

Intelligent Systems

Escuela Técnica Superior de Informática

Universitat Politècnica de València

Block 2 Chapter 4:

Clustering

Unsupervised learning: K-means algorithm

Índice

- 1 Introduction ▷ 2
- 2 Partitional clustering ▷ 4
- 3 C-means algorithm ▷ 13

Índice

- 1 *Introduction* ▷ 2
- 2 Partitional clustering ▷ 4
- 3 C-means algorithm ▷ 13

Introduction

Now we shall investigate ***unsupervised learning*** or ***clustering***, which use *unlabeled samples*.

According to [Anderberg,1973], the aim of clustering is:

To group objects in classes so the objects belonging to the same class have a high degree of ***natural association***, while the other classes are relatively different. The purpose is to create classes that are relatively different from each other.

In other words:

To find ***natural groupings*** of a set of input patterns so that the descriptions of these objects can be done in terms of classes or groups with strong internal similarities.

We need a distance (similarity) function to assess the quality of a clustering result (high similarity within a cluster and low similarity between clusters).

Two types of clustering: ***Partitional*** and ***Hierarchical***.

Índice

- 1 Introduction ▷ 2
- 2 *Partitional clustering* ▷ 4
- 3 C-means algorithm ▷ 13

Partitional Clustering

General problem:

Given a set of N objects or observations (x_1, x_2, \dots, x_N) , where each object is a d -dimensional real vector, the objective is to find a partition Π of the N objects into C **classes or clusters**; i.e.: $\Pi = \{X_1, \dots, X_C\}$, where X_1, \dots, X_C denote the C clusters. In order to do so, we assume there is available a **critierion function** J to evaluate the quality of a partition Π . Thus, the clustering problem can be seen as a search problem:

$$\Pi^* = \arg \min_{\Pi = \{X_1, \dots, X_C\}} J(\Pi) \quad (\text{best partition, the one that minimizes } J \text{ for each partition}) \quad (1)$$

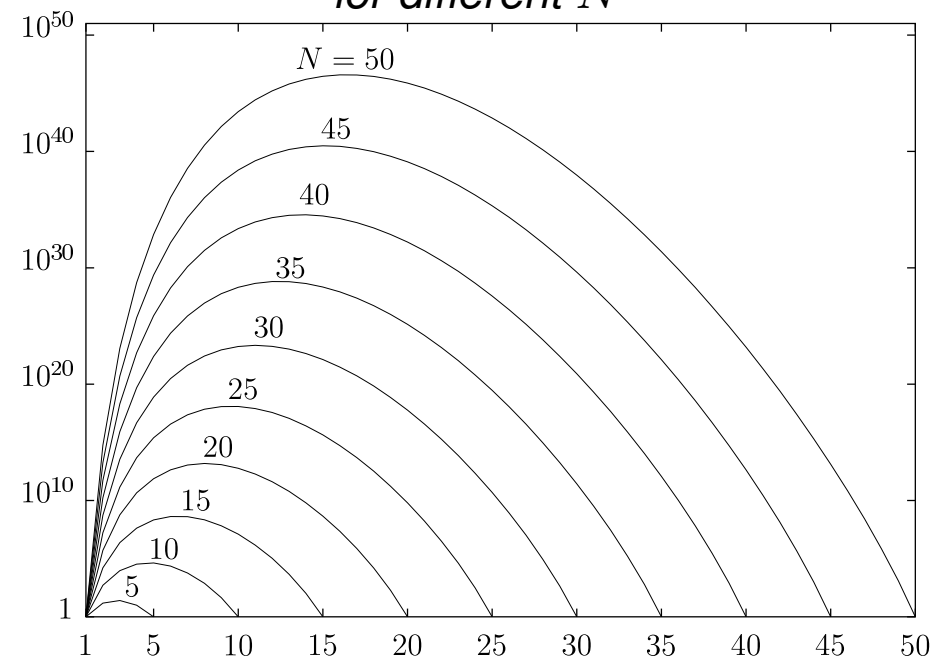
Difficulty:

The number of partitions to evaluate is too high even for small values of N and C (see right). Searching for global optimal solutions through enumeration techniques (explicitly or implicitly) is not feasible except for a few particular cases.

Solution:

Search for suboptimal solutions obtained by approximate algorithms.

Number of partitions with respect to C for different N



Partitional clustering: “Sum of Squared Errors” (SSE) criterion

The SSE is a criterion to select a partition Π of N objects into C clusters. The objective is to select the partition that minimizes SSE.

The SSE value of a partition Π of N objects (x_1, x_2, \dots, x_N) into C clusters, $\Pi = \{X_1, \dots, X_C\}$, is computed as follows:

- Each cluster X_c contains a subset of the N objects (disjoint sets)
- For each cluster X_c , compute the squared errors for all $x \in X_c$: $SE_x = \|x - m_c\|^2$.
 m_c is the mean of points in X_c :

$$m_c = \frac{1}{|X_c|} \sum_{x \in X_c} x \quad (2)$$

m_c is interpreted as the natural prototype of X_c . Each object $x \in X_c$, is interpreted as a “distorted version” of m_c and the distortion of x is modeled by the *error vector* $x - m_c$.

- For each cluster X_c , J_c is the function for cluster X_c . $J_c = WCSE_c$ (within-cluster squared errors)

$$J_c = WCSE_c = \sum_{x \in X_c} \|x - m_c\|^2 \quad (3)$$

- The SSE value of a partition Π is computed as:

$$SSE_{\Pi} = J(X_1, \dots, X_C) = \sum_c J_c = \sum_c WCSE_c \quad (4)$$

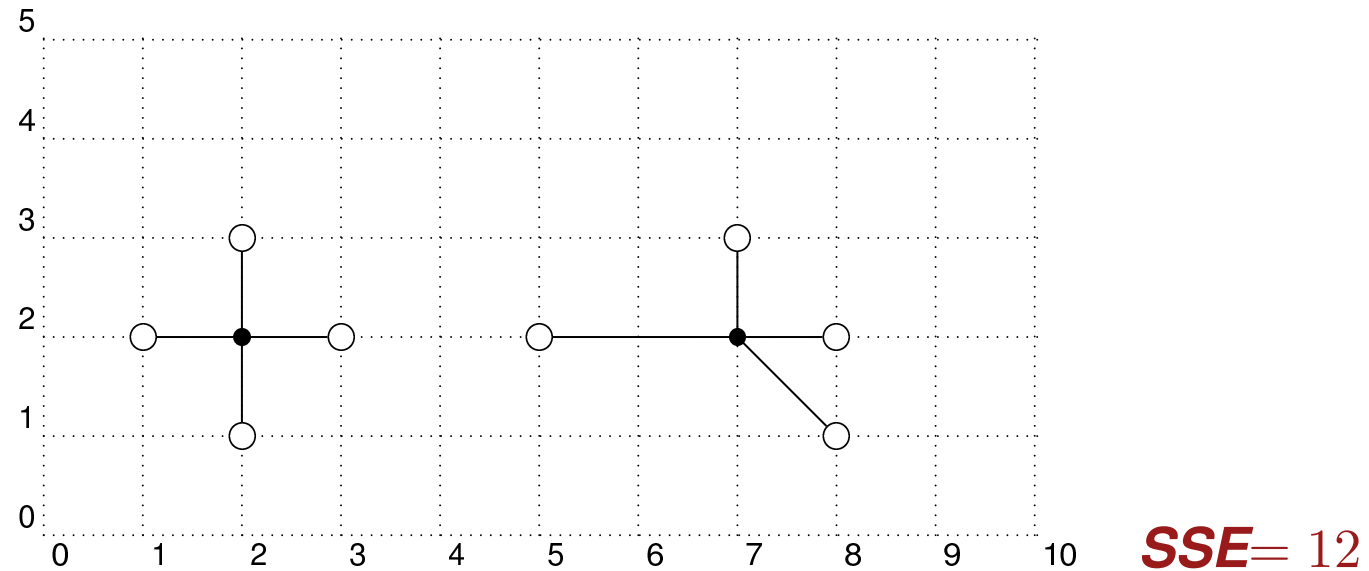
SSE criterion (continuation)

Now, we would have to calculate the SSE value for every partition (if they are explicitly given) and take the one that minimizes SSE (Π^* , best partition)

SSE measures the sum (or average) of the squares of the magnitudes of the error vectors, and, obviously, it is a criterion to minimize.

The *mean* of each cluster, m_c , is the point that represents the object of this cluster with the lowest SSE. m_c is the center of the cluster X_c (called *centroid*).

Partitional clustering example (partition I)



Partition I (Π_1): 8 samples (x_1, \dots, x_N) , 2 clusters in partition $\Pi_1 = \{X_1, X_2\}$

$$X_1 = \{x_1, \dots, x_4\} : x_1 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, x_2 = \begin{pmatrix} 2 \\ 1 \end{pmatrix}, x_3 = \begin{pmatrix} 2 \\ 3 \end{pmatrix}, x_4 = \begin{pmatrix} 3 \\ 2 \end{pmatrix}$$

$$X_2 = \{x_5, \dots, x_8\} : x_5 = \begin{pmatrix} 5 \\ 2 \end{pmatrix}, x_6 = \begin{pmatrix} 7 \\ 3 \end{pmatrix}, x_7 = \begin{pmatrix} 8 \\ 1 \end{pmatrix}, x_8 = \begin{pmatrix} 8 \\ 2 \end{pmatrix}$$

$$m_1 = \frac{1}{4} \sum_{x \in X_1} x = \frac{1}{4} \begin{pmatrix} 8 \\ 8 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix} \quad m_2 = \frac{1}{4} \sum_{x \in X_2} x = \frac{1}{4} \begin{pmatrix} 28 \\ 8 \end{pmatrix} = \begin{pmatrix} 7 \\ 2 \end{pmatrix}$$

Partitional clustering example (partition I, continuation)

Calculations for cluster X_1 :

$$SE_{x_1} = \|\mathbf{x}_1 - \mathbf{m}_1\|^2 = \left(\begin{pmatrix} 1 \\ 2 \end{pmatrix} - \begin{pmatrix} 2 \\ 2 \end{pmatrix} \right)^2 = \begin{pmatrix} -1 \\ 0 \end{pmatrix}^2 = 1$$

$$SE_{x_2} = \|\mathbf{x}_2 - \mathbf{m}_1\|^2 = \left(\begin{pmatrix} 2 \\ 1 \end{pmatrix} - \begin{pmatrix} 2 \\ 2 \end{pmatrix} \right)^2 = \begin{pmatrix} 0 \\ -1 \end{pmatrix}^2 = 1$$

$$SE_{x_3} = \|\mathbf{x}_3 - \mathbf{m}_1\|^2 = \dots = 1$$

$$SE_{x_4} = \|\mathbf{x}_4 - \mathbf{m}_1\|^2 = \dots = 1$$

$$J_{X_1} = WCSE_{X_1} = \sum_{\mathbf{x} \in X_1} \|\mathbf{x} - \mathbf{m}_1\|^2 = 4 \quad (\text{squared errors within cluster } X_1)$$

Partitional clustering example (partition I, continuation)

Calculations for cluster X_2 :

$$SE_{x5} = \|\mathbf{x5} - \mathbf{m}_2\|^2 = \left(\begin{pmatrix} 5 \\ 2 \end{pmatrix} - \begin{pmatrix} 7 \\ 2 \end{pmatrix} \right)^2 = \begin{pmatrix} -2 \\ 0 \end{pmatrix}^2 = 4$$

$$SE_{x6} = \|\mathbf{x6} - \mathbf{m}_2\|^2 = \left(\begin{pmatrix} 7 \\ 3 \end{pmatrix} - \begin{pmatrix} 7 \\ 2 \end{pmatrix} \right)^2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}^2 = 1$$

$$SE_{x7} = \|\mathbf{x7} - \mathbf{m}_2\|^2 = \dots = 2$$

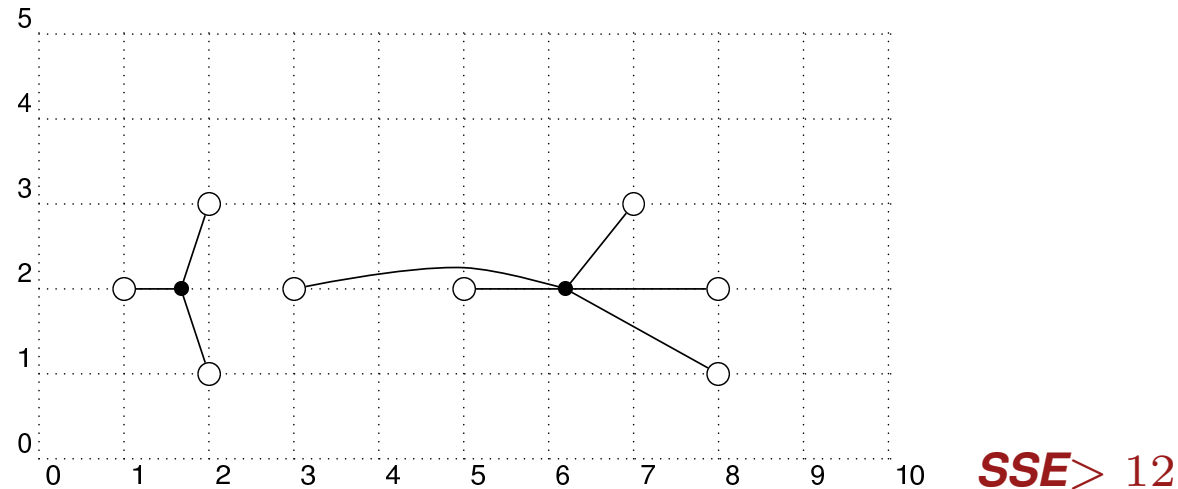
$$SE_{x8} = \|\mathbf{x8} - \mathbf{m}_2\|^2 = \dots = 1$$

$$J_{X_2} = WCSE_{X_2} = \sum_{\mathbf{x} \in X_2} \|\mathbf{x} - \mathbf{m}_2\|^2 = 8 \quad (\text{squared errors within cluster } X_2)$$

Sum of squared errors for all clusters in partition Π_1 :

$$J(\Pi_1) = J(X_1, X_2) = SSE_{\Pi_1} = \sum_c WCSE_c = 12$$

Partitional clustering example (partition II)



Partition II (Π_2): 8 samples (x_1, \dots, x_N) , 2 clusters in partition $\Pi_2 = \{X_1, X_2\}$

$$X_1 = \{x_1, \dots, x_3\} : x_1 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, x_2 = \begin{pmatrix} 2 \\ 1 \end{pmatrix}, x_3 = \begin{pmatrix} 2 \\ 3 \end{pmatrix}$$

$$X_2 = \{x_4, \dots, x_8\} : x_4 = \begin{pmatrix} 3 \\ 2 \end{pmatrix}, x_5 = \begin{pmatrix} 5 \\ 2 \end{pmatrix}, x_6 = \begin{pmatrix} 7 \\ 3 \end{pmatrix}, x_7 = \begin{pmatrix} 8 \\ 1 \end{pmatrix}, x_8 = \begin{pmatrix} 8 \\ 2 \end{pmatrix}$$

$$\mathbf{m}_1 = \frac{1}{3} \sum_{x \in X_1} x = \frac{1}{3} \begin{pmatrix} 5 \\ 6 \end{pmatrix} = \begin{pmatrix} 1.66 \\ 2 \end{pmatrix} \quad \mathbf{m}_2 = \frac{1}{5} \sum_{x \in X_2} x = \frac{1}{5} \begin{pmatrix} 31 \\ 10 \end{pmatrix} = \begin{pmatrix} 6.2 \\ 2 \end{pmatrix}$$

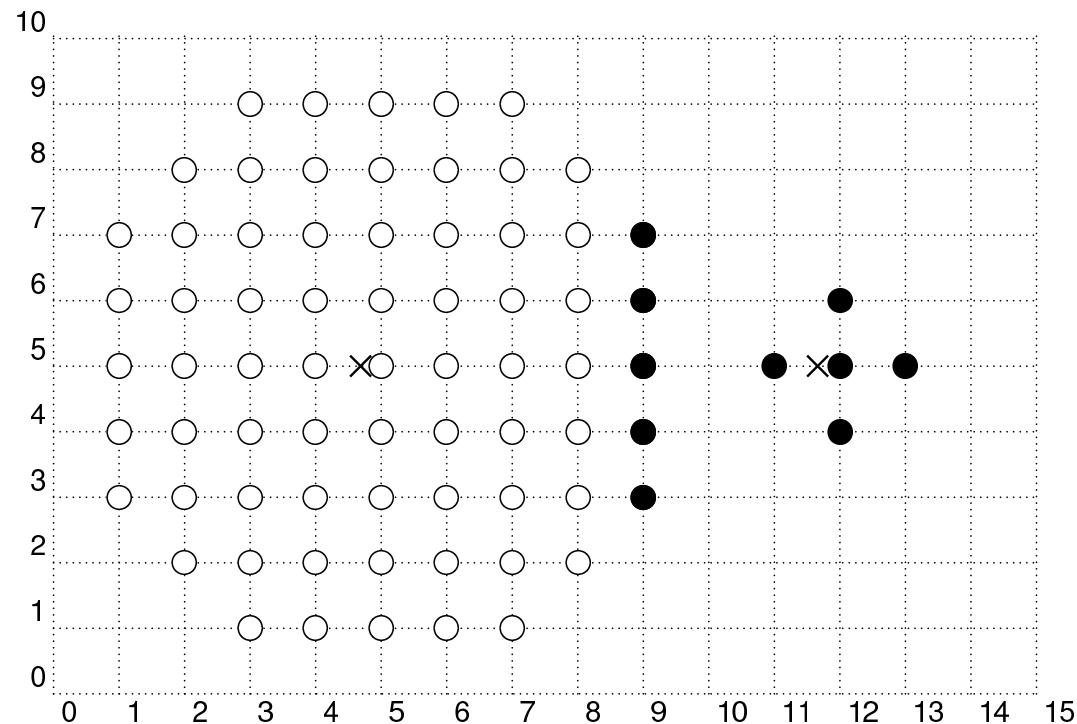
Exercise: calculate the values for $WCSE_{X_1}$ and $WCSE_{X_2}$ in partition Π_2 and check that the value of SSE_{Π_2} is higher than 12.

Goodness of the SSE criterion

The SSE criterion is appropriate only when the objects form *hyper-spherics clusters of similar size*.

If the sizes of the clusters are too different, it may happen that the natural grouping does not give the minimum SSE and so we would not find the natural clusters.

Therefore, the centroid representation alone works well if the clusters are of the hyper-spherical shape. If clusters are elongated or are of other shapes, centroids are not sufficient.



Índice

- 1 Introduction ▷ 2
- 2 Partitional clustering ▷ 4
- 3 *C-means algorithm* ▷ 13

Incremental computation of SSE after moving x from cluster X_i to X_j

$$X'_i = X_i - \{x\}$$

$$X'_j = X_j + \{x\}$$

$$\mathbf{m}'_i = \mathbf{m}_i - \frac{\mathbf{x} - \mathbf{m}_i}{n_i - 1}$$

$$\mathbf{m}'_j = \mathbf{m}_j + \frac{\mathbf{x} - \mathbf{m}_j}{n_j + 1}$$

$$J'_i = J_i - \frac{n_i}{n_i - 1} \|\mathbf{x} - \mathbf{m}_i\|^2$$

$$J'_j = J_j + \frac{n_j}{n_j + 1} \|\mathbf{x} - \mathbf{m}_j\|^2$$

$$\Delta J = \frac{n_j}{n_j + 1} \|\mathbf{x} - \mathbf{m}_j\|^2 - \frac{n_i}{n_i - 1} \|\mathbf{x} - \mathbf{m}_i\|^2$$

The movement will be successful if the increment of SSE is negative; that is:

$$\frac{n_j}{n_j + 1} \|\mathbf{x} - \mathbf{m}_j\|^2 < \frac{n_i}{n_i - 1} \|\mathbf{x} - \mathbf{m}_i\|^2 \quad (5)$$

***These equations allow to minimize the SSE
by using successive refinements from an initial partition.***

SSE optimization: algorithm C -means **(K -means in English bibliography)**

Input:

- $X = (x_1, x_2, \dots, x_N)$, the N objects/observations
- C , the number of clusters
- $\Pi_0 = \{X_1, \dots, X_C\}$, a randomly initial partition that classifies the N objects into C clusters; these will be the initial centroids (cluster centers)

Output:

- Π^* , the best partition
- $\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_C$, the cluster centers for Π^*
- J (optionally, it depends on the C -means version)

Goal:

The data given by X is clustered by C -means algorithm, which aims to partition the N objects into C clusters such that the sum of the squared errors from points to the assigned cluster is minimized.

Algorithm *K-means* (“correct” version [Duda & Hart])

Input: $X = (x_1, x_2, \dots, x_N)$; C ; $\Pi_0 = \{X_1, \dots, X_C\}$; (random initial partition)

Output: $\Pi^* = \{X_1, \dots, X_C\}$; $\mathbf{m}_1, \dots, \mathbf{m}_C$; $J(\Pi^*)$

for $c = 1$ **to** C **do** $\mathbf{m}_c = \frac{1}{n_c} \sum_{\mathbf{x} \in X_c} \mathbf{x}$ **endfor** (compute initial centroids for the clusters in input partition)

compute $J = J(\Pi_0)$

repeat

$transfers = \text{false}$

forall $\mathbf{x} \in X$ (let $i : \mathbf{x} \in X_i$) **do** (for all the objects ... assume object \mathbf{x} belongs to cluster X_i)

if $n_i > 1$ **then** (if cluster X_i has more than one object, not only \mathbf{x})

$j^* = \arg \min_{j \neq i} \frac{n_j}{n_j + 1} \|\mathbf{x} - \mathbf{m}_j\|^2$ (study \mathbf{x} in any other cluster different from X_i ; take the min. cluster X_{j^*})

$\Delta J = \frac{n_{j^*}}{n_{j^*} + 1} \|\mathbf{x} - \mathbf{m}_{j^*}\|^2 - \frac{n_i}{n_i - 1} \|\mathbf{x} - \mathbf{m}_i\|^2$ (increment of SSE, difference bt. X_{j^*} and current cluster)

if $\Delta J < 0$ **then**

$transfers = \text{true}$ (move \mathbf{x} to X_{j^*})

$\mathbf{m}_i = \mathbf{m}_i - \frac{\mathbf{x} - \mathbf{m}_i}{n_i - 1}$ $\mathbf{m}_{j^*} = \mathbf{m}_{j^*} + \frac{\mathbf{x} - \mathbf{m}_{j^*}}{n_{j^*} + 1}$

$X_i = X_i - \{\mathbf{x}\}$ $X_{j^*} = X_{j^*} + \{\mathbf{x}\}$

$J^* = J + \Delta J$ $J = J^*$ (J value of the best partition so far ...)

endif

endif

endforall

until $\neg transfers$ Cost per iteration: $O(N \cdot C \cdot D)$, $N = |X|$, $D = \text{dimension}$

Example of C -means application (Duda & Hart) (1)

$$X = \{1, 3, 4.5\} \subset \mathbb{R}; \quad x_1 = (1), x_2 = (3), x_3 = (4.5)$$

$$C = 2 \text{ clusters}; \quad X_1 = \{x_1, x_2\}, \quad X_2 = \{x_3\} \quad \Pi^0 = \{\{1, 3\}, \{4.5\}\}$$

$$\mathbf{m}_1 = 2 \quad J_1 = 2 \quad \mathbf{m}_2 = 4.5 \quad J_2 = 0 \quad J(\Pi^0) = 2$$

$x = x_1 = 1$ (Cluster $X_1, i = 1, n_1 = 2$: two objects in X_1)

A) we have only one more cluster (X_2) so $\arg \min_{j \neq i}$ only analyzes cluster $j = 2$ ($n_j = n_2 = 1$).

B)

$$j^* = \arg \min_{j \neq i} \frac{n_j}{n_j + 1} \|\mathbf{x} - \mathbf{m}_j\|^2 = \arg \min_{j \neq i} \frac{n_2}{n_2 + 1} \|\mathbf{x}_1 - \mathbf{m}_2\|^2 =$$

$$\frac{1}{2} \|1 - 4.5\|^2 = 6.125$$

C)

$$\Delta J = \frac{n_{j^*}}{n_{j^*} + 1} \|\mathbf{x} - \mathbf{m}_{j^*}\|^2 - \frac{n_i}{n_i + 1} \|\mathbf{x} - \mathbf{m}_i\|^2 =$$

$$\frac{n_2}{n_2 + 1} \|\mathbf{x}_1 - \mathbf{m}_2\|^2 - \frac{n_1}{n_1 + 1} \|\mathbf{x}_1 - \mathbf{m}_1\|^2 =$$

$$6.125 - \frac{2}{1} \|1 - 2\|^2 = 6.125 - 2 = 4.125$$

D) x_1 is not transferred to X_2

Example of C -means application (Duda & Hart) (2)

$x = x_2 = 3$ (Cluster X_1 , $i = 1$, $n_1 = 2$: two objects in X_1)

A) we have only one more cluster (X_2) so $\arg \min_{j \neq i}$ only analyzes cluster $j = 2$ ($n_j = n_2 = 1$).

B)

$$j^* = \arg \min_{j \neq i} \frac{n_j}{n_j + 1} \|\mathbf{x} - \mathbf{m}_j\|^2 = \arg \min_{j \neq i} \frac{n_2}{n_2 + 1} \|\mathbf{x}_2 - \mathbf{m}_2\|^2 = \frac{1}{2} \|3 - 4.5\|^2 = 1.125$$

C)

$$\begin{aligned} \Delta J &= \frac{n_{j^*}}{n_{j^*} + 1} \|\mathbf{x} - \mathbf{m}_{j^*}\|^2 - \frac{n_i}{n_i - 1} \|\mathbf{x} - \mathbf{m}_i\|^2 = \\ &= \frac{n_2}{n_2 + 1} \|\mathbf{x}_2 - \mathbf{m}_2\|^2 - \frac{n_1}{n_1 - 1} \|\mathbf{x}_2 - \mathbf{m}_1\|^2 = \\ &= 1.125 - \frac{2}{1} \|3 - 2\|^2 = 1.125 - 2 = -0.875 \end{aligned}$$

D) *transfers = true*, we move x_2 to cluster X_2

E)

$$\begin{aligned} \mathbf{m}_i &= \mathbf{m}_i - \frac{\mathbf{x} - \mathbf{m}_i}{n_i - 1} = \mathbf{m}_1 - \frac{\mathbf{x}_2 - \mathbf{m}_1}{n_1 - 1} = 2 - \frac{3 - 2}{2 - 1} = 1 \quad (m_1 = 1) \\ \mathbf{m}_{j^*} &= \mathbf{m}_{j^*} + \frac{\mathbf{x} - \mathbf{m}_{j^*}}{n_{j^*} + 1} = \mathbf{m}_2 + \frac{\mathbf{x}_2 - \mathbf{m}_2}{n_2 + 1} = 4.5 + \frac{3 - 4.5}{1 + 1} = 3.75 \quad (m_2 = 3.75) \end{aligned}$$

Example of C -means application (Duda & Hart) (3)

F) We update the clusters: $X_1 = \{x_1\}$, $X_2 = \{x_2, x_3\}$

G) $J = J + \Delta J = 2 - 0.875 = 1.125$

$x = x_3 = 4.5$ (Cluster X_2 , $i = 2$, $n_2 = 2$: two objects in X_2)

A) we have only one more cluster (X_1) so $\arg \min_{j \neq i}$ only analyzes cluster $j = 1$ ($n_j = n_1 = 1$).

B)

$$j^* = \arg \min_{j \neq i} \frac{n_j}{n_j + 1} \|\mathbf{x} - \mathbf{m}_j\|^2 = \arg \min_{j \neq i} \frac{n_1}{n_1 + 1} \|\mathbf{x}_3 - \mathbf{m}_1\|^2 =$$

$$\frac{1}{2} \|4.5 - 1\|^2 = 6.125$$

C)

$$\Delta J = \frac{n_{j^*}}{n_{j^*} + 1} \|\mathbf{x} - \mathbf{m}_{j^*}\|^2 - \frac{n_i}{n_i + 1} \|\mathbf{x} - \mathbf{m}_i\|^2 =$$

$$\frac{n_1}{n_1 + 1} \|\mathbf{x}_3 - \mathbf{m}_1\|^2 - \frac{n_2}{n_2 + 1} \|\mathbf{x}_3 - \mathbf{m}_2\|^2 =$$

$$6.125 - \frac{2}{1} \|4.5 - 3.75\|^2 = 6.125 - 1.125 = 5$$

D) x_3 is not transferred to X_1

Example of C -means application (Duda & Hart) (4)

The algorithm makes one more iteration of the loop 'repeat ... until' because *transfers* is true; in the next iteration of the algorithm no modifications are produced.

Final result: $\Pi^* = \{X_1, X_2\}$, $X_1 = \{x_1\}$, $X_2 = \{x_2, x_3\}$, $J(\Pi^*) = SSE_{\Pi^*} = 1.125$

Algorithm *C-means* (“popular” version)

Output: $\Pi^* = \{X_1, \dots, X_C\}; m_1, \dots, m_C$

transfers = false

forall $x \in X$ (let $i : x \in X_i$) **do**

$$j^* = \arg \min_{1 \leq j \leq C} d(\mathbf{x}, \mathbf{m}_j) \quad (\text{Euclidean distance from } x \text{ to each centroid; } d(\mathbf{x}, \mathbf{m}_j) = \|\mathbf{x} - \mathbf{m}_j\|)$$

if $j^* \neq i$ then (if the min. cluster is different from the current cluster, move x to the min. cluster)

$$X_i = X_i - \{\mathbf{x}\}; X_{j^*} = X_{j^*} + \{\mathbf{x}\}$$
endif

until $\neg transfers$ (Cost per iteration: $O(N \cdot C \cdot D)$, $N = |X|$, $D = \text{cost of } d(\cdot, \cdot)$)

Página Block 2 Chapter 4.21

Optimality of algorithms *C*-means

Exercise: Consider the same example in page 17. Apply the popular version of the *C*-means algorithm. The results are:

$\Pi^* = \{X_1, X_2\}$, $X_1 = \{x_1, x_2\}$, $X_2 = \{x_3\}$, $J(\Pi^*) = SSE_{\Pi^*} = 2.0$; that is:

$\Pi^* = \Pi^0$, $J(\Pi^*) = J(\Pi^0)$

- None of the two *C-means* versions guarantee a global minimum of the SSE
- The version presented in the book Duda & Hart obtains a *local* minimum
- The “popular” version does not guarantee the local minimization in some cases

The SSE criterion and Vectorial Quantification

A common way to represent clusters is the centroid representation: use the centroid of each cluster as the representative of the cluster $\mathbf{r}_1, \dots, \mathbf{r}_C$.

The next criteria to minimize are equivalent:

$$J(X_1, \dots, X_C) = \sum_c \sum_{\mathbf{x} \in X_c} \|\mathbf{x} - \mathbf{m}_c\|^2 \quad (6)$$

$$J(X_1, \dots, X_C; \mathbf{r}_1, \dots, \mathbf{r}_C) = \sum_c \sum_{\mathbf{x} \in X_c} \|\mathbf{x} - \mathbf{r}_c\|^2 \quad (7)$$

$$J(\mathbf{r}_1, \dots, \mathbf{r}_C) = \sum_{\mathbf{x}} \min_c \|\mathbf{x} - \mathbf{r}_c\|^2 \quad (8)$$

Justification:

- (7) is equivalent to (6) because for all partitions X_1, \dots, X_C , the cluster representatives $\mathbf{r}_1, \dots, \mathbf{r}_C$ that minimize (7) are the mean of the clusters.
- (7) is equivalent to (8) because for all representative set $\mathbf{r}_1, \dots, \mathbf{r}_C$, the partition that minimizes (7) is the one in which each object is assigned the nearest cluster representative.
- (8) is known as the problem of *vectorial quantification design* in *information theory*

Another interpretation of SSE

SSE can also be defined without using the *mean* of the clusters:

$$J(X_1, \dots, X_C) = \frac{1}{2} \sum_c n_c \bar{s}_c \quad (9)$$

where n_c is the number of objects in X_c and \bar{s}_c is the mean of the squared Euclidean distances for each pair of objects in the cluster.

$$\bar{s}_c = \frac{1}{n_c^2} \sum_{\mathbf{x}, \mathbf{x}' \in X_c} \|\mathbf{x} - \mathbf{x}'\|^2 \quad (10)$$

Thus, SSE can also be interpreted as a weighted average of intra-cluster squared distances.

Under this interpretation, it is possible to redefine \bar{s}_c and obtain criteria similar to SSE (valid even for *non-vectorial* data):

$$\bar{s}_c = \frac{1}{n_c^2} \sum_{x, x' \in X_c} d(x, x') \quad \bar{s}_c = \frac{1}{n_c^2} \max_{x, x' \in X_c} d(x, x') \quad (11)$$