

Module: Data Mining  
 Nature de document: TD2-Solution  
 Niveau: L3-STID

Année: 2022-2023

M<sup>me</sup> N.BERMAD

**Exercice1:**

**1. L'algorithme des centres mobiles**

	$X_1$	$X_2$
$W_1$	-2	2
$W_2$	-2	-1
$W_3$	0	-1
$W_4$	2	2
$W_5$	-2	3
$W_6$	3	0

**Etape0 :**

La matrice des distances entre les individus et les centres  $C_1^0$  de coordonnées (1, 1) et  $C_2^0$  de coordonnées (2,3) est :

	$W_1$	$W_2$	$W_3$	$W_4$	$W_5$	$W_6$
$C_1^0$	3.16	3.16	2.24	1.41	3.60	2.24
$C_2^0$	4.12	5.66	4.47	1	4	3.16

En utilisant la distance euclidienne :

$$d(w_1, c_1^0) = \sqrt{(-2 - 1)^2 + (2 - 1)^2} = 3.16$$

D'où les deux groupes :

$$A = \{W_1, W_2, W_3, W_6, W_5\}, B = \{W_4\}$$

**Etape01:**

On considère deux nouveaux centres  $C_1^1, C_2^1$  lesquels sont les centres de gravité des deux groupes A et B. Donc :

$$C_1^1 = \left( \frac{-2 + (-2) + 0 + (-2) + 3}{5}, \frac{2 + (-1) + (-1) + 0}{5} \right) = (-0.6, 0.6)$$

$$C_2^1 = (2, 2)$$

Le tableau des distances entre les individus et ces centres est :

Module: Data Mining  
 Nature de document: TD2-Solution  
 Niveau: L3-STID

Année: 2022-2023

M<sup>me</sup> N.BERMAD

	W1	W2	W3	W4	W5	W6
$C_1^1$	1.98	2.13	1.71	2.95	2.78	3.65
$C_2^1$	4	5	3.60	0	4.12	2.24

D'où les deux groupes :

$$A = \{W1, W2, W3, W5\}, \quad B = \{W4, W6\}$$

**Etape02 :**

On considère deux nouveaux centres  $C_1^2, C_2^2$  lesquels sont les centres de gravité des deux groupes A et B. Donc

$$C_1^2 = \left( \frac{-2 + (-2) + 0 + (-2)}{4}, \frac{2 + (-1) + (-1) + 3}{4} \right) = (-1.5, 0.75)$$

$$C_2^2 = \left( \frac{2 + 3}{2}, \frac{2 + 0}{2} \right) = (2.5, 1)$$

Le tableau des distances entre les individus et ces centres est

	W1	W2	W3	W4	W5	W6
$C_1^2$	1.34	1.82	2.30	3.72	2.30	4.56
$C_2^2$	4.61	4.92	3.20	1.12	4.92	1.12

D'où les deux groupes :

$$A = \{W1, W2, W3, W5\}, \quad B = \{W4, W6\}$$

On retrouve la même classification que l'étape précédente, on arrête l'algorithme.

**2. a. L'algorithme CAH**

**Etape0:**

On va regrouper les individus avec l'algorithme CAH et la méthode du voisin le plus proche munie de la distance euclidienne. Au départ chaque individu est un cluster, donc :

$$P_0 = \{\{w1\}, \{w2\}, \{w3\}, \{w4\}, \{w5\}\}$$

Module: Data Mining  
 Nature de document: TD2-Solution  
 Niveau: L3-STID

Année: 2022-2023

M<sup>me</sup> N.BERMAD

- La matrice de distance associée à  $P_0 = \{\{w1\}, \{w2\}, \{w3\}, \{w4\}, \{w5\}\}$  est :

	w1	w2	w3	w4	w5
w1	0				
w2	5.85	0			
w3	<b>1.41</b>	4.60	0		
w4	3.35	7.07	3.20	0	
w5	4.47	1.50	3.16	5.59	0

Les éléments (individus) w1 et w3 ont l'écart le plus petit : ce sont les éléments les plus proches. On les rassemble pour former le groupe :  $A = \{w1, w3\}$ . On a une nouvelle partition de P :

$$P_1 = \{\{w2\}, \{w4\}, \{w5\}, A\}$$

**Étape1:**

- La matrice de distance associée à  $P_1$  est :

	w2	w4	w5	A
w2	0			
w4	7.07	0		
w5	<b>1.50</b>	5.59	0	
A	4.60	3.20	3.16	0

On a:

$$e(w2,A) = \text{Min}(e(w2,w1), e(w2,w3)) = \text{Min}(5.85, 4.60) = 4.60$$

$$e(w4,A) = \text{Min}(e(w4,w1), e(w4,w3)) = \text{Min}(3.35, 3.20) = 3.20$$

$$e(w5,A) = \text{Min}(e(w5,w1), e(w5,w3)) = \text{Min}(4.47, 3.16) = 3.16$$

Les éléments (individus) w2 et w5 sont les plus proches. On les rassemble pour former le groupe :  $B = \{w2, w5\}$ . On a une nouvelle partition de P :

$$P_2 = \{\{w4\}, A, B\}$$

Module: Data Mining  
 Nature de document: TD2-Solution  
 Niveau: L3-STID

Année: 2022-2023

M<sup>me</sup> N.BERMAD

**Etape2:**

- La matrice de distance associée à  $P_2$  est :

	<b>w4</b>	<b>A</b>	<b>B</b>
<b>w4</b>	0		
<b>A</b>	3.20	0	
<b>B</b>	5.59	<b>3.16</b>	0

On a:

$$e(B,w4)=\text{Min}(e(w2,w4),e(w5,w4))=\text{Min}(7.07,5.59)=5.59$$

Et

$$e(B,A)=\text{Min}(e(w2,A),e(w5,A))=\text{Min}(4.60,3.16)=3.16$$

Les éléments (individus) A et B sont les plus proches. On les rassemble pour former le groupe :  $C = \{A, B\} = \{w1, w3, w2, w5\}$ . On a une nouvelle partition de P :

$$P_3 = \{w4, C\}$$

**Etape3 :**

- La matrice de distance associée à  $P_3$  est :

	<b>w4</b>	<b>C</b>
<b>w4</b>	0	
<b>C</b>	<b>3.20</b>	0

On a :

$$e(C, w4)=\text{Min}(e(A,w4), e(B,w4))=\text{Min}(3.20, 5.59)=3.20$$

Il ne reste plus que 2 éléments, w4 et C ; on les regroupe. On obtient la partition  $P_4 = \{w1, w2, w3, w4, w5\} = P$ . Cela termine l'algorithme de CAH.

**2. b. Construction du dendrogramme**

- Les individus  $\{w1\}$ ,  $\{w3\}$  ont été regroupés avec un écart de **1.41**

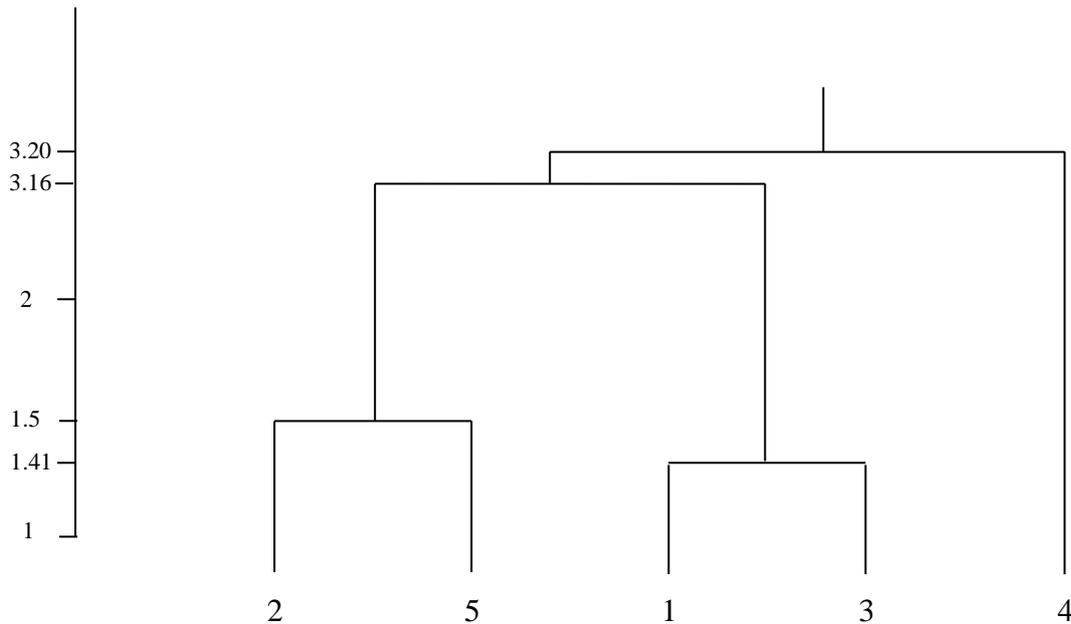
Module: Data Mining  
 Nature de document: TD2-Solution  
 Niveau: L3-STID

Année: 2022-2023

M<sup>me</sup> N.BERMAD

- Les individus  $\{w_2\}$ ,  $\{w_5\}$  ont été regroupés avec un écart de **1.50**
- Les individus  $A = \{w_1, w_3\}$  et  $B = \{w_2, w_5\}$  ont été regroupés avec un écart de **3.16**
- Les individus  $C = \{A, B\}$  et  $w_4$  ont été regroupés avec un écart de **3.20**

On peut donc construire le dendrogramme associé:



**Exercice 2 :**

1. La table de contingence (TDC) associée :

	F	H	R	U	C	M	A
W <sub>1</sub>	0	1	0	1	1	0	0
W <sub>2</sub>	1	0	0	1	1	0	0
W <sub>3</sub>	1	0	1	0	0	1	0
W <sub>4</sub>	1	0	0	1	0	0	1
W <sub>5</sub>	0	1	1	0	0	1	0
W <sub>6</sub>	0	1	1	0	0	0	1

Module: Data Mining  
 Nature de document: TD2-Solution  
 Niveau: L3-STID

Année: 2022-2023

M<sup>me</sup> N.BERMAD

$d(w_1, w_2)$ :

		$W_2$		
		1	0	sum
$W_1$	1	$a_{1,2}=2$	$b_{1,2}=1$	$a+b=3$
	0	$c_{1,2}=1$	$d_{1,2}=3$	$c+d=4$
	sum	$a+c=3$	$b+d=4$	

On a:

$$a_{1,2} = 2, b_{1,2} = 1, c_{1,2} = 1, \text{ et } d_{1,2} = 3.$$

Donc

$$d(w_1, w_2) = \frac{b_{1,2} + c_{1,2}}{a_{1,2} + b_{1,2} + c_{1,2}} = \frac{2}{4} = 0.5 // \text{ La distance de Jaccard mesure la}$$

dissimilarité entre les individus ( $\text{sim}(w_1, w_2) = 1 - d(w_1, w_2) = 1 - 0.5 = 0.5$ )

$d(w_3, w_6)$ :

On a:

$$a_{3,6} = 1, b_{3,6} = 2, c_{3,6} = 2, \text{ et } d_{3,6} = 2.$$

Donc

$$d(w_3, w_6) = \frac{b_{3,6} + c_{3,6}}{a_{3,6} + b_{3,6} + c_{3,6}} = \frac{4}{5} = 0.8$$

2. Comme

$$d(w_1, w_2) < d(w_3, w_6)$$

$w_1$  est plus proche de  $w_2$  que  $w_3$  de  $w_6$ .