}}{\sigma_j^{(t-1)}} \cdot \exp\left(-\frac{(x_i - \mu_j^{(t-1)})^2}{2(\sigma_j^{(t-1)})^2}\right)}{\frac{\pi_1^{(t-1)}}{\sigma_1^{(t-1)}} \cdot \exp\left(-\frac{(x_i - \mu_1^{(t-1)})^2}{2(\sigma_1^{(t-1)})^2}\right) + \frac{\pi_2^{(t-1)}}{\sigma_2^{(t-1)}} \cdot \exp\left(-\frac{(x_i - \mu_2^{(t-1)})^2}{2(\sigma_2^{(t-1)})^2}\right)}$$

The likelihood of a “complete” instance  $(x_i, z_{i1}, z_{i2})$ :

$$\begin{aligned}
 & P(X_i = x_i, Z_{i1} = z_{i1}, Z_{i2} = z_{i2} \mid \mu, \sigma, \pi) \\
 &= P(X_i = x_i \mid Z_{i1} = z_{i1}, Z_{i2} = z_{i2}, \mu_i, \sigma_i, \pi_i) \cdot P(Z_{i1} = z_{i1}, Z_{i2} = z_{i2} \mid \mu_i, \sigma_i, \pi_i) \\
 &= \frac{1}{\sqrt{2\pi}\sigma_j} \cdot \exp\left(-\frac{(x_i - \mu_j)^2}{2\sigma_j^2}\right) \cdot \pi_j, \text{ where } z_{ij} = 1 \text{ and } z_{ij'} = 0 \text{ for } j' \neq j \\
 &= \frac{1}{\sqrt{2\pi} \sigma_1^{z_{i1}} \sigma_2^{z_{i2}}} \cdot \exp\left(-\frac{1}{2} \sum_{j \in \{1,2\}} z_{ij} \frac{(x_i - \mu_j)^2}{\sigma_j^2}\right) \cdot \pi_1^{z_{i1}} \pi_2^{z_{i2}}
 \end{aligned}$$

The log-likelihood of the same “complete” instance will be:

$$\begin{aligned}
 & \ln P(X_i = x_i, Z_{i1} = z_{i1}, Z_{i2} = z_{i2} \mid \mu, \sigma, \pi) \\
 &= -\frac{1}{2} \ln(2\pi) - \sum_{j=1}^2 z_{ij} \ln \sigma_j - \frac{1}{2} \sum_{j=1}^2 z_{ij} \frac{(x_i - \mu_j)^2}{\sigma_j^2} + \sum_{j=1}^2 z_{ij} \ln \pi_j
 \end{aligned}$$

Given the dataset  $X = \{x_1, \dots, x_n\}$ , the log-likelihood function will be:

$$\begin{aligned}
 l(\mu, \sigma, \pi) &\stackrel{\text{def.}}{=} \ln P(X, Z_1, Z_2 \mid \mu, \sigma, \pi) \stackrel{i.i.d.}{=} \ln \prod_{i=1}^n P(X_i = x_i, Z_{i1}, Z_{i2} \mid \mu, \sigma, \pi) \\
 &= \sum_{i=1}^n \ln P(X_i = x_i, Z_{i1}, Z_{i2} \mid \mu, \sigma, \pi) \\
 &= -\frac{n}{2} \ln(2\pi) - \sum_{i=1}^n \sum_{j=1}^2 Z_{ij} \ln \sigma_j - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^2 Z_{ij} \frac{(x_i - \mu_j)^2}{\sigma_j^2} + \sum_{i=1}^n \sum_{j=1}^2 Z_{ij} \ln \pi_j
 \end{aligned}$$

The expectation of the log-likelihood function:

$$E[\ln P(X, Z_1, Z_2 \mid \mu, \sigma, \pi)] =$$

$$-\frac{n}{2} \ln(2\pi) - \sum_{i=1}^n \sum_{j=1}^2 E[Z_{ij}] \ln \sigma_j - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^2 E[Z_{ij}] \frac{(x_i - \mu_j)^2}{\sigma_j^2} + \sum_{i=1}^n \sum_{j=1}^2 E[Z_{ij}] \ln \pi_j$$

Here above, the probability function w.r.t. which the expectation was computed was left unspecified. Now we will make it explicit:

$$Q(\mu, \sigma, \pi \mid \mu^{(t)}, \sigma^{(t)}, \pi^{(t)}) \stackrel{not.}{=} E_{Z \mid X, \mu^{(t)}, \sigma^{(t)}, \pi^{(t)}}[\ln P(X, Z_1, Z_2 \mid X, \mu, \sigma, \pi)]$$

$$= -\frac{n}{2} \ln(2\pi) - \sum_{i=1}^n \sum_{j=1}^2 E[Z_{ij} \mid X_i, \mu^{(t)}, \sigma^{(t)}, \pi^{(t)}] \ln \sigma_j$$

$$- \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^2 E[Z_{ij} \mid X_i, \mu^{(t)}, \sigma^{(t)}, \pi^{(t)}] \frac{(x_i - \mu_j)^2}{\sigma_j^2}$$

$$+ \sum_{i=1}^n \sum_{j=1}^2 E[Z_{ij} \mid X_i, \mu^{(t)}, \sigma^{(t)}, \pi^{(t)}] \ln \pi_j$$

$$p_{ij}^{(t)} \stackrel{not.}{=} E[Z_{ij} \mid X, \mu^{(t-1)}, \sigma^{(t-1)}, \pi^{(t-1)}] \Rightarrow$$

$$Q(\mu, \sigma, \pi \mid \mu^{(t)}, \sigma^{(t)}, \pi^{(t)}) =$$

$$-\frac{n}{2} \ln(2\pi) - \sum_{i=1}^n \sum_{j=1}^2 p_{ij}^{(t)} \ln \sigma_j - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^2 p_{ij}^{(t)} \frac{(x_i - \mu_j)^2}{\sigma_j^2} + \sum_{i=1}^n \sum_{j=1}^2 p_{ij}^{(t)} \ln \pi_j$$

Since  $K = 2$  and  $\pi_1 + \pi_2 = 1$ , we get

$$Q(\mu, \sigma, \pi \mid \mu^{(t)}, \sigma^{(t)}, \pi^{(t)}) =$$

$$-\frac{n}{2} \ln 2\pi - \sum_{i=1}^n \sum_{j=1}^2 p_{ij}^{(t)} \ln \sigma_j - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^2 p_{ij}^{(t)} \frac{(x_i - \mu_j)^2}{\sigma_j^2} + \sum_{i=1}^n (p_{i1}^{(t)} \ln \pi_1 + p_{i2}^{(t)} \ln(1 - \pi_1))$$

**Maximization (M) Step:**

**[For  $K = 2$ :]**

$$\begin{aligned} \frac{\partial}{\partial \pi_1} Q(\mu, \sigma, \pi \mid \mu^{(t)}, \sigma^{(t)}, \pi^{(t)}) = 0 &\Leftrightarrow \frac{1}{\pi_1} \sum_{i=1}^n p_{i1}^{(t)} = \frac{1}{1 - \pi_1} \sum_{i=1}^n p_{i2}^{(t)} \Leftrightarrow \\ \sum_{i=1}^n p_{i1}^{(t)} = \pi_1 \left( \sum_{i=1}^n p_{i1}^{(t)} + \sum_{i=1}^n p_{i2}^{(t)} \right) &\Leftrightarrow \sum_{i=1}^n p_{i1}^{(t)} = \pi_1 \sum_{i=1}^n \underbrace{(p_{i1}^{(t)} + p_{i2}^{(t)})}_1 \Leftrightarrow \sum_{i=1}^n p_{i1}^{(t)} = n\pi_1 \\ \Rightarrow \pi_1^{(t+1)} &\leftarrow \frac{1}{n} \sum_{i=1}^n p_{i1}^{(t)} \end{aligned}$$

**Taking into account that  $\pi_1^{(t+1)} + \pi_2^{(t+1)} = 1$  and  $p_{i1}^{(t)} + p_{i2}^{(t)} = 1$  for  $i = 1, \dots, n$ ,**

$$\Rightarrow \pi_2^{(t+1)} \leftarrow \frac{1}{n} \sum_{i=1}^n p_{i2}^{(t)}$$

$$\frac{\partial}{\partial \mu_1} Q(\mu, \sigma, \pi \mid \mu^{(t)}, \sigma^{(t)}, \pi^{(t)}) = 0 \Leftrightarrow \frac{1}{\sigma_1^2} \sum_{i=1}^n p_{i1}^{(t)} (x_i - \mu_1) = 0 \Leftrightarrow \sum_{i=1}^n p_{i1}^{(t)} (x_i - \mu_1) = 0$$

$$\Rightarrow \mu_1^{(t+1)} \leftarrow \frac{\sum_{i=1}^n p_{i1}^{(t)} x_i}{\sum_{i=1}^n p_{i1}^{(t)}}$$

**Similarly,**  $\mu_2^{(t+1)} \leftarrow \frac{\sum_{i=1}^n p_{i2}^{(t)} x_i}{\sum_{i=1}^n p_{i2}^{(t)}}$

$$\frac{\partial}{\partial \sigma_1} Q(\mu, \sigma, \pi \mid \mu^{(t)}, \sigma^{(t)}, \pi^{(t)}) = 0 \Leftrightarrow -\frac{1}{\sigma_1} \sum_{i=1}^n p_{i1}^{(t)} + \frac{1}{\sigma_1^3} \sum_{i=1}^n p_{i1}^{(t)} (x_i - \mu_1)^2 = 0,$$

$$\Rightarrow \left( \sigma_1^{(t+1)} \right)^2 \leftarrow \frac{\sum_{i=1}^n p_{i1}^{(t)} (x_i - \mu_1^{(t+1)})^2}{\sum_{i=1}^n p_{i1}^{(t)}}$$

$$\text{Similarly, } \left( \sigma_2^{(t+1)} \right)^2 \leftarrow \frac{\sum_{i=1}^n p_{i2}^{(t)} (x_i - \mu_2^{(t+1)})^2}{\sum_{i=1}^n p_{i2}^{(t)}}$$

**Note:** One could relatively easily prove that these solutions (namely,  $\pi^{(t+1)}, \mu^{(t+1)}, \sigma^{(t+1)}$ ) of the partial derivatives of the *auxiliary function*  $Q$  designate the values for which  $Q$  reaches its *maximum*.



## Generalization to $K > 2$

In this case, the Bernoulli distribution is replaced by a categorical one. The only one change needed in the above proof concerns updating the parameters of this distribution.

Since  $\pi_1 + \dots + \pi_K = 1$ , we must solve the following *constraint optimization problem*:

$$\begin{aligned} & \max_{\pi, \mu, \sigma} Q(\pi, \mu, \sigma | \pi^{(t)}, \mu^{(t)}, \sigma^{(t)}) \\ & \text{subject to } \sum_{i=1}^K \pi_j = 1 \text{ and } \pi_j \geq 0, \forall j = 1, \dots, K. \end{aligned}$$

By letting asside the  $\geq$  constraints, and using the *Lagrangian multiplier*  $\lambda \in \mathbb{R}$ , this problem becomes:

$$\max_{\pi, \mu, \sigma} \left( Q(\pi, \mu, \sigma | \pi^{(t)}, \mu^{(t)}, \sigma^{(t)}) + \lambda \left( 1 - \sum_{i=1}^K \pi_j \right) \right).$$

**For  $j = 1, \dots, K$ :**

$$\frac{\partial}{\partial \pi_j} Q(\mu, \sigma, \pi \mid \mu^{(t)}, \sigma^{(t)}, \pi^{(t)}) = 0 \Leftrightarrow \sum_{i=1}^n p_{ij}^{(t)} \frac{1}{\pi_j} = \lambda \Leftrightarrow \pi_j^{(t+1)} = \frac{1}{\lambda} \sum_{i=1}^n p_{ij}^{(t)}.$$

**Because  $\sum_{j=1}^K \pi_j^{(t+1)} = 1$ , it follows that**

$$\lambda = \sum_{j=1}^K \sum_{i=1}^n \pi_j^{(t+1)} = \sum_{i=1}^n \underbrace{\sum_{j=1}^K \pi_j^{(t+1)}}_1 = \sum_{i=1}^n 1 = n.$$

**Therefore,**

$$\pi_j^{(t+1)} \leftarrow \frac{1}{n} \sum_{i=1}^n p_{ij}^{(t)}.$$

***Note*** that indeed  $\pi_j^{(t+1)} \geq 0$ , because the  $p_{ij}^{(t)}$  terms designate some probabilities (see E-step).

## To summarize:

**E Step:**

$$p_{ij}^{(t)} \stackrel{\text{not.}}{=} P(z_{ij} = 1 \mid x_i; \mu^{(t)}, (\sigma^2)^{(t)}, \pi^{(t)}) = \frac{N(x_i \mid \mu_j^{(t)}, (\sigma_j^2)^{(t)}) \cdot \pi_j^{(t)}}{\sum_{l=1}^K N(x_i \mid \mu_l^{(t)}, (\sigma_l^2)^{(t)}) \cdot \pi_l^{(t)}}$$

**where**  $N(x_i \mid \mu_j, \sigma_j^2) \stackrel{\text{def.}}{=} \frac{1}{\sqrt{2\pi}\sigma_j} \cdot \exp\left(-\frac{(x_i - \mu_j)^2}{2\sigma_j^2}\right).$

**M Step:**

$$\pi_j^{(t+1)} \leftarrow \frac{1}{n} \sum_{i=1}^n p_{ij}^{(t)}$$

$$\mu_j^{(t+1)} \leftarrow \frac{\sum_{i=1}^n p_{ij}^{(t)} x_i}{\sum_{i=1}^n p_{ij}^{(t)}}$$

$$\left(\sigma_j^{(t+1)}\right)^2 \leftarrow \frac{\sum_{i=1}^n p_{ij}^{(t)} (x_i - \mu_j^{(t+1)})^2}{\sum_{i=1}^n p_{ij}^{(t)}}$$

## Example: Modelling the waiting and eruption times for the Old Faithful geyser, (Yellowstone Park, USA)

Michael Eichler (University of Chicago, Statistics course (24600) - Spring 2004)

### R code

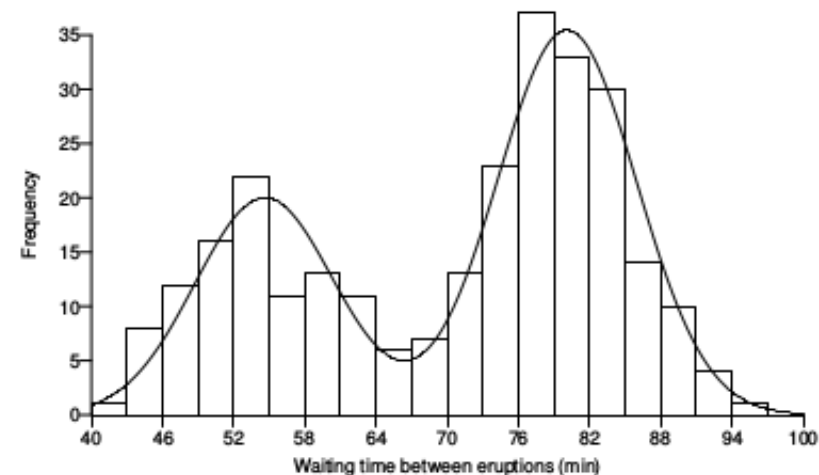
```
p<-c(0.5, 40, 90, 20, 20)
emstep<-function(Y, p) {
  EZ<-p[1] * dnorm(Y, p[2], sqrt(p[4]))/
    (p[1] * dnorm(Y, p[2], sqrt(p[4]))
    +(1 - p[1]) * dnorm(Y, p[3], sqrt(p[5]))))
  p[1]<-mean(EZ)
  p[2]<-sum(EZ * Y) / sum(EZ)
  p[3]<-sum((1 - EZ)*Y) / sum(1-EZ)
  p[4]<-sum(EZ * (Y - p[2])^2) / sum(EZ)
  p[5]<-sum((1 - EZ) * (Y - p[3])^2) / sum(1 - EZ)
  p
}
emiteration<-function(Y, p, n=10) {
  for (i in (1:n)) {
    p<-emstep(Y, p)
  }
  p
}
p<-c(0.5, 40, 90, 20, 20)
p<-emiteration(Y, p, 20)
p
p<-emstep(Y, p)
p
}
```

### Starting values:

$$p^{(0)} = 0.4$$

$$\mu_1^{(0)} = 40, \sigma_1^{(0)} = 4$$

$$\mu_1^{(0)} = 90, \sigma_1^{(0)} = 4$$



$k$	$p^{(k)}$	$\mu_1^{(k)}$	$\mu_2^{(k)}$	$\sigma_1^{(k)}$	$\sigma_2^{(k)}$
1	0.3508	54.22	79.91	5.465	5.999
2	0.3539	54.38	79.94	5.671	6.013
3	0.3562	54.46	79.99	5.744	5.969
4	0.3578	54.51	80.02	5.787	5.935
5	0.3588	54.55	80.05	5.815	5.912
6	0.3595	54.57	80.06	5.834	5.897
7	0.3600	54.59	80.07	5.846	5.887
8	0.3603	54.60	80.08	5.855	5.880
9	0.3605	54.60	80.08	5.860	5.876
10	0.3606	54.61	80.09	5.864	5.873
11	0.3607	54.61	80.09	5.866	5.871
12	0.3608	54.61	80.09	5.868	5.870
13	0.3608	54.61	80.09	5.869	5.869
14	0.3608	54.61	80.09	5.870	5.869
15	0.3609	54.61	80.09	5.870	5.868
20	0.3609	54.61	80.09	5.871	5.868
25	0.3609	54.61	80.09	5.871	5.868

Exemplifying

some **methodological issues** regarding the application of  
the **EM algorithmic schema**

(using a simple EM/GMM algorithm on data from  $\mathbb{R}$  ( $\pi_1 = \pi_2 = 1/2$ ))

CMU, 2007 spring, Eric Xing, final exam, pr. 1.8

A long time ago there was a village amidst hundreds of lakes. Two types of fish lived in the region, but only one type in each lake.

These types of fish both looked exactly the same, smelled exactly the same when cooked, and had the exact same delicious taste – except one was poisonous and would kill any villager who ate it. The only other difference between the fish was their effect on the pH (acidity) of the lake they occupy.

The pH for lakes occupied by the non-poisonous type of fish was distributed according to a Gaussian with unknown mean ( $\mu_{safe}$ ) and variance ( $\sigma_{safe}^2$ ) and the pH for lakes occupied by the poisonous type was distributed according to a different Gaussian with unknown mean ( $\mu_{deadly}$ ) and variance ( $\sigma_{deadly}^2$ ). (Poisonous fish tended to cause slightly more acidic conditions).

Naturally, the villagers turned to machine learning for help. However, there was much debate about the right way to apply EM to their problem. For each of the following procedures, indicate whether it is an accurate implementation of Expectation-Maximization and will provide a reasonable estimate for parameters  $\mu$  and  $\sigma^2$  for each class.

a.

Guess initial values of  $\mu$  and  $\sigma^2$  for each class.

(1) For each lake, find the most likely class of fish for the lake.

(2) Update the  $\mu$  and  $\sigma^2$  values using their maximum likelihood estimates based on these predictions.

Iterate (1) and (2) until convergence.

b.

For each lake, guess an initial probability that it is safe.

(1) Using these probabilities, find the maximum likelihood estimates for the  $\mu$  and  $\sigma$  values for each class.

(2) Use these estimates of  $\mu$  and  $\sigma$  to reestimate lake safety probabilities.

Iterate (1) and (2) until convergence.

c.

Compute the mean and variance of the pH levels across all lakes.

Use these values for the  $\mu$  and  $\sigma^2$  value of each class of fish.

(1) Use the  $\mu$  and  $\sigma^2$  values of each class to compute the belief that each lake contains poisonous fish.

(2) Find the maximum likelihood values for  $\mu$  and  $\sigma^2$ .

Iterate (1) and (2) until convergence.



## Solution

- a. It'll do ok if we give sensible enough  $\mu$  and  $\sigma^2$  initial values.
- b. Ok, this is the same as *a* after the first M-step.  
(See the general EM algorithmic schema on the next slide.)
- c. This will be stuck at the initial  $\mu$  and  $\sigma^2$ :

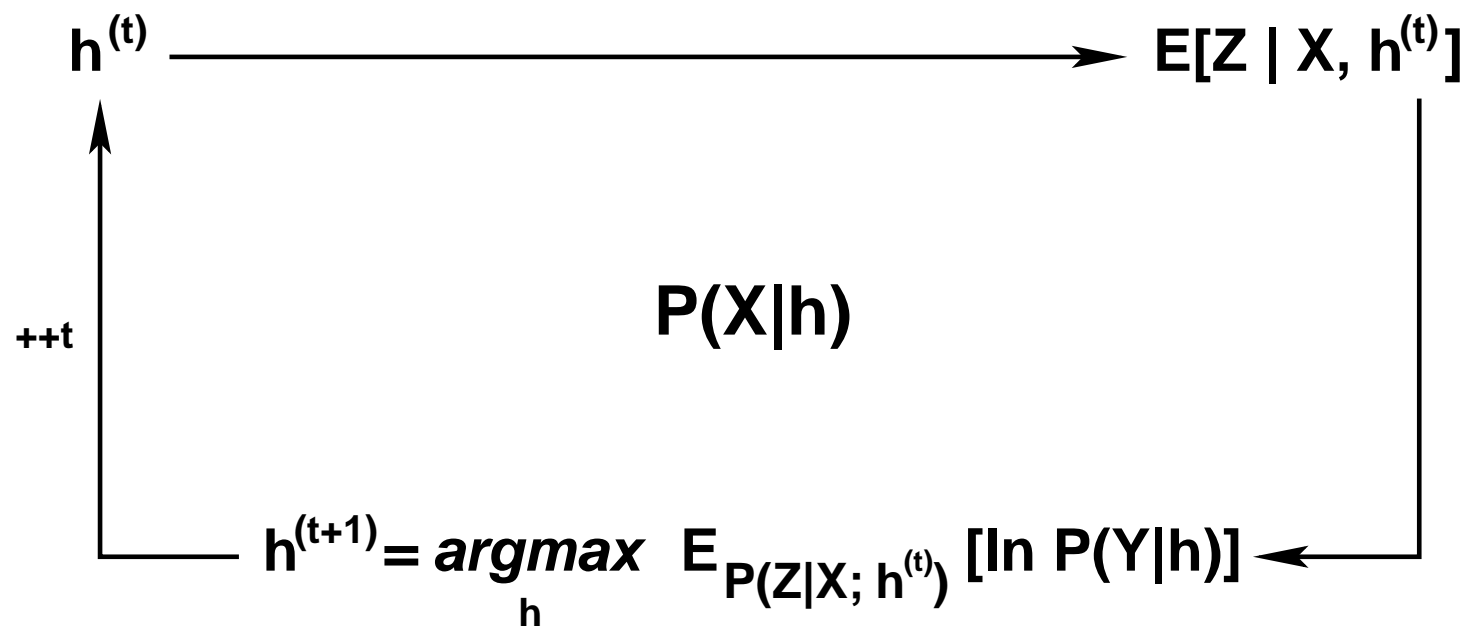
In the E-step we'll get:

$$P(\text{safe}|x) = \frac{P(x|\text{safe}) \cdot P(\text{safe})}{P(x|\text{safe}) \cdot P(\text{safe}) + P(x|\text{deadly}) \cdot P(\text{deadly})} = \frac{1}{2}$$

since on one side we assume  $P(\text{safe}) = P(\text{poison}) = \frac{1}{2}$  and on the other side  $P(x|\text{safe}) = P(x|\text{poison})$  because  $\mu_{\text{safe}} = \mu_{\text{deadly}}$  and  $\sigma_{\text{safe}}^2 = \sigma_{\text{deadly}}^2$ .

In the M-step  $\mu$  and  $\sigma^2$  will not change since we are again letting them be calculated from all lakes (weighted equally).

## The [general] EM algorithmic schema



The EM algorithm for modeling  
mixtures of multi-variate Gaussians  
with diagonal covariance matrices

Utah University, Piyush Rai  
ML course (CS5350/6350), 2009 fall,  
lecture notes, *Gaussian Mixture Models*

[adapted by Liviu Ciortuz]

Fie mixtura de distribuții gaussiene

$$gmm(x) = \sum_{j=1}^K \pi_j \mathcal{N}(x; \mu_j, \sigma^2 I),$$

unde

$$x \in \mathbb{R}^d,$$

probabilitățile a priori de selecție  $\pi_j \in \mathbb{R}$  satisfac (ca de obicei) restricțiile  $\pi_j \geq 0$  pentru  $j = 1, \dots, K$  și  $\sum_{j=1}^K \pi_j = 1$ ;

mediile gaussianelor  $j = 1, \dots, K$  sunt vectorii  $\mu_j \in \mathbb{R}^d$ , iar

matricele de covarianță ale acestor gaussiene sunt identice, ba chiar au forma particulară  $\sigma^2 I$ , cu  $\sigma \in \mathbb{R}$  și  $\sigma > 0$ , matricea  $I$  fiind matricea identitate.

Se consideră instanțele  $x_1, \dots, x_n \in \mathbb{R}^d$  generate cu distribuția probabilistă de mai sus ( $gmm$ ).

În acest exercițiu veți deduce *regulile de actualizare* din cadrul [pasului E și al pasului M al] algoritmului EM care face estimarea parametrilor  $\pi \stackrel{not.}{=} (\pi_1, \dots, \pi_K)$ ,  $\mu \stackrel{not.}{=} (\mu_1, \dots, \mu_K)$  și  $\sigma$ .

a. Se știe că expresia funcției de densitate a distribuției gaussiene multivariate ( $d$ -dimensionale) de medie  $\mu \in \mathbb{R}^d$  și matrice de covarianță  $\Sigma \in \mathbb{R}^{d \times d}$  este

$$\frac{1}{(\sqrt{2\pi})^d \sqrt{|\Sigma|}} \exp \left( -\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right),$$

unde  $x$  și  $\mu$  sunt considerați vectori-coloană din  $\mathbb{R}^d$ , iar operatorul  $\top$  desemnează operația de transpunere a vectorilor/matricelor.

Aduceți expresia de mai sus la forma cea mai simplă pentru cazul  $\Sigma = \sigma^2 I$ .

Vă recomandăm să folosiți faptul că  $\|x - \mu\|^2 = (x - \mu)^\top (x - \mu) = (x - \mu) \cdot (x - \mu)$ , unde operatorul  $\cdot$  desemnează produsul scalar al vectorilor.

*Observație:* Prima din ultimele două egalități implică un ușor *abuz* (sau, mai degrabă, o convenție) *de notație*: o matrice reală de dimensiune  $1 \times 1$  este identificată cu un număr real, care este chiar singurul ei element. Același tip de *abuz/convenție* a intervenit și în scrierea expresiei  $\exp(\dots)$  de mai sus.

## Answer

$$\Sigma = \sigma^2 I \Rightarrow \Sigma^{-1} = \frac{1}{\sigma^2} I,$$

and also

$$|\Sigma| = (\sigma^2)^d \Rightarrow \sqrt{|\Sigma|} = \sigma^d \text{ since } \sigma > 0.$$

Therefore,

$$\begin{aligned} \mathcal{N}(x; \mu, \Sigma = \sigma^2 I) &= \frac{1}{(\sqrt{2\pi})^d \sqrt{|\Sigma|}} \exp \left( -\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right) \\ &= \frac{1}{(\sqrt{2\pi})^d \sigma^d} \exp \left( -\frac{1}{2} (x - \mu)^\top \frac{1}{\sigma^2} (x - \mu) \right) \\ &= \frac{1}{(\sqrt{2\pi}\sigma)^d} \exp \left( -\frac{1}{2\sigma^2} \|x - \mu\|^2 \right) \end{aligned}$$

b. Vom asocia fiecărei instanțe  $x_i$  un vector-indicator (mai precis, un vector de variabile aleatoare  $z_i \in \{0, 1\}^K$ , cu  $z_{ij} = 1$  dacă și numai dacă  $x_i$  a fost generat de gaussiană  $\mathcal{N}(x; \mu_j, \sigma^2 I)$ ).

Pentru *pasul E* al algoritmului EM, veți demonstra mai întâi că media  $E[z_{ij}] \stackrel{\text{not.}}{=} E[z_{ij}|x_i; \pi, \mu, \sigma]$ , unde  $x_i = (x_{i,1}, \dots, x_{i,d}) \in \mathbb{R}^d$ ,  $\mu \stackrel{\text{not.}}{=} (\mu_1, \dots, \mu_K) \in (\mathbb{R}^d)^K$  și  $\sigma \in \mathbb{R}_+$ , are valoarea  $P(z_{ij} = 1|x_i; \pi, \mu, \sigma)$ , iar apoi veți elabora formula de calcul a acestei probabilități, folosind teorema lui Bayes.

## Answer

For the sake of simplicity, we will designate  $E[z_{ij}|x_i; \pi, \mu, \sigma]$  as  $E[z_{ij}]$ .  
So,

$$\begin{aligned}
 E[z_{ij}] &\stackrel{not.}{=} E[z_{ij}|x_i; \pi, \mu, \sigma] \stackrel{def.}{=} 0 \cdot P(z_{ij} = 0|x_i; \pi, \mu, \sigma) + 1 \cdot P(z_{ij} = 1|x_i; \pi, \mu, \sigma) \\
 &= P(z_{ij} = 1|x_i; \pi, \mu, \sigma) \\
 &\stackrel{Bayes F.}{=} \frac{P(x_i|z_{ij} = 1; \pi, \mu, \sigma) \cdot P(z_{ij} = 1; \pi, \mu, \sigma)}{P(x_i; \pi, \mu, \sigma)} \\
 &= \frac{P(x_i|z_{ij} = 1; \pi, \mu, \sigma) \cdot P(z_{ij} = 1; \pi, \mu, \sigma)}{\sum_{j'=1}^K P(x_i|z_{ij'} = 1; \pi, \mu, \sigma) P(z_{ij'} = 1; \pi, \mu, \sigma)} \\
 &\stackrel{a.}{=} \frac{\frac{1}{(\sqrt{2\pi}\sigma)^d} \exp\left(-\frac{1}{2\sigma^2}\|x_i - \mu_j\|^2\right) \pi_j}{\sum_{j'=1}^K \frac{1}{(\sqrt{2\pi}\sigma)^d} \exp\left(-\frac{1}{2\sigma^2}\|x_i - \mu_{j'}\|^2\right) \pi_{j'}} \\
 &= \frac{\pi_j \exp\left(-\frac{1}{2\sigma^2}\|x_i - \mu_j\|^2\right)}{\sum_{j'=1}^K \pi_{j'} \exp\left(-\frac{1}{2\sigma^2}\|x_i - \mu_{j'}\|^2\right)} \tag{4}
 \end{aligned}$$



c. Arătați că expresia funcției de log-verosimilitate a datelor „complete“ în raport cu parametrii  $\pi, \mu$  și  $\sigma$  este

$$\ln p(x, z | \pi, \mu, \sigma) = \sum_{i=1}^n \sum_{j=1}^K z_{ij} (\ln \pi_j + \ln \mathcal{N}(x_i; \mu_j, \sigma^2 I)),$$

unde  $x \stackrel{not.}{=} (x_1, \dots, x_n)$  și  $z \stackrel{not.}{=} (z_1, \dots, z_n)$ .

Deduceți apoi expresia „funcției auxiliare“

$$Q(\pi, \mu, \sigma | \pi^{(t)}, \mu^{(t)}, \sigma^{(t)}) \stackrel{def.}{=} E[\ln p(x, z | \pi, \mu, \sigma)],$$

cu precizarea că media aceasta este calculată în raport cu distribuția / distribuțiile  $P(z_{ij} | x_i, \pi^{(t)}, \mu^{(t)}, \sigma^{(t)})$ .

## Answer

The log-verosimilarity of the complete data is

$$\begin{aligned}
 \ln p(x, z | \pi, \mu, \sigma) &\stackrel{\text{def.}}{=} \ln p((x_1, z_1), \dots, (x_n, z_n) | \pi, \mu, \sigma) \stackrel{i.i.d.}{=} \ln \prod_{i=1}^n p(x_i, z_i | \pi, \mu, \sigma) \\
 &= \sum_{i=1}^n \ln p(x_i, z_i | \pi, \mu, \sigma) \stackrel{\text{mult. rule}}{=} \sum_{i=1}^n \ln p(x_i | z_i; \pi, \mu, \sigma) \cdot \underbrace{p(z_i | \pi, \mu, \sigma)}_{\pi_j} \\
 &= \sum_{i=1}^n \sum_{j=1}^K z_{ij} [\ln \mathcal{N}(x_i | \mu_j, \sigma) + \ln \pi_j],
 \end{aligned}$$

since  $z_i = (z_{i,1}, \dots, z_{i,j}, \dots, z_{i,K})$ , with  $z_{i,j} = 1$  and  $z_{i,j'} = 0$  for all  $j' \neq j$ .

Furthermore, using the result we got at part a, we can write

$$\ln p(x, z | \pi, \mu, \sigma) = \sum_{i=1}^n \sum_{j=1}^K z_{ij} \left[ -\frac{d}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|x_i - \mu_j\|^2 + \ln \pi_j \right].$$

Finally, using the linearity of expectation,

$$\begin{aligned}
 Q(\pi, \mu, \sigma | \pi^{(t)}, \mu^{(t)}, \sigma^{(t)}) &\stackrel{\text{def.}}{=} E[\ln p(x, z | \pi, \mu, \sigma)] \\
 &= \sum_{i=1}^n \sum_{j=1}^K E[z_{ij}] \left[ -\frac{d}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|x_i - \mu_j\|^2 + \ln \pi_j \right].
 \end{aligned}$$

Pentru *Pasul M*, în contextul precizat în enunț, vom avea de rezolvat următoarea *problemă de optimizare*:

$$(\pi^{(t+1)}, \mu^{(t+1)}, \sigma^{(t+1)}) = \operatorname{argmax}_{\pi, \mu, \sigma} Q(\pi, \mu, \sigma | \pi^{(t)}, \mu^{(t)}, \sigma^{(t)}), \quad (5)$$

$$\text{cu restricțiile } \pi_j^{(t+1)} \geq 0 \text{ pentru } j = 1, \dots, K \text{ și } \sum_{j=1}^K \pi_j^{(t+1)} = 1.$$

Această problemă se rezolvă optimizând funcția ei obiectiv în mod separat în raport cu variabilele  $\pi, \mu$  și  $\sigma$ .

d. Aplicați metoda multiplicatorilor lui Lagrange pentru a rezolva problema de optimizare cu restricții (5) în raport (doar) cu variabilele  $\pi$ .

## Answer

For now, we will ignore the constraints  $\pi_j \geq 0$ . Thus, the *Lagrangian functional* associated to our optimisation problem is:

$$(\pi^{(t+1)}, \mu^{(t+1)}, \sigma^{(t+1)}) = \underset{\pi, \mu, \sigma, \lambda}{\operatorname{argmax}} (Q(\pi, \mu, \sigma | \pi^{(t)}, \mu^{(t)}, \sigma^{(t)}) - \lambda(1 - \sum_{j=1}^K \pi_j)), \quad (6)$$

with  $\lambda \in \mathbb{R}$  playing the role of *Lagrangian multiplier*. (After solving this optimisation problem it will be seen that indeed  $\pi_j \geq 0$  for all  $j = 1, \dots, K$ .)

By taking the partial derivative of the Lagrangian functional with respect to  $\pi_j$  (with  $j \in \{1, \dots, K\}$ ) and then solving for it,

$$\begin{aligned} \frac{\partial}{\partial \pi_j} (Q(\pi, \mu, \sigma | \pi^{(t)}, \mu^{(t)}, \sigma^{(t)}) - \lambda(1 - \sum_{j=1}^K \pi_j)) &= 0 \Leftrightarrow \\ \frac{\partial}{\partial \pi_j} \left( \sum_{i=1}^n \sum_{j=1}^K E[z_{ij}] \left[ -\frac{d}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|x_i - \mu_j\|^2 + \ln \pi_j \right] \right) - \lambda &= 0 \Leftrightarrow \sum_{i=1}^n E[z_{ij}] \frac{1}{\pi_j} = \lambda \end{aligned}$$

we will obtain the solution

$$\pi_j^{(t+1)} = \frac{1}{\lambda} \sum_{i=1}^n E[z_{ij}]. \quad (7)$$

Moreover, since these solutions should satisfy the constraint  $\sum_{j=1}^K \pi_j^{(t+1)} = 1$ , it follows that

$$\begin{aligned} \sum_{j=1}^K \frac{1}{\lambda} \sum_{i=1}^n E[z_{ij}] = 1 &\Leftrightarrow \frac{1}{\lambda} \sum_{j=1}^K \sum_{i=1}^n E[z_{ij}] = 1 \stackrel{b}{\Leftrightarrow} \frac{1}{\lambda} \sum_{i=1}^n \sum_{j=1}^K P(z_{ij} = 1 | x_i, \pi, \mu, \sigma) = 1 \Leftrightarrow \\ &\underbrace{\frac{1}{\lambda} \sum_{i=1}^n \sum_{j=1}^K P(z_{ij} = 1 | x_i, \pi, \mu, \sigma)}_1 = 1 \Leftrightarrow \frac{1}{\lambda} \sum_{i=1}^n 1 = 1 \Leftrightarrow \lambda = \frac{1}{n}. \end{aligned}$$

By replacing this value of  $\lambda$  back into relation (7), we will get

$$\pi_j^{(t+1)} = \frac{1}{n} \sum_{i=1}^n E[z_{ij}]. \quad (8)$$

It is obvious that  $\pi_j^{(t+1)} \geq 0$ .

e. Optimizați funcția  $Q$  (i.e., rezolvați problema de optimizare (5)) în raport cu variabilele  $\mu$ .

**Answer:**

**Remember the notation**  $\mu = (\mu_1, \dots, \mu_K) \in (\mathbb{R}^d)^K$ , **with**  $\mu_j \in \mathbb{R}^d$  **for**  $j = 1, \dots, K$ .

$$\begin{aligned}
 & \nabla_{\mu_j} Q(\pi, \mu, \sigma | \pi^{(t)}, \mu^{(t)}, \sigma^{(t)}) \\
 &= \nabla_{\mu_j} \sum_{i=1}^n \sum_{j'=1}^K E[z_{ij'}] \left[ -\frac{d}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|x_i - \mu_{j'}'\|^2 + \ln \pi_{j'} \right] \\
 &= \sum_{i=1}^n E[z_{ij}] \left[ -\frac{1}{2\sigma^2} \nabla_{\mu_j} (x_i - \mu_j)^2 \right] = -\frac{1}{2\sigma^2} \sum_{i=1}^n E[z_{ij}] 2(x_i - \mu_j)(-1) \\
 &= \frac{1}{\sigma^2} \sum_{i=1}^n E[z_{ij}] (x_i - \mu_j) = \frac{1}{\sigma^2} \left[ \left( \sum_{i=1}^n E[z_{ij}] x_i \right) - \left( \sum_{i=1}^n E[z_{ij}] \right) \mu_j \right]
 \end{aligned}$$

After equating this expression to zero (in fact, the column-vector  $(0, \dots, 0)^\top \in \mathbb{R}^d$ ), we get the solution:

$$\mu^{(t+1)} = \frac{\sum_{i=1}^n E[z_{ij}] x_i}{\sum_{i=1}^n E[z_{ij}]} . \tag{9}$$

f. Optimizați funcția  $Q$  (i.e., rezolvați problema de optimizare (5)) în raport cu variabila  $\sigma$ .

**Answer:**

The partial derivative of  $Q(\pi, \mu, \sigma | \pi^{(t)}, \mu^{(t)}, \sigma^{(t)})$  with respect to  $\sigma$  is:

$$\begin{aligned}
 & \frac{\partial}{\partial \sigma} Q(\pi, \mu, \sigma | \pi^{(t)}, \mu^{(t)}, \sigma^{(t)}) \\
 &= \frac{\partial}{\partial \sigma} \sum_{i=1}^n \sum_{j=1}^K E[z_{ij}] \left[ -\frac{d}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|x_i - \mu_j\|^2 + \ln \pi_j \right] \\
 &= \sum_{i=1}^n \sum_{j=1}^K E[z_{ij}] \frac{\partial}{\partial \sigma} \left[ -\frac{d}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|x_i - \mu_j\|^2 \right] \\
 &= \sum_{i=1}^n \sum_{j=1}^K E[z_{ij}] \left[ -\frac{d}{2} \cdot \frac{2\sigma}{\sigma^2} - \frac{-2}{2\sigma^3} \|x_i - \mu_j\|^2 \right] = \sum_{i=1}^n \sum_{j=1}^K E[z_{ij}] \left[ -\frac{d}{\sigma} + \frac{1}{\sigma^3} \|x_i - \mu_j\|^2 \right] \\
 &= \frac{1}{\sigma^3} \left( -d\sigma^2 \underbrace{\sum_{i=1}^n \sum_{j=1}^K E[z_{ij}]}_1 + \sum_{i=1}^n \sum_{j=1}^K E[z_{ij}] \|x_i - \mu_j\|^2 \right) = \frac{1}{\sigma^3} \left( -nd\sigma^2 + \sum_{i=1}^n \sum_{j=1}^K E[z_{ij}] \|x_i - \mu_j\|^2 \right)
 \end{aligned}$$

By equating  $\frac{\partial}{\partial \sigma} Q(\pi, \mu, \sigma | \pi^{(t)}, \mu^{(t)}, \sigma^{(t)})$  to 0 we get the solution:

$$(\sigma^{(t+1)})^2 = \frac{1}{nd} \sum_{i=1}^n \sum_{j=1}^K E[z_{ij}] \|x_i - \mu_j^{(t+1)}\|^2 \geq 0. \quad (10)$$

It is not too difficult to see that this solution leads to the maximization of  $Q$ . (Note that  $\sigma > 0$  and the expression  $-nd\sigma^2 + \sum_{i=1}^n \sum_{j=1}^K E[z_{ij}] \|x_i - \mu_j\|^2$  as a function of  $\sigma^2$ , is linear and decreasing.)



g. Sumarizați rezultatele obținute la punctele de mai sus (*b* și *d-f*), redactând în pseudo-cod algoritmul EM pentru rezolvarea mixturii de gaussiene din enunț.

**Answer:**

By gathering the relations (4), (8), (9) and (10), we are now able to write the pseudo-code of our algorithm:

- Initialize the a priori probabilities  $\pi$ , the means  $\mu$ , and the variance  $\sigma$ .
- Iterate until a certain *termination condition* is met:

*Step E:* Compute the expectation (i.e., a posteriori probabilities) of  $z$  variables:

$$p_{ij}^{(t+1)} \stackrel{not.}{=} E[z_{ij}|x_i; \pi^{(t)}, \mu^{(t)}, \sigma^{(t)}] = \frac{\pi_j^{(t)} \mathcal{N}(x; \mu_j^{(t)}, (\sigma^{(t)})^2 I)}{\sum_{j'=1}^K \pi_{j'}^{(t)} \mathcal{N}(x; \mu_{j'}^{(t)}, (\sigma^{(t)})^2 I)}$$

*Step M:* Compute new values for  $\pi, \mu$  și  $\sigma$ :

$$\pi_j^{(t+1)} = \frac{1}{n} \sum_{i=1}^n p_{ij}^{(t)} \quad \mu_j^{(t+1)} = \frac{\sum_{i=1}^n p_{ij}^{(t)} x_i}{\sum_{i=1}^n p_{ij}^{(t)}} \quad \left(\sigma^{(t+1)}\right)^2 = \frac{1}{dn} \sum_{i=1}^n \sum_{j=1}^K p_{ij}^{(t)} \|x_i - \mu_j^{(t)}\|^2$$

## Important Remark: a slight *generalization*

Suppose that our mixture model is made of Gaussians with unrestricted diagonal covariance matrices, i.e.,

$$Z_i \sim \text{Categorical}(p_1, \dots, p_K)$$

$$X_i | Z_i = j \sim \mathcal{N} \left( \begin{bmatrix} \mu_{j,1} \\ \vdots \\ \mu_{j,d} \end{bmatrix}, \begin{bmatrix} (\sigma_{j,1})^2 & 0 & \dots & 0 \\ 0 & (\sigma_{j,2})^2 & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \dots & 0 & (\sigma_{j,d})^2 \end{bmatrix} \right)$$

It can be easily be proven (by going along the lines of the above proof) that the only *update relation* that changes in the formulation of the EM algorithm is (10). It becomes:

$$(\sigma_{jk}^{(t+1)})^2 = \frac{\sum_{i=1}^n E[z_{ij}] (x_{i,k} - \mu_{j,k}^{(t+1)})^2}{\sum_{i=1}^n E[z_{ij}]} \geq 0 \text{ for } j = 1, \dots, K \text{ and } k = 1, \dots, d. \quad (11)$$

This is indeed what we would expect, given on one hand the fact that the Gaussian components are mutually independent, and on the other hand the updating formulas [that we have already obtained] for the EM algorithm for solving mixtures of uni-variate Gaussians [of unrestricted form].

The EM algorithm for modeling  
mixtures of multi-variate Gaussians

Stanford University, Prof. Andrew Ng  
ML course, 2009, lecture notes, parts VIII and IX  
[adapted by Liviu Ciortuz]

Suppose that we are given the *instances*  $x_1, \dots, x_n \in \mathbb{R}^d$  (all seen as column-vectors). We wish to *model* these data by specifying a joint distribution  $p(x_i, z_i) = p(x_i|z_i) \cdot p(z_i)$ . Here,

$z_i \sim \text{Categorical}(\pi)$ ,

$K$  denotes the number of values that the  $z_i$ 's can take on, namely  $\pi_j \stackrel{\text{not.}}{=} p(z_i = j)$  for  $j = 1, \dots, K$ , with  $\sum_{j=1}^K \pi_j = 1$ , and the [conditional] distribution  $x_i|z_i = j$  is a Gaussian of mean vector  $\mu_j$  and covariance matrix  $\Sigma_j$ .

Thus, our model posits that each  $x_i$  was *generated* by randomly choosing  $z_i$  from  $\{1, \dots, K\}$ , and then  $x_i$  was drawn from one of the  $K$  Gaussians, depending on  $z_i$ . This is called *the mixture of [multi-variate] Gaussians* model. Remember that

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

*Note* that the  $z_i$ 's are *latent* random variables, meaning that they're hidden/unobserved.

[Use the EM general scheme (see Tom Mitchell's *Machine Learning* book, 1997, pag. 194-195) to] prove that the EM algorithm for estimating the parameters  $\pi, \mu$  and  $\Sigma$  of our mixture of multi-variate Gaussian distributions has the following *update rules*:

**E-step:**

$$w_{ij} \stackrel{not.}{=} E[z_i = j | x_i; \pi', \mu', \Sigma'] = \frac{\mathcal{N}(x_i; \mu', \Sigma') \pi_j}{\sum_{l=1}^K \mathcal{N}(x_i; \mu', \Sigma') \pi_l} \quad (12)$$

**M-step:**

$$\pi_j = \frac{1}{n} \sum_{i=1}^n w_{ij}, \quad (13)$$

$$\mu_j = \frac{\sum_{i=1}^n w_{ij} x_i}{\sum_{i=1}^n w_{ij}}. \quad (14)$$

$$\Sigma_j = \frac{\sum_{i=1}^n w_{ij} (x_i - \mu_j)(x_i - \mu_j)^\top}{\sum_{i=1}^n w_{ij}}. \quad (15)$$

where  $\pi', \mu'$  and  $\Sigma'$  represent the values of our parameters at initialization, and respectively the previous iteration of the EM algorithm.

**Hint:** You may find useful the following formulas (from *Matrix Identities*, by Sam Roweis, 1999):

$$(1e) \quad (A^{-1})^\top = (A^\top)^{-1}$$

$$(2b) \quad |A^{-1}| = \frac{1}{|A|}$$

$$(4a) \quad \frac{\partial}{\partial X} |AXB| = |AXB|(X^{-1})^\top = |AXB|(X^\top)^{-1}$$

$$(4b) \quad \frac{\partial}{\partial X} \ln |X| = (X^{-1})^\top = (X^\top)^{-1}$$

$$(5a) \quad \frac{\partial}{\partial X} a^\top X = \frac{\partial}{\partial X} X^\top a = a$$

$$(5b) \quad \frac{\partial}{\partial X} X^\top AX = (A + A^\top)X$$

$$(5c) \quad \frac{\partial}{\partial X} a^\top Xb = ab^\top$$

$$(5e) \quad \frac{\partial}{\partial X} a^\top Xa = \frac{\partial}{\partial X} a^\top X^\top a = aa^\top$$

$$(5g) \quad \frac{\partial}{\partial X} (Xa + b)^\top C(Xa + b) = (C + C^\top)(Xa + b)a^\top$$

## Solution

The *E-step* is easy (use Bayes rule):

$$\begin{aligned}
 w_{ij} &\stackrel{not.}{=} E[z_i = j | x_i; \pi', \mu', \Sigma'] = p(z_i = j | x_i; \pi', \mu', \Sigma') = \\
 &= \frac{p(x_i | z_i = j; \mu', \Sigma') p(z_i = j; \pi')}{\sum_{l=1}^K p(x_i | z_i = l; \mu', \Sigma') p(z_i = l; \pi')} = \frac{\mathcal{N}(x_i; \mu', \Sigma') \pi'_j}{\sum_{l=1}^K \mathcal{N}(x_i; \mu', \Sigma') \pi_l}
 \end{aligned}$$

We will now concentrate on the *M-step*:

According to *the general EM scheme*, we need to maximize, with respect to our parameters  $\pi, \mu, \Sigma$ , the “auxiliary” function

$$\begin{aligned}
 Q(\pi, \mu, \Sigma | \pi', \mu', \Sigma') &\stackrel{\text{def.}}{=} E_{p(z_i=j|x_i; \pi', \mu', \Sigma')} \ln p(x, z; \mu, \Sigma, \pi) \\
 &= \sum_{i=1}^n \sum_{j=1}^K p(z_i = j | x_i; \pi', \mu', \Sigma') \ln p(x_i, z_i = j; \mu, \Sigma, \pi) \\
 &= \sum_{i=1}^n \sum_{j=1}^K p(z_i = j | x_i; \pi', \mu', \Sigma') \ln(p(x_i | z_i = j; \mu, \Sigma) p(z_i = j; \pi)) \\
 &= \sum_{i=1}^n \sum_{j=1}^K w_{ij} \ln \frac{1}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} \cdot \exp\left(-\frac{1}{2}(x_i - \mu_j)^\top \Sigma_j^{-1} (x_i - \mu_j)\right) \cdot \pi_j \\
 &= \sum_{i=1}^n \sum_{j=1}^K w_{ij} \left[ -\ln((2\pi)^{d/2} |\Sigma_j|^{1/2}) - \frac{1}{2}(x_i - \mu_j)^\top \Sigma_j^{-1} (x_i - \mu_j) + \ln \pi_j \right] \quad (16)
 \end{aligned}$$



First, let's derive the M-step update rule for  $\mu_l$ , with  $l = 1, \dots, K$ .

We have to maximize (16) with respect to  $\mu_l$ , so let's compute the corresponding derivative:

$$\begin{aligned}
& \frac{\partial}{\partial \mu_l} \sum_{i=1}^n \sum_{j=1}^K w_{ij} \left[ -\ln((2\pi)^{d/2} |\Sigma_j|^{1/2}) - \frac{1}{2} (x_i - \mu_j)^\top \Sigma_j^{-1} (x_i - \mu_j) + \ln \pi_j \right] \\
&= -\frac{\partial}{\partial \mu_l} \sum_{i=1}^n \sum_{j=1}^K w_{ij} \frac{1}{2} (x_i - \mu_j)^\top \Sigma_j^{-1} (x_i - \mu_j) = -\frac{1}{2} \sum_{i=1}^n w_{ij} \frac{\partial}{\partial \mu_l} \sum_{i=1}^n (x_i - \mu_j)^\top \Sigma_j^{-1} (x_i - \mu_j) \\
&\stackrel{(5g)}{=} -\frac{1}{2} \sum_{i=1}^n w_{ij} \frac{\partial}{\partial \mu_l} \sum_{i=1}^n (\Sigma_j^{-1} + (\Sigma_j^{-1})^\top) (x_i - \mu_j) \\
&\stackrel{(1e)}{=} \frac{1}{2} \sum_{i=1}^n w_{ij} \sum_{i=1}^n (\Sigma_j^{-1} + \underbrace{(\Sigma_j^\top)^{-1}}_{\Sigma_j}) (x_i - \mu_j) = \frac{1}{2} \sum_{i=1}^n w_{ij} \sum_{i=1}^n 2 \Sigma_j^{-1} (x_i - \mu_j) \\
&= \sum_{i=1}^n w_{il} (\Sigma_l^{-1} x_i - \Sigma_l^{-1} \mu_l) = \sum_{i=1}^n w_{il} \Sigma_l^{-1} x_i - \sum_{i=1}^n w_{il} \Sigma_l^{-1} \mu_l.
\end{aligned}$$

Setting this to zero and solving for  $\mu_l$  therefore yields the update rule

$$\mu_l = \frac{\sum_{i=1}^n w_{il} x_i}{\sum_{i=1}^n w_{il}}.$$

Secondly, we'll derive the M-step updates to  $\Sigma_j$ , for  $j = 1, \dots, K$ . Grouping together only the terms that depend on  $\Sigma_j$  in (16), we find that we need to maximize

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=1}^K w_{ij} \left[ \ln \frac{1}{|\Sigma_j|^{1/2}} - \frac{1}{2} (x_i - \mu_j)^\top \Sigma_j^{-1} (x_i - \mu_j) \right] \\ & \stackrel{(2b)}{=} \sum_{i=1}^n \sum_{j=1}^K w_{ij} \left[ \ln |\Sigma_j^{-1}|^{1/2} - \frac{1}{2} (x_i - \mu_j)^\top \Sigma_j^{-1} (x_i - \mu_j) \right]. \end{aligned}$$

We use the usual trick of working with the precision matrix  $\Lambda_j \stackrel{not.}{=} \Sigma_j^{-1}$ , where  $\Sigma_j$  is assumed invertible.

When maximizing the above quantity with respect to  $\Lambda_j$  by taking derivatives, we find:

$$\begin{aligned}
& \frac{\partial}{\partial \Lambda_j} \sum_{i=1}^n w_{ij} \left[ \ln |\Lambda_j|^{1/2} - \frac{1}{2} (x_i - \mu_j)^\top \Lambda_j (x_i - \mu_j) \right] \\
&= \frac{1}{2} \sum_{i=1}^n w_{ij} \frac{\partial}{\partial \Lambda_j} \ln |\Lambda_j| - \frac{1}{2} \sum_{i=1}^n w_{ij} \frac{\partial}{\partial \Lambda_j} [(x_i - \mu_j)^\top \Lambda_j (x_i - \mu_j)] \\
&\stackrel{(4b),(5c)}{=} \frac{1}{2} \sum_{i=1}^n w_{ij} \Lambda_j^{-1} - \frac{1}{2} \sum_{i=1}^n w_{ij} (x_i - \mu_j)(x_i - \mu_j)^\top \\
&= \frac{1}{2} \Lambda_j^{-1} \sum_{i=1}^n w_{ij} - \frac{1}{2} \sum_{i=1}^n w_{ij} (x_i - \mu_j)(x_i - \mu_j)^\top.
\end{aligned}$$

Setting this to zero and solving, we get:

$$\Sigma_j = \Lambda_j^{-1} = \frac{\sum_{i=1}^n w_{ij} (x_i - \mu_j)(x_i - \mu_j)^\top}{\sum_{i=1}^n w_{ij}}.$$

Finally, let's derive the M-step update for the parameters  $\pi_j$ . Grouping together only the terms that depend on  $\pi_j$  in (16), we find that we need to maximize

$$\sum_{i=1}^n \sum_{j=1}^K w_{ij} \ln \pi_j.$$

However, there is an additional constraint that the  $\pi_j$ 's sum to 1, since they represent the probabilities  $\pi_j = p(z_i = j; \pi)$ . To deal with the constraint that  $\sum_{j=1}^K \pi_j = 1$ , we construct the Lagrangian

$$\mathcal{L}(\pi) = \sum_{i=1}^n \sum_{j=1}^K w_{ij} \ln \pi_j + \beta \left( \sum_{j=1}^K \pi_j - 1 \right),$$

where  $\beta$  is the Lagrange multiplier.

*Note:* We don't need to worry about the constraint that  $\pi_j \geq 0$ , because as we'll shortly see, the solution we'll find from this derivation will automatically satisfy that anyway.

Taking derivatives of  $\mathcal{L}(\pi)$ , we find:

$$\frac{\partial}{\partial \pi_j} \mathcal{L}(\pi) = \sum_{i=1}^n \frac{w_{ij}}{\pi_j} + \beta = \frac{1}{\pi_j} \sum_{i=1}^n w_{ij} + \beta.$$

Setting this to zero and solving, we get  $\pi_j = \frac{\sum_{i=1}^n w_{ij}}{-\beta}$ .

By using the constraint  $\sum_j \pi_j = 1$ , and given the fact that  $\sum_j w_{ij} = 1$  since  $w_{ij} \stackrel{\text{not.}}{=} p(z_i = j | x_i; \pi', \mu', \Sigma')$ , we easily find

$$-\beta = \sum_{i=1}^n \sum_{j=1}^K w_{ij} = \sum_{i=1}^n 1 = n.$$

We therefore have our M-step derivation for the parameters  $\pi_j$ :

$$\pi_j = \frac{1}{n} \sum_{i=1}^n w_{ij},$$

and, obviously,  $\pi_j \geq 0$ .

## Remarks

1. Let's contrast the update rules in the M-step with the formulas we would have when the  $z_i$ s were known exactly (see the MLE of the parameters of a single multi-variate Gaussian distribution, CMU, 2010 fall, Aarti Singh, HW1, pr. 3.2.1):

$$\begin{aligned}\pi_j &= \frac{1}{n} \sum_{i=1}^n 1\{z_i = j\}, \\ \mu_j &= \frac{\sum_{i=1}^n 1\{z_i = j\} x_i}{\sum_{i=1}^n 1\{z_i = j\}}, \\ \Sigma_j &= \frac{\sum_{i=1}^n 1\{z_i = j\} (x_i - \mu_j)(x_i - \mu_j)^\top}{\sum_{i=1}^n 1\{z_i = j\}},\end{aligned}$$

with  $1\{z_i = j\}$  (“indicator functions”) indicating from which Gaussian each datapoint had come.

They are identical, except that instead of the indicator functions  $1\{z_i = j\}$  indicating from which Gaussian each datapoint had come, we now have the  $w_{ij}$ 's.

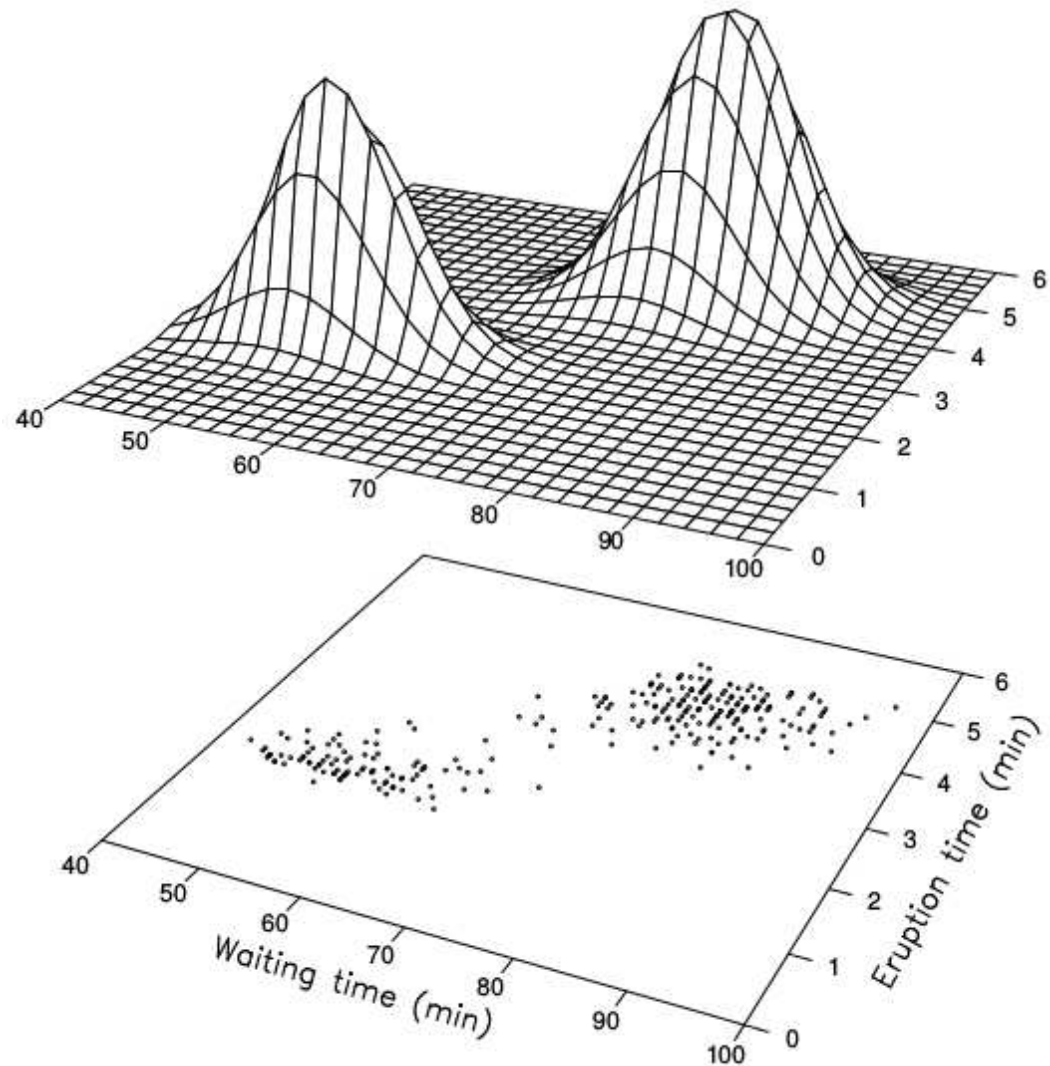
## Remarks (cont'd)

2. The EM-algorithm is reminiscent of the  $K$ -means clustering algorithm, except that instead of the “hard” cluster assignments  $c(i)$ , we have the “soft” assignments  $w_{ij}$ .
3. Similar to  $K$ -means, the EM algorithm is also susceptible to local optima, so reinitializing at several different initial parameters may be a good idea.
4. It's clear that the EM algorithm has a very natural interpretation of repeatedly trying to guess the unknown  $z_i$ 's.

## Example:

Modelling the waiting and eruption times for  
*the Old Faithful geyser*,  
(Yellowstone Park, USA)

Michael Eichler  
University of Chicago  
Statistics course (24600),  
Spring 2004





A link between  
*K*-means and EM/GMM (the multi-variate case)

CMU, 2008 fall, Eric Xing, HW4, pr. 2.2

(see also CMU, 2010 fall, Aarti Singh, HW4, pr. 1.2)

Given  $N$  data points  $x_i$ , ( $i = 1, \dots, N$ ),  $K$ -means will group them into  $K$  clusters by minimizing the *distortion* function

$$J = \sum_{i=1}^N \sum_{j=1}^K \gamma_{ij} \|x_i - \mu_j\|^2,$$

where  $\mu_j$  is the centroid of the  $j$ -th cluster, and  $\gamma_{ij} = 1$  if  $x_i$  belongs to the  $j$ -th cluster and  $\gamma_{ij} = 0$  otherwise.

In this exercise, we will use the following procedure for  $K$ -means:

- Initialize [randomly] the cluster centroids  $\mu_j$ ,  $j = 1, \dots, K$ ;
- Iterate until *convergence*:
  - for every data point  $x_n$ , update its *cluster assignment*:  $\gamma_{ij} = 1$  if  $j = \arg \min_k \|x_i - \mu_k\|^2$ , and  $\gamma_{ij} = 0$  otherwise.
  - for each cluster  $j$ , update its centroid:  $\mu_j = \frac{\sum_{i=1}^N \gamma_{ij} x_i}{\sum_{i=1}^N \gamma_{ij}}$

Remember that in GMM,  $p(x) = \sum_{j=1}^K \pi_j \mathcal{N}(x|\mu_j, \Sigma_j)$ , where  $\pi_j$  is the prior [probability] for the  $j^{\text{th}}$  component,  $\mu_j$  and  $\Sigma_j$  are the mean and covariance matrix for the  $j^{\text{th}}$  component respectively. In the E-step of the EM algorithm, we will update

$$p(z_{ij} = 1|x_i) = \frac{\pi_j \mathcal{N}(x_i|\mu_j, \Sigma_j)}{\sum_{k=1}^K \pi_k \mathcal{N}(x_i|\mu_k, \Sigma_k)}$$

Now, suppose that

- i.  $\Sigma_k = \sigma^2 I$ , for some  $\sigma > 0$ , and for all  $k = 1, \dots, K$
- ii.  $\pi_k \neq 0$  for  $k = 1, \dots, K$  [LC: at any iteration of the EM algorithm], and
- iii.  $\|x_i - \mu_{k'}\| \neq \|x_i - \mu_k\|$  for any  $k' \neq k$  [at any iteration of the EM algorithm].

Under the above assumptions, prove that when  $\sigma \rightarrow 0_+$  we will get  $p(z_{ij} = 1|x_i) \rightarrow \gamma_{ij}$ , where  $\gamma_{ij}$  is the *cluster assignment* used in  $K$ -means.

**Answer:**

$$\begin{aligned}
 p(z_{ij} = 1|x_i) &= \frac{\pi_j \mathcal{N}(x_i|\mu_j, \Sigma_j)}{\sum_{k=1}^K \pi_k \mathcal{N}(x_i|\mu_k, \Sigma_k)} = \frac{\pi_j \exp\left(-\frac{1}{2\sigma^2}\|x_i - \mu_j\|^2\right)}{\sum_{k=1}^K \pi_k \exp\left(-\frac{1}{2\sigma^2}\|x_i - \mu_k\|^2\right)} \\
 &= \frac{1}{1 + \sum_{k \neq j} \frac{\pi_k}{\pi_j} \exp\left(\frac{1}{2\sigma^2}(\|x_i - \mu_j\|^2 - \|x_i - \mu_k\|^2)\right)}
 \end{aligned}$$

**Case 1:**

**If  $\|x_i - \mu_j\| = \min_k \|x_i - \mu_k\|$ , then for each  $k \neq j$  we have  $\|x_i - \mu_j\|^2 - \|x_i - \mu_k\|^2 < 0$ . Since  $\sigma \rightarrow 0_+$ , it will follow that  $\exp\left(\frac{1}{2\sigma^2}(\|x_i - \mu_j\|^2 - \|x_i - \mu_k\|^2)\right) \rightarrow 0$ . So,  $p(z_{ij} = 1|x_i) \rightarrow 1$ .**

**Case 2:**

**If  $\|x_i - \mu_j\| \neq \min_k \|x_i - \mu_k\|$ , then**

- **for all  $k$  such that  $\|x_i - \mu_j\| < \|x_i - \mu_k\|$  it will follow (exactly as above) that  $\exp\left(\frac{1}{2\sigma^2}(\|x_i - \mu_j\|^2 - \|x_i - \mu_k\|^2)\right) \rightarrow 0$ ;**
- **for all  $k$  such that  $\|x_i - \mu_j\| > \|x_i - \mu_k\|$  we will have  $\exp\left(\frac{1}{2\sigma^2}(\|x_i - \mu_j\|^2 - \|x_i - \mu_k\|^2)\right) \rightarrow +\infty$ .**

**Therefore,  $p(z_{ij} = 1|x_i) \rightarrow \frac{1}{1 + \infty} = 0$ .**