

## Indicative exercises for clustering algorithms

### Solved exercises

**Exercise 1 (proximity measures):** Prove that the **Euclidean distance** between two  $l$ -dimensional vectors  $\mathbf{x} = [x_1, \dots, x_l]^T$  and  $\mathbf{y} = [y_1, \dots, y_l]^T$ ,  $d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^l (x_i - y_i)^2}$ , is a **metric dissimilarity measure**.

*Hint:* The *Minkowski inequality* states that for two  $l$ -dimensional vectors  $\mathbf{x} = [x_1, \dots, x_l]^T$  and  $\mathbf{y} = [y_1, \dots, y_l]^T$  and a positive integer  $p$ , it holds  $(\sum_{i=1}^l |x_i + y_i|^p)^{1/p} \leq (\sum_{i=1}^l |x_i|^p)^{1/p} + (\sum_{i=1}^l |y_i|^p)^{1/p}$

**Solution:** In order to prove that  $d(\mathbf{x}, \mathbf{y})$  is a metric dissimilarity measure, we need to prove that it possesses all the following properties:

- (a)  $\exists d_0 \in R$ , so that  $d_0 \leq d(\mathbf{x}, \mathbf{y}) < +\infty$ , for all  $\mathbf{x}, \mathbf{y} \in R^l$
- (b)  $d(\mathbf{x}, \mathbf{x}) = d_0$ , for all  $\mathbf{x} \in R^l$
- (c)  $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$ , for all  $\mathbf{x}, \mathbf{y} \in R^l$
- (d)  $d(\mathbf{x}, \mathbf{y}) = d_0 \Leftrightarrow \mathbf{x} = \mathbf{y}$
- (e)  $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$ , for all  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in R^l$

For (a): It is straightforward to see that for any two vectors  $\mathbf{x}$  and  $\mathbf{y}$ , it is  $0 \leq d(\mathbf{x}, \mathbf{y})$ . Therefore, in this case, it is  $d_0 = 0$ .

For (b): It is  $d(\mathbf{x}, \mathbf{x}) = \sqrt{\sum_{i=1}^l (x_i - x_i)^2} = 0 \equiv d_0$

For (c): It is  $d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^l (x_i - y_i)^2} = \sqrt{\sum_{i=1}^l (y_i - x_i)^2} = d(\mathbf{y}, \mathbf{x})$

For (d): It is  $d(\mathbf{x}, \mathbf{y}) = d_0 \Leftrightarrow \sqrt{\sum_{i=1}^l (x_i - y_i)^2} = 0 \Leftrightarrow x_i = y_i, \text{ for all } i = 1, \dots, l \Leftrightarrow \mathbf{x} = \mathbf{y}$

For (e): Here, the Minkowski inequality will be utilized. It is

$$\begin{aligned} d(\mathbf{x}, \mathbf{z}) &= \left( \sum_{i=1}^l |x_i - z_i|^2 \right)^{1/2} = \left( \sum_{i=1}^l |(x_i - y_i) + (y_i - z_i)|^2 \right)^{1/2} \\ &\leq \left( \sum_{i=1}^l |x_i - y_i|^2 \right)^{1/2} + \left( \sum_{i=1}^l |y_i - z_i|^2 \right)^{1/2} = d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z}) \end{aligned}$$

Therefore, since the Euclidean distance possesses all the properties from (a) to (e), it follows that it is a metric dissimilarity measure. Q.E.D.

**Exercise 2 (proximity measures):** Prove that the **distance** between for two  $l$ -dimensional vectors  $\mathbf{x} = [x_1, \dots, x_l]^T$  and  $\mathbf{y} = [y_1, \dots, y_l]^T$ , defined as  $d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^l (x_i^2 - y_i^2)^2}$  is a **dissimilarity measure** but **not a metric**.

**Solution:** In order to prove that  $d(\mathbf{x}, \mathbf{y})$  is a dissimilarity measure, we need to prove that it possesses all the following three properties:

(a)  $\exists d_0 \in \mathbb{R}$ , so that  $d_0 \leq d(\mathbf{x}, \mathbf{y}) < +\infty$ , for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^l$

(b)  $d(\mathbf{x}, \mathbf{x}) = d_0$ , for all  $\mathbf{x} \in \mathbb{R}^l$

(c)  $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$ , for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^l$

For (a): It is straightforward to see that for any two vectors  $\mathbf{x}$  and  $\mathbf{y}$ , it is  $0 \leq d(\mathbf{x}, \mathbf{y})$ . Therefore, in this case, it is  $d_0 = 0$ .

For (b): It is  $d(\mathbf{x}, \mathbf{x}) = \sqrt{\sum_{i=1}^l (x_i^2 - x_i^2)^2} = 0 \equiv d_0$

For (c): It is  $d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^l (x_i^2 - y_i^2)^2} = \sqrt{\sum_{i=1}^l (y_i^2 - x_i^2)^2} = d(\mathbf{y}, \mathbf{x})$

Therefore, since the distance  $d(\cdot, \cdot)$  possesses all the properties from (a) to (c), it follows that it is a dissimilarity measure.

In order to prove that it is not a metric, we need to prove that at least one of the following two properties are not possessed by the distance under study:

(d)  $d(\mathbf{x}, \mathbf{y}) = d_0 \Leftrightarrow \mathbf{x} = \mathbf{y}$

(e)  $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$ , for all  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^l$

For (d): It is  $d(\mathbf{x}, \mathbf{y}) = d_0 \Leftrightarrow \sqrt{\sum_{i=1}^l (x_i^2 - y_i^2)^2} = 0 \Leftrightarrow x_i^2 = y_i^2, \text{ for all } i = 1, \dots, l \Leftrightarrow x_i = \pm y_i, \text{ for all } i = 1, \dots, l$

This implies that there exist vectors that although they are different, they have the minimum possible distance value (for example, for  $l=3$ , the vectors  $\mathbf{x} = [2, 3, 4]^T$ ,  $\mathbf{y} = [2, -3, -4]^T$ , although they are different, their distance is equal to 0). Therefore the distance under study, is not a metric. Q.E.D.

**Exercise 3 (k-means cost function optimization algorithm):**

Consider the data set  $Y = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5\}$ , where  $\mathbf{x}_1 = [0, 0]^T$ ,  $\mathbf{x}_2 = [3, 0]^T$ ,  $\mathbf{x}_3 = [0, 6]^T$ ,  $\mathbf{x}_4 = [0, 7]^T$ ,  $\mathbf{x}_5 = [-3, 7]^T$ .

(a) Run the k-means clustering algorithm, for two representatives,  $\theta_1$  and  $\theta_2$ , whose initial positions are  $\theta_1(0) = [1, 6]^T$  and  $\theta_2(0) = [0, 8]^T$ , respectively. Report the formed clusters, along with their respective representatives.

- (b) What would be the clustering result for the case where  $\theta_2(0) = [0, 20]^T$ ?  
(c) How many clusters will be obtained if three representatives were employed?

**Solution:** We remind that, at each iteration of the k-means algorithm, two processing steps are involved. At the first one the data vectors are assigned to the clusters (each data vector  $\mathbf{x}$  is assigned to the cluster whose representative is closest to  $\mathbf{x}$ , in terms of the squared Euclidean distance) and at the second one each cluster representative  $\theta_j$  is re-estimated (as the mean of the data vectors that belong to the respective cluster).

Based on this, we proceed as follows:

### 1<sup>st</sup> iteration

1<sup>st</sup> step: The **squared Euclidean distances** of each data point from the two representatives are given in the following table

	$\theta_1(0) = [1, 6]^T$	$\theta_2(0) = [0, 8]^T$
$\mathbf{x}_1 = [0, 0]^T$	$(0 - 1)^2 + (0 - 6)^2 = \mathbf{37}$	$(0 - 0)^2 + (0 - 8)^2 = \mathbf{64}$
$\mathbf{x}_2 = [3, 0]^T$	$(3 - 1)^2 + (0 - 6)^2 = \mathbf{40}$	$(3 - 0)^2 + (0 - 8)^2 = \mathbf{73}$
$\mathbf{x}_3 = [0, 6]^T$	$(0 - 1)^2 + (6 - 6)^2 = \mathbf{1}$	$(0 - 0)^2 + (6 - 8)^2 = \mathbf{4}$
$\mathbf{x}_4 = [0, 7]^T$	$(0 - 1)^2 + (7 - 6)^2 = \mathbf{2}$	$(0 - 0)^2 + (7 - 8)^2 = \mathbf{1}$
$\mathbf{x}_5 = [-3, 7]^T$	$(-3 - 1)^2 + (7 - 6)^2 = \mathbf{17}$	$(-3 - 0)^2 + (7 - 8)^2 = \mathbf{10}$

Since:

$$\|\mathbf{x}_1 - \theta_1(0)\|^2 = \min_{j=1,2} \|\mathbf{x}_1 - \theta_j(0)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_1.^1$$

$$\|\mathbf{x}_2 - \theta_1(0)\|^2 = \min_{j=1,2} \|\mathbf{x}_2 - \theta_j(0)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_1.$$

$$\|\mathbf{x}_3 - \theta_1(0)\|^2 = \min_{j=1,2} \|\mathbf{x}_3 - \theta_j(0)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_1.$$

$$\|\mathbf{x}_4 - \theta_2(0)\|^2 = \min_{j=1,2} \|\mathbf{x}_4 - \theta_j(0)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_2.$$

$$\|\mathbf{x}_5 - \theta_2(0)\|^2 = \min_{j=1,2} \|\mathbf{x}_5 - \theta_j(0)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_2.$$

Thus,  $C_1 = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$  and  $C_2 = \{\mathbf{x}_4, \mathbf{x}_5\}$ .

2<sup>nd</sup> step: Letting  $n_1 = 3$  and  $n_2 = 2$  be the cardinalities of  $C_1$  and  $C_2$ , respectively, the cluster representatives are re-estimated as

$$\theta_1 \equiv \theta_1(1) = \frac{1}{n_1} \cdot (\mathbf{x}_1 + \mathbf{x}_2 + \mathbf{x}_3) = \frac{1}{3} \cdot \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 3 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 6 \end{bmatrix} \right) = \frac{1}{3} \cdot \begin{bmatrix} 3 \\ 6 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

and

<sup>1</sup> Note that cluster  $C_1$  ( $C_2$ ) is associated with the representative  $\theta_1$  ( $\theta_2$ ).

$$\theta_2 \equiv \theta_2(1) = \frac{1}{n_2} \cdot (\mathbf{x}_4 + \mathbf{x}_5) = \frac{1}{2} \cdot \left( \begin{bmatrix} 0 \\ 7 \end{bmatrix} + \begin{bmatrix} -3 \\ 7 \end{bmatrix} \right) = \frac{1}{2} \cdot \begin{bmatrix} -3 \\ 14 \end{bmatrix} = \begin{bmatrix} -1.5 \\ 7 \end{bmatrix}$$

### 2<sup>nd</sup> iteration

1<sup>st</sup> step: The squared Euclidean distances of each data point from the two representatives are given in the following table

	$\theta_1(1) = [1, 2]^T$	$\theta_2(1) = [-1.5, 7]^T$
$\mathbf{x}_1 = [0, 0]^T$	$(0 - 1)^2 + (0 - 2)^2 = \mathbf{5}$	$(0 - (-1.5))^2 + (0 - 7)^2 = \mathbf{51.25}$
$\mathbf{x}_2 = [3, 0]^T$	$(3 - 1)^2 + (0 - 2)^2 = \mathbf{8}$	$(3 - (-1.5))^2 + (0 - 7)^2 = \mathbf{69.25}$
$\mathbf{x}_3 = [0, 6]^T$	$(0 - 1)^2 + (6 - 2)^2 = \mathbf{17}$	$(0 - (-1.5))^2 + (6 - 7)^2 = \mathbf{3.25}$
$\mathbf{x}_4 = [0, 7]^T$	$(0 - 1)^2 + (7 - 2)^2 = \mathbf{26}$	$(0 - (-1.5))^2 + (7 - 7)^2 = \mathbf{2.25}$
$\mathbf{x}_5 = [-3, 7]^T$	$(-3 - 1)^2 + (7 - 2)^2 = \mathbf{41}$	$(-3 - (-1.5))^2 + (7 - 7)^2 = \mathbf{2.25}$

Since:

$$\|\mathbf{x}_1 - \theta_1(1)\|^2 = \min_{j=1,2} \|\mathbf{x}_1 - \theta_j(1)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_1.$$

$$\|\mathbf{x}_2 - \theta_1(1)\|^2 = \min_{j=1,2} \|\mathbf{x}_2 - \theta_j(1)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_1.$$

$$\|\mathbf{x}_3 - \theta_2(1)\|^2 = \min_{j=1,2} \|\mathbf{x}_3 - \theta_j(1)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_2.$$

$$\|\mathbf{x}_4 - \theta_2(1)\|^2 = \min_{j=1,2} \|\mathbf{x}_4 - \theta_j(1)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_2.$$

$$\|\mathbf{x}_5 - \theta_2(1)\|^2 = \min_{j=1,2} \|\mathbf{x}_5 - \theta_j(1)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_2.$$

Thus,  $C_1 = \{\mathbf{x}_1, \mathbf{x}_2\}$  and  $C_2 = \{\mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5\}$ .

2<sup>nd</sup> step: Noting that the cardinalities of  $C_1$  and  $C_2$  are  $n_1 = 2$  and  $n_2 = 3$ , respectively, the cluster representatives are re-estimated as

$$\theta_1 \equiv \theta_1(2) = \frac{1}{n_1} \cdot (\mathbf{x}_1 + \mathbf{x}_2) = \frac{1}{2} \cdot \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 3 \\ 0 \end{bmatrix} \right) = \frac{1}{2} \cdot \begin{bmatrix} 3 \\ 0 \end{bmatrix} = \begin{bmatrix} 1.5 \\ 0 \end{bmatrix}$$

and

$$\theta_2 \equiv \theta_2(2) = \frac{1}{n_2} \cdot (\mathbf{x}_3 + \mathbf{x}_4 + \mathbf{x}_5) = \frac{1}{3} \cdot \left( \begin{bmatrix} 0 \\ 6 \end{bmatrix} + \begin{bmatrix} 0 \\ 7 \end{bmatrix} + \begin{bmatrix} -3 \\ 7 \end{bmatrix} \right) = \frac{1}{3} \cdot \begin{bmatrix} -3 \\ 20 \end{bmatrix} = \begin{bmatrix} -1 \\ 6.7 \end{bmatrix}$$

### 3<sup>rd</sup> iteration

1<sup>st</sup> step: The squared Euclidean distances of each data point from the two representatives are given in the following table

	$\theta_1(2) = [1.5, 0]^T$	$\theta_2(2) = [-1, 6.7]^T$
$\mathbf{x}_1 = [0, 0]^T$	$(0 - 1.5)^2 + (0 - 0)^2 = \mathbf{2.25}$	$(0 - (-1))^2 + (0 - 6.7)^2 = \mathbf{45.89}$
$\mathbf{x}_2 = [3, 0]^T$	$(3 - 1.5)^2 + (0 - 0)^2 = \mathbf{2.25}$	$(3 - (-1))^2 + (0 - 6.7)^2 = \mathbf{60.89}$
$\mathbf{x}_3 = [0, 6]^T$	$(0 - 1.5)^2 + (6 - 0)^2 = \mathbf{38.25}$	$(0 - (-1))^2 + (6 - 6.7)^2 = \mathbf{1.49}$
$\mathbf{x}_4 = [0, 7]^T$	$(0 - 1.5)^2 + (7 - 0)^2 = \mathbf{51.25}$	$(0 - (-1))^2 + (7 - 6.7)^2 = \mathbf{1.09}$
$\mathbf{x}_5 = [-3, 7]^T$	$(-3 - 1.5)^2 + (7 - 0)^2 = \mathbf{69.25}$	$(-3 - (-1))^2 + (7 - 6.7)^2 = \mathbf{4.09}$

Since:

$$\|\mathbf{x}_1 - \theta_1(2)\|^2 = \min_{j=1,2} \|\mathbf{x}_1 - \theta_j(2)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_1.$$

$$\|\mathbf{x}_2 - \theta_1(2)\|^2 = \min_{j=1,2} \|\mathbf{x}_2 - \theta_j(2)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_1.$$

$$\|\mathbf{x}_3 - \theta_2(2)\|^2 = \min_{j=1,2} \|\mathbf{x}_3 - \theta_j(2)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_2.$$

$$\|\mathbf{x}_4 - \theta_2(2)\|^2 = \min_{j=1,2} \|\mathbf{x}_4 - \theta_j(2)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_2.$$

$$\|\mathbf{x}_5 - \theta_2(2)\|^2 = \min_{j=1,2} \|\mathbf{x}_5 - \theta_j(2)\|^2, \mathbf{x}_1 \text{ is assigned to cluster } C_2.$$

Thus,  $C_1 = \{\mathbf{x}_1, \mathbf{x}_2\}$  and  $C_2 = \{\mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5\}$ .

2<sup>nd</sup> step: Noting that the cardinalities of  $C_1$  and  $C_2$  are  $n_1 = 2$  and  $n_2 = 3$ , respectively, the cluster representatives are re-estimated as

$$\theta_1 \equiv \theta_1(3) = \frac{1}{n_1} \cdot (\mathbf{x}_1 + \mathbf{x}_2) = \frac{1}{2} \cdot \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 3 \\ 0 \end{bmatrix} \right) = \frac{1}{2} \cdot \begin{bmatrix} 3 \\ 0 \end{bmatrix} = \begin{bmatrix} 1.5 \\ 0 \end{bmatrix}$$

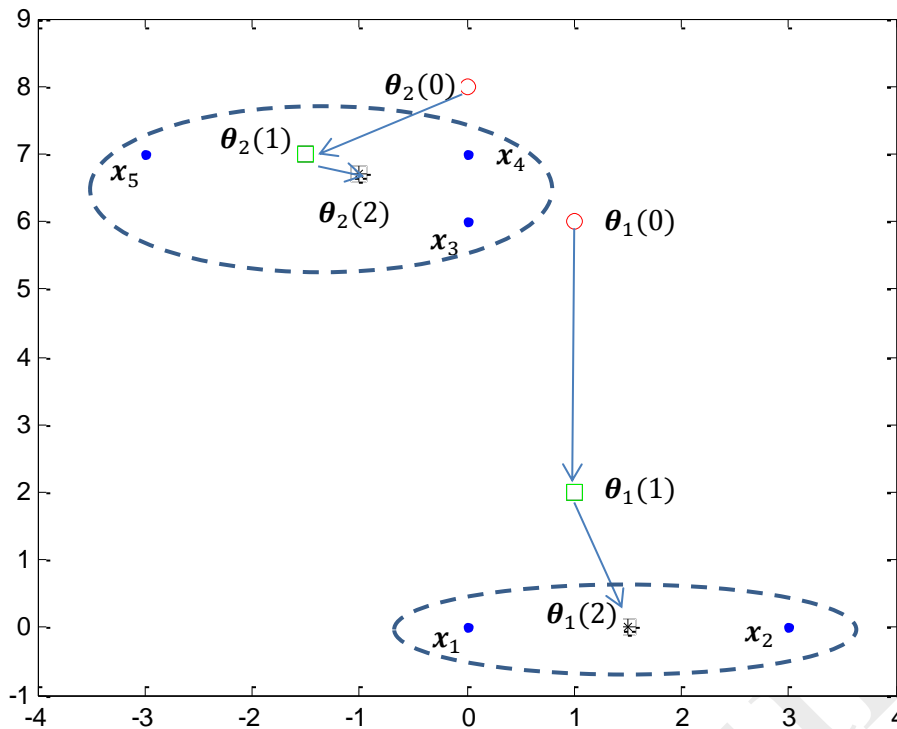
and

$$\theta_2 \equiv \theta_2(3) = \frac{1}{n_2} \cdot (\mathbf{x}_3 + \mathbf{x}_4 + \mathbf{x}_5) = \frac{1}{3} \cdot \left( \begin{bmatrix} 0 \\ 6 \end{bmatrix} + \begin{bmatrix} 0 \\ 7 \end{bmatrix} + \begin{bmatrix} -3 \\ 7 \end{bmatrix} \right) = \frac{1}{3} \cdot \begin{bmatrix} -3 \\ 20 \end{bmatrix} = \begin{bmatrix} -1 \\ 20/3 \end{bmatrix}.$$

Since the values of  $\theta_j$ 's,  $j = 1, 2$ , remain unaltered for two successive iterations (2<sup>nd</sup> and 3<sup>rd</sup>), the algorithm terminates (equivalently, we can say that since the clustering of the data points remain unaltered for two successive iterations, the algorithm terminates). The resulting clusters are  $C_1 = \{\mathbf{x}_1, \mathbf{x}_2\}$  and  $C_2 =$

$\{\mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5\}$  and the respective cluster representatives are  $\theta_1 = \begin{bmatrix} 1.5 \\ 0 \end{bmatrix}$  and  $\theta_2 = \begin{bmatrix} -1 \\ 20/3 \end{bmatrix}$ .

The evolution of the algorithm is shown in the following figure.



(b) Applying the k-means algorithm for this initialization scenario, we will see that all data points lie closer to  $\theta_1(0)$  than  $\theta_2(0)$  (alternatively, we can say that  $\theta_2(0)$  does not “win” on any data point). Thus, the k-means will return a single cluster containing all the data points, or, strictly speaking,  $C_1 = \{x_1, x_2, x_3, x_4, x_5\}$  and  $C_2 = \emptyset$ . Cluster  $C_1$  is represented by

$$\theta_1 \equiv \theta_1(1) = \frac{1}{n_1} \cdot (x_1 + x_2 + x_3 + x_4 + x_5) = \frac{1}{5} \cdot \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 3 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 6 \end{bmatrix} + \begin{bmatrix} 0 \\ 7 \end{bmatrix} + \begin{bmatrix} -3 \\ 7 \end{bmatrix} \right) = \frac{1}{5} \cdot \begin{bmatrix} 0 \\ 20 \end{bmatrix} = \begin{bmatrix} 0 \\ 4 \end{bmatrix}.$$

(d) In general, the data set will be split into three clusters, provided that each representative “wins” on at least one data point.

**Exercise 4 (Matrix theory-based hierarchical clustering):** Consider the following dissimilarity matrix:

$$P_0 = \begin{bmatrix} 0 & 1 & 2 & 26 & 37 \\ 1 & 0 & 3 & 25 & 36 \\ 2 & 3 & 0 & 16 & 25 \\ 26 & 25 & 16 & 0 & 1.5 \\ 37 & 36 & 25 & 1.5 & 0 \end{bmatrix}$$

where the corresponding squared Euclidean distance is adopted. Run the seven agglomerative matrix-based clustering algorithms and determine the resulting clustering hierarchy, as well as the dissimilarity levels where the clusterings are produced.

**Solution:** As one can easily observe, the first three vectors,  $\mathbf{x}_1$ ,  $\mathbf{x}_2$ , and  $\mathbf{x}_2$ , are very close to each other and far away from the others. Likewise,  $\mathbf{x}_4$  and  $\mathbf{x}_5$  lie very close to each other and far away from the first three vectors.

For this problem all seven algorithms discussed before result in the same dendrogram. The only difference is that each clustering is formed at a different dissimilarity level. Of course, the initial clustering is  $\mathfrak{R}_0 = \{\{\mathbf{x}_1\}, \{\mathbf{x}_2\}, \{\mathbf{x}_3\}, \{\mathbf{x}_4\}, \{\mathbf{x}_5\}\}$ . Let us first consider the single link algorithm. Since  $P_0$  is symmetric, we consider only the upper diagonal elements. The smallest of these elements equals 1 and occurs at position (1,2) of  $P_0$ . Thus,  $\mathbf{x}_1$  and  $\mathbf{x}_2$  come into the same cluster and  $\mathfrak{R}_1 = \{\{\mathbf{x}_1, \mathbf{x}_2\}, \{\mathbf{x}_3\}, \{\mathbf{x}_4\}, \{\mathbf{x}_5\}\}$  is produced. In the sequel, the dissimilarities among the newly formed cluster and the remaining ones have to be computed. This can be achieved via Eq. (2)<sup>2</sup>. The resulting proximity matrix,  $P_1$ , is

$$P_1 = \begin{bmatrix} 0 & 2 & 25 & 36 \\ 2 & 0 & 16 & 25 \\ 25 & 16 & 0 & 1.5 \\ 36 & 25 & 1.5 & 0 \end{bmatrix}$$

Its first row and column correspond to the cluster  $\{\mathbf{x}_1, \mathbf{x}_2\}$ . The smallest of the upper diagonal elements of  $P_1$  equals 1.5. This means that at the next stage, the clusters  $\{\mathbf{x}_4\}$  and  $\{\mathbf{x}_5\}$  will stick together into a single cluster, producing  $\mathfrak{R}_2 = \{\{\mathbf{x}_1, \mathbf{x}_2\}, \{\mathbf{x}_3\}, \{\mathbf{x}_4, \mathbf{x}_5\}\}$ . Employing Eq. (13.4), we obtain

$$P_2 = \begin{bmatrix} 0 & 2 & 25 \\ 2 & 0 & 16 \\ 25 & 16 & 0 \end{bmatrix}$$

where the first row (column) corresponds to  $\{\mathbf{x}_1, \mathbf{x}_2\}$ , and the second and third rows (columns) correspond to  $\{\mathbf{x}_3\}$  and  $\{\mathbf{x}_4, \mathbf{x}_5\}$ , respectively. Proceeding as before, at the next stage  $\{\mathbf{x}_1, \mathbf{x}_2\}$  and  $\{\mathbf{x}_3\}$  will get together in a single cluster and  $\mathfrak{R}_3 = \{\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}, \{\mathbf{x}_4, \mathbf{x}_5\}\}$  is produced. The new proximity matrix,  $P_3$ , becomes

$$P_3 = \begin{bmatrix} 0 & 16 \\ 16 & 0 \end{bmatrix}$$

where the first and the second row (column) correspond to  $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$  and  $\{\mathbf{x}_4, \mathbf{x}_5\}$  clusters, respectively. Finally,  $\mathfrak{R}_4 = \{\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5\}\}$  will be formed at dissimilarity level equal to 16.

Working in a similar fashion, we can apply the remaining six algorithms to  $P_0$ .<sup>3</sup>

However, care must be taken when we apply UPGMA, UPGMC, and Ward's method. In these cases, when a merging takes place the parameters  $a_i$ ,  $a_j$ ,  $b$ , and  $c$  in

$$d(C_q, C_s) = a_i d(C_i, C_s) + a_j (d(C_j, C_s) + b d(C_i, C_j) + c |d(C_i, C_s) - d(C_j, C_s)|)$$

must be properly adjusted. The proximity levels at which each clustering is formed for each algorithm are shown in Table 1.

<sup>2</sup> All references are referred to the slides of the 9<sup>th</sup> lecture.

<sup>3</sup> Note that in the case of Ward's algorithm, the initial dissimilarity matrix should be  $12 \cdot P_0$ , due to the definition of the distance  $d'_{ij}$  between  $C_i$  and  $C_j$  is defined as

$$d'_{ij} = \frac{n_i n_j}{n_i + n_j} \|\mathbf{m}_i - \mathbf{m}_j\|^2$$

where  $\mathbf{m}_i$  and  $\mathbf{m}_j$  are the mean vectors associated with  $C_i$  and  $C_j$ , respectively.

The considered task is a nice problem with two well-defined compact clusters lying away from each other. The preceding example demonstrates that in such “easy” cases all algorithms work satisfactorily (as happens with most of the clustering algorithms proposed in the literature). The particular characteristics of each algorithm are revealed when more demanding situations are faced.

	<i>SL</i>	<i>CL</i>	<i>WPGMA</i>	<i>UPGMA</i>	<i>WPGMC</i>	<i>UPGMC</i>	<i>Ward</i>
$\mathcal{R}_0$	0	0	0	0	0	0	0
$\mathcal{R}_1$	1	1	1	1	1	1	0.5
$\mathcal{R}_2$	1.5	1.5	1.5	1.5	1.5	1.5	0.75
$\mathcal{R}_3$	2	3	2.5	2.5	2.25	2.25	1.5
$\mathcal{R}_4$	16	37	25.75	27.5	24.69	26.46	31.75

**Table 4.1:** The Results Obtained with the Seven Algorithms Discussed when they are applied to the proximity matrix of Example 1.

**Exercise 5** (*graph theory-based hierarchical clustering*): Consider the following dissimilarity matrix:

$$P(X) = \begin{bmatrix} 0 & 1 & 9 & 18 & 19 & 20 & 21 \\ 1 & 0 & 8 & 13 & 14 & 15 & 16 \\ 9 & 8 & 0 & 17 & 10 & 11 & 12 \\ 18 & 13 & 17 & 0 & 5 & 6 & 7 \\ 19 & 14 & 10 & 5 & 0 & 3 & 4 \\ 20 & 15 & 11 & 6 & 3 & 0 & 2 \\ 21 & 16 & 12 & 7 & 4 & 2 & 0 \end{bmatrix}$$

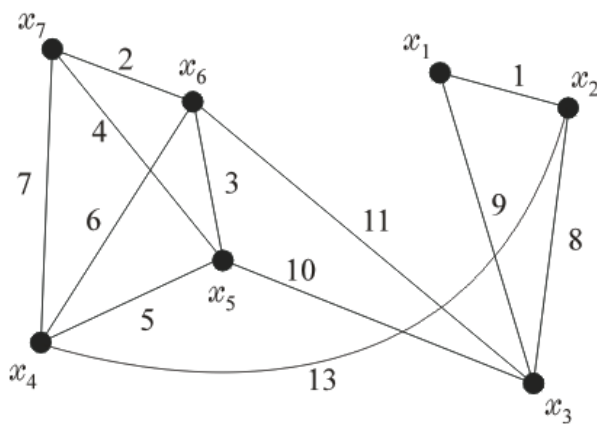
Let  $h(k)$  be the **node degree property** with  $k = 2$ ; that is, it is required that each node has at least two incident edges. Derive the associated dissimilarity dendrogram.

Figure 1 shows the  $G(13)$  proximity graph produced by this dissimilarity matrix. Then the obtained threshold dendrogram is shown in Figure 2a. At dissimilarity level 1,  $\mathbf{x}_1$  and  $\mathbf{x}_2$  form a single cluster. This happens because  $\{\mathbf{x}_1\} \cup \{\mathbf{x}_2\}$  is complete at  $G(1)$ , despite the fact that property  $h(2)$  is not satisfied (remember the disjunction between conditions (b1) and (b2) in Eq. (4)). Similarly,  $\{\mathbf{x}_6\} \cup \{\mathbf{x}_7\}$  forms a cluster at dissimilarity level 2. The next clustering is formed at level 4, since  $\{\mathbf{x}_5\} \cup \{\mathbf{x}_6, \mathbf{x}_7\}$  becomes complete in  $G(4)$ .

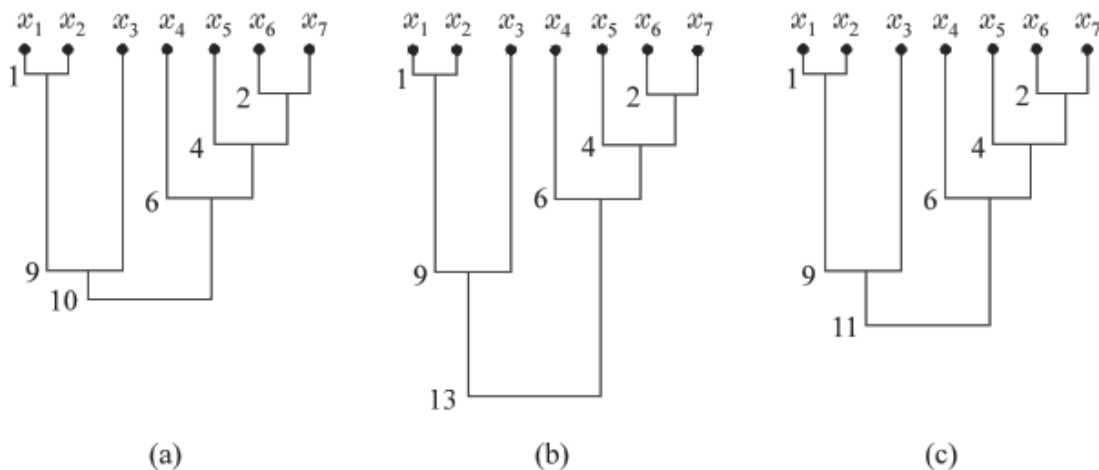
At level 6,  $\mathbf{x}_4$ ,  $\mathbf{x}_5$ ,  $\mathbf{x}_6$ , and  $\mathbf{x}_7$  lie for the first time in the same cluster. Although this subgraph is not complete, it does satisfy  $h(2)$ . Finally, at level 9,  $\mathbf{x}_1$ ,  $\mathbf{x}_2$ , and  $\mathbf{x}_3$  come into the same cluster. Note that, although all nodes in the graph have node degree equal to 2, the final clustering will be formed at level 10 because at level 9 the graph is not connected. Assume now that  $h(k)$  is the node connectivity property, with



$k = 2$ ; that is, all pairs of nodes in a connected subgraph are joined by at least two paths having no nodes in common. The dissimilarity dendrogram produced in this case is shown in Figure 2b. Finally, the dissimilarity dendrogram produced when the edge connectivity property with  $k = 2$  is employed is shown in Figure 2c.



**Figure 5.1:** The proximity graph  $G(13)$  derived by the dissimilarity matrix  $P$  given in Exercise 5.



**Figure 5.2:** Dissimilarity dendrograms related to Exercise 5. (a) Dissimilarity dendrogram produced when  $h(k)$  is the node degree property, with  $k = 2$ . (b) Dissimilarity dendrogram produced when  $h(k)$  is the node connectivity property, with  $k = 2$ . (c) Dissimilarity dendrogram produced when  $h(k)$  is the edge connectivity property, with  $k = 2$ .

**Exercise 6 (Spectral clustering):** Consider the data set  $X = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5\}$ , where  $\mathbf{x}_1 = [0, 0]^T$ ,  $\mathbf{x}_2 = [0, 1]^T$ ,  $\mathbf{x}_3 = [5, 0]^T$ ,  $\mathbf{x}_4 = [5, 1]^T$ ,  $\mathbf{x}_5 = [4, 1]^T$ . Perform spectral clustering using the 1-NN for the construction of the similarity graph (initial phase). The weight of the edges in the graph will be computed via the equation  $s(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2)$ .

**Solution: Construction of the similarity matrix:** It is  $s(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\left\|\begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 1 \end{bmatrix}\right\|^2\right) = \exp(-1) = 0.3679 \approx 0.4$ . In the same spirit, we compute the similarities between any pair of points and we end up with the following similarity matrix<sup>4</sup>

$$S = \begin{bmatrix} 1.0 & 0.4 & 0.0 & 0.0 & 0.0 \\ 0.4 & 1.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 1.0 & 0.4 & 0.1 \\ 0.0 & 0.0 & 0.4 & 1.0 & 0.4 \\ 0.0 & 0.0 & 0.1 & 0.4 & 1.0 \end{bmatrix}$$

**Construction of the similarity graph:** From the first row of matrix  $S$ , it follows that the nearest neighbor of  $\mathbf{x}_1$  is  $\mathbf{x}_2$ ,

the nearest neighbor of  $\mathbf{x}_2$  is  $\mathbf{x}_1$ ,

the nearest neighbor of  $\mathbf{x}_3$  is  $\mathbf{x}_4$ ,

the nearest neighbor of  $\mathbf{x}_4$  is  $\mathbf{x}_3$  (we could take  $\mathbf{x}_5$  instead),

the nearest neighbor of  $\mathbf{x}_5$  is  $\mathbf{x}_4$ .

Thus, the similarity graph  $G = (V, E)$ , consists of five vertices, i.e.,  $V = \{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4, \mathbf{v}_5\}$ , each one corresponding to a data point, while the set of edges is  $E = \{e_{12}, e_{34}, e_{45}\}$ . The associated **weighted adjacency matrix** is

$$W = \begin{bmatrix} 1.0 & 0.4 & 0.0 & 0.0 & 0.0 \\ 0.4 & 1.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 1.0 & 0.4 & 0.0 \\ 0.0 & 0.0 & 0.4 & 1.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.4 & 1.0 \end{bmatrix}$$

**Construction of the Laplacian matrix:** The degree matrix and the (unnormalized) Laplacian matrices are, respectively

$$D = \begin{bmatrix} 1.4 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 1.4 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 1.4 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 1.4 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 1.4 \end{bmatrix} \text{ and } L = D - W = \begin{bmatrix} +0.4 & -0.4 & 0.0 & 0.0 & 0.0 \\ -0.4 & +0.4 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & +0.4 & -0.4 & 0.0 \\ 0.0 & 0.0 & -0.4 & +0.4 & 0.0 \\ 0.0 & 0.0 & 0.0 & -0.4 & +0.4 \end{bmatrix}$$

**Determination of the zero eigenvalues and the associated eigenvectors of  $L$ :** It is

$$\begin{aligned} \det(L - \lambda I) &= \begin{vmatrix} 0.4 - \lambda & -0.4 & 0.0 & 0.0 & 0.0 \\ -0.4 & 0.4 - \lambda & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.4 - \lambda & -0.4 & 0.0 \\ 0.0 & 0.0 & -0.4 & 0.4 - \lambda & 0.0 \\ 0.0 & 0.0 & 0.0 & -0.4 & 0.4 - \lambda \end{vmatrix} = -\lambda \begin{vmatrix} 1 & -0.4 & 0.0 & 0.0 & 0.0 \\ 1 & 0.4 - \lambda & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.4 - \lambda & -0.4 & 0.0 \\ 0.0 & 0.0 & -0.4 & 0.4 - \lambda & 0.0 \\ 0.0 & 0.0 & 0.0 & -0.4 & 0.4 - \lambda \end{vmatrix} \\ &= (-\lambda)^2 \begin{vmatrix} 1 & -0.4 & 0.0 & 0.0 & 0.0 \\ 1 & 0.4 - \lambda & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 1 & -0.4 & 0.0 \\ 0.0 & 0.0 & 1 & 0.4 - \lambda & 0.0 \\ 0.0 & 0.0 & 1 & -0.4 & 0.4 - \lambda \end{vmatrix} \end{aligned}$$

<sup>4</sup> The precision is up to the first decimal.

The equation  $\det(L - \lambda I) = 0$  has **two zero** eigenvalues<sup>5</sup>.

The eigenvectors  $\mathbf{x}$  corresponding to a zero eigenvalue of  $L$  should satisfy  $L\mathbf{x} = 0\mathbf{x}$ . It is clear that for the vectors  $\mathbf{u}_1 = [1,1,0,0,0]^T$ ,  $\mathbf{u}_2 = [0,0,1,1,1]^T$  it is

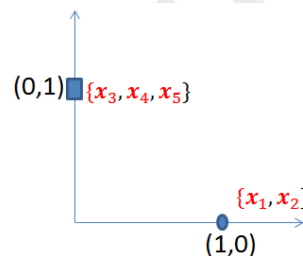
$$L \cdot \mathbf{u}_1 = \begin{bmatrix} +0.4 & -0.4 & 0.0 & 0.0 & 0.0 \\ -0.4 & +0.4 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & +0.4 & -0.4 & 0.0 \\ 0.0 & 0.0 & -0.4 & +0.4 & 0.0 \\ 0.0 & 0.0 & 0.0 & -0.4 & +0.4 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = 0 \cdot \mathbf{u}_1$$

$$L \cdot \mathbf{u}_2 = \begin{bmatrix} +0.4 & -0.4 & 0.0 & 0.0 & 0.0 \\ -0.4 & +0.4 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & +0.4 & -0.4 & 0.0 \\ 0.0 & 0.0 & -0.4 & +0.4 & 0.0 \\ 0.0 & 0.0 & 0.0 & -0.4 & +0.4 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = 0 \cdot \mathbf{u}_2$$

Thus, these are the eigenvectors associated with the zero eigenvalues.

Construction of the U matrix: It is

$$U = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \equiv \begin{matrix} \mathbf{y}_1 \rightarrow \mathbf{x}_1 \\ \mathbf{y}_2 \rightarrow \mathbf{x}_2 \\ \mathbf{y}_3 \rightarrow \mathbf{x}_3 \\ \mathbf{y}_4 \rightarrow \mathbf{x}_4 \\ \mathbf{y}_5 \rightarrow \mathbf{x}_5 \end{matrix}$$



The  $i$ -th data vector is mapped to a vector in a new two-dim. space whose coordinates are the  $i$ -th associated coordinates of the two eigenvectors. Thus, the final clustering in the transformed space consists of the clusters  $C'_1 = \{\mathbf{y}_1, \mathbf{y}_2\}$  and  $C'_2 = \{\mathbf{y}_3, \mathbf{y}_4, \mathbf{y}_5\}$ . Thus the clustering of the original data consists of the clusters  $C_1 = \{\mathbf{x}_1, \mathbf{x}_2\}$  and  $C_2 = \{\mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5\}$ .

### Indicative questions

**Question 1:** Consider a clustering task where the involved  $N$  entities are represented in a two-dimensional feature space associated with the real-valued features  $x_1$  and  $x_2$ . The  $x_1$  values of the entities are ranged in  $[0,1]$ , while the  $x_2$  values of the entities are ranged in  $[0,1000]$  (assume also that the  $N$   $x_1$  values are uniformly arranged in  $[0,1]$  and the  $N$   $x_2$  values are uniformly arranged in  $[0,1000]$ ). Propose a transformation of the original feature space, so that the distance between any two vectors to be equally influenced by both feature values.

Hint: If  $a \leq x \leq b$ , then  $c \leq c + \frac{x-a}{b-a}(d-c) \leq d$

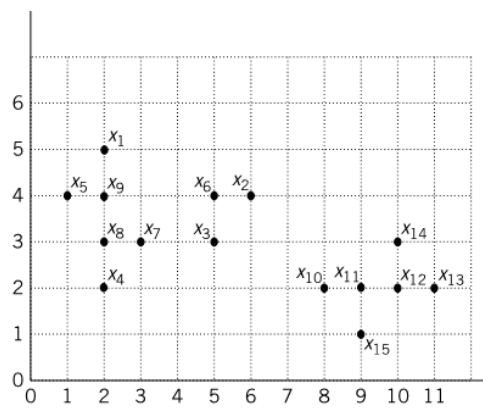
**Note 1:** You should be able to compute the overlap distance between two discrete-valued data vectors.

<sup>5</sup> The determinant left has no additional zero eigenvalues, since if we set  $\lambda=0$  to it, all the columns of the determinant are independent.

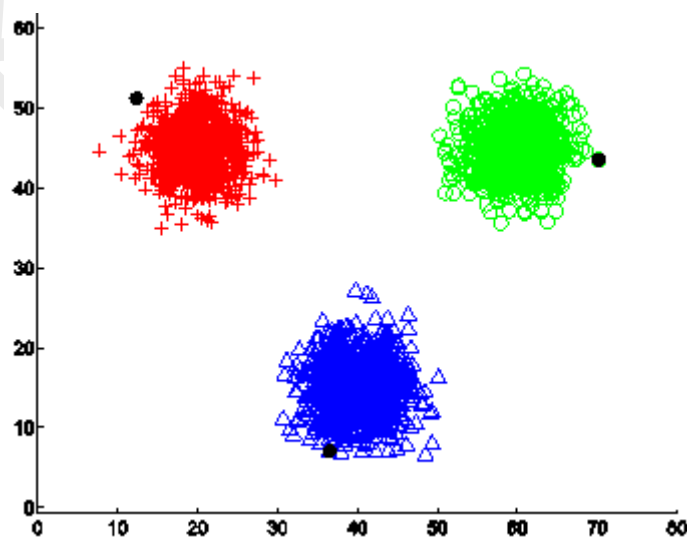
**Question 2:** What are the differences between the mean vector, the mean center and the median center of cluster?

**Question 3:** In the case where elongated clusters are formed by the data vectors of a data set, what kind of representatives you would use in the establishment a relevant clustering algorithm? Is there any case where the adoption of point representatives would work satisfactory in this case?

**Question 4:** In the data set the is depicted graphically below run the BSAS algorithm for  $\theta = 2$  and  $q = 2$  starting from  $x_1$  and then from  $x_{15}$ . What is the resulting clusterings for the two cases? Are they identical?



**Question 5:** Consider the following data set. What will be the result of (a) the k-means, (b) the fuzzy c-means and (c) the possibilistic algorithms for (a) three, (b) four and (c) five representatives, starting from different initial conditions?



**Question 6:** Derive a cost function optimization clustering algorithm based on a certain cost function  $J(\theta, U)$ , where  $\theta$  is the set of the parameter vectors associated with the clusters and  $U$  is a matrix whose  $(i,j)$  quantifies the relationship between the  $i$ -th vector and the  $j$ -th cluster.

*Hint:* An iterative two-step procedure will be followed: At the first step  $U$  will be updated for fixed  $\theta$  and at the second step  $\theta$  will be updated for fixed  $U$ . The kind of representatives will be decided based on the shape of the clusters formed by the data.

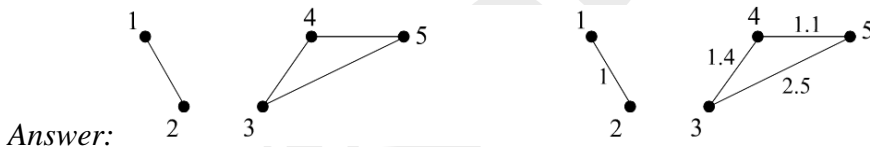
**Question 7:** Propose a way for estimating the true number of clusters, using an algorithm that (a) does not take as input the number of clusters and (b) it does require knowledge of the number of clusters.

*Hint:* Remember the plots with (a) the large “plateau” and (b) the significant “knee”.

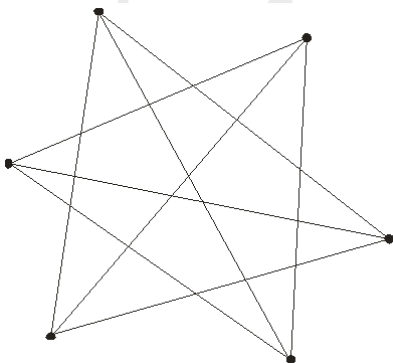
**Question 8:** Is it possible for the k-medoids algorithm to be derived using tools from mathematical analysis (e.g., gradients etc)?

**Question 9:** Determine the threshold graph  $G(3)$  and the dissimilarity graph  $G_p(3)$  for the data set whose associated proximity matrix is

$$P(X) = \begin{bmatrix} 0 & 1 & 5 & 6.4 & 7.4 \\ 1 & 0 & 4.2 & 5.7 & 6.7 \\ 5 & 4.2 & 0 & 1.4 & 2.5 \\ 6.4 & 5.7 & 1.4 & 0 & 1.1 \\ 7.4 & 6.7 & 2.5 & 1.1 & 0 \end{bmatrix}$$



**Question 10:** Determine the node connectivity, the edge connectivity and the node degree of the graph.



*Hint:* See the slides.

**Question 11:** Propose ways for determining the natural clustering from a hierarchy of clusterings that best describes the data.

*Hint:* See slides (e.g., intrinsic and extrinsic methods)

**Question 12:** What is the main approach in dealing with large data sets, in the hierarchical clustering algorithms case?

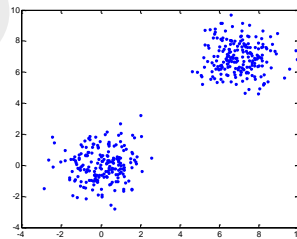
*Hint:* Sampling.

**Question 13:** If the neighboring edges of a given edge (of weight 30) in the MST associated with a data set have weights 2, 3, 5, 1 and  $q = 2$ , is the above edge “unusually large”?

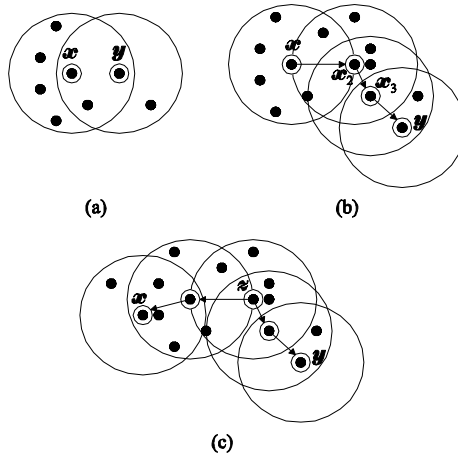
*Hint:* See slides (graph-algorithm based on the MST).

**Question 14:** How many data points are in the region of influence defined by two data points (in the algorithm based on regions of influence).

**Question 15:** Consider the two-cluster task shown below. Consider the basic competitive and the leaky learning scheme and assume that two representatives are considered. The first one lies in the middle of the two clusters and the other one far away from both. What will be the result of each one of the clusters?



**Question 16:** Is the point  $x$  (a) directly density reachable, (b) density reachable and (c) density connected with  $y$  (DBSCAN  $q = 5$ )?



**Question 17:** We apply different clustering algorithms on a certain data set and we end up with  $K$  different clusterings. We would like to combine all of them in order to take a single representative clustering. Define the corresponding co-association matrix and run a hierarchical algorithm.

**Question 18:** Subspace clustering algorithms have been designed for identifying clusters that live in the same subspace of the original feature space (yes/no).

Answer: No

**Question 19:** Consider the Rand measure that measures the agreement between two clustering structures. When it takes its minimum and maximum values and what are they?

Answer: Maximum value (1) – perfect match, Minimum value (0) – perfect mismatch