

TP: Clustering

Exercice I:

$$X_1=1; X_2=2; X_3=9; X_4=12; X_5=20$$

1) a- $K=2$ $g_1=1$ et $g_2=20$ (K-means)

1^{ère} itération :

	X_1	X_2	X_3	X_4	X_5
	1	2	9	12	20
$g_1=1$ ($X_1=1$)	0	1	8	11	19
$g_2=20$ ($X_5=20$)	19	18	11	8	0

$$C_1 = \{X_1, X_2, X_3\} \rightarrow g_1 = \frac{X_1 + X_2 + X_3}{3} = \frac{1+2+9}{3} = 4$$

$$C_2 = \{X_4, X_5\} \rightarrow g_2 = 16$$

2^{ème} itération :

	X_1	X_2	X_3	X_4	X_5
	1	2	9	12	20
$g_1=4$	3	2	5	8	16
$g_2=16$	15	14	7	4	4

$$C_1 = \{X_1, X_2, X_3\}$$

$$C_2 = \{X_4, X_5\}$$

les classes n'ont pas changé, alors on arrête l'algorithme.

①

• Centre de gravité $g = \frac{1+2+9+12+20}{5} = 8,8$

• Inertie totale: $\frac{1}{5} \sum_{i=1}^5 d^2(x_i, g)$

$$= \frac{1}{5} \left[\underbrace{(1-8,8)^2}_{\substack{\downarrow \\ x_1}} + \underbrace{(2-8,8)^2}_{\substack{\downarrow \\ x_2}} + \underbrace{(9-8,8)^2}_{\substack{\downarrow \\ x_3}} + \underbrace{(12-8,8)^2}_{\substack{\downarrow \\ g}} + \underbrace{(20-8,8)^2}_{\substack{\downarrow \\ g}} \right]$$

$= 48,56$

• Inertie Inter classe: $\frac{1}{5} \sum_{k=1}^2 n_k \cdot d^2(g_k, g)$ → on a 2 classes

$$= \frac{1}{5} \left[\underbrace{(3)}_{\substack{\swarrow \\ \text{eff du} \\ \text{classe } C_1}} \times \underbrace{(4)}_{\substack{\downarrow \\ g_1}} \cdot \underbrace{(8,8)^2}_{\substack{\downarrow \\ g}} + \underbrace{(2)}_{\substack{\swarrow \\ \text{eff classe } C_2}} \times \underbrace{(16)}_{\substack{\downarrow \\ g_2}} \cdot \underbrace{(8,8)^2}_{\substack{\downarrow \\ g}} \right] = 34,56$$

• Inertie Intra-classe = Inertie totale - Inertie inter-classe
 $= 48,56 - 34,56 = 14.$

• pourcentage d'inertie expliquée par les classes:

$$R^2 = \frac{\text{Inertie Inter-classe}}{I_{\text{total}}} = \frac{34,56}{48,56} = 0.71$$

(2)

b) 1^{re} Iteration :

	X_1	X_2	X_3	X_4	X_5
	1	2	9	12	20
$g_1 = 2$	1	0	7	10	18
$g_2 = 9$	8	7	0	3	11

$$C_1 = \{X_1, X_2\} \rightarrow g_1 = \frac{1+2}{2} = 1.5$$

$$C_2 = \{X_3, X_4, X_5\} \rightarrow g_2 = \frac{9+12+20}{3} = 13.66$$

2^{ème} Iteration :

	X_1	X_2	X_3	X_4	X_5
	1	2	9	12	20
$g_1 = 1.5$	0.5	0.5	7.5	10.5	18.5
$g_2 = 13.66$	12.66	11.66	4.66	1.66	6.34

$$C_1 = \{X_1, X_2\}$$

$C_2 = \{X_3, X_4, X_5\}$; Les classes n'ont pas changé, alors on arrête l'algorithme.

• Centre de gravité : $g = 8.8$

• Inertie inter-classes :
$$I_{\text{inter}} = \frac{1}{5} [2 \times (1.5 - 8.8)^2 + 3 \times (13.66 - 8.8)^2] = 35.48$$

• Inertie totale = 48.56

(3)

$$R^2 = \frac{35,48}{48,56} = 0,73$$

c) $K=3$; $g_1=1$; $g_2=9$; $g_3=12$

1^{re} itération:

	X_1	X_2	X_3	X_4	X_5
	1	2	9	12	20
$g_1=1$	0	1	8	11	19
$g_2=9$	8	7	0	3	11
$g_3=12$	11	10	3	0	8

$$C_1 = \{X_1, X_2\} \rightarrow g_1 = \frac{1+2}{2} = 1,5$$

$$C_2 = \{X_3\} \rightarrow g_2 = 9$$

$$C_3 = \{X_4, X_5\} \rightarrow g_3 = \frac{12+20}{2} = 16$$

2^e itération:

	X_1	X_2	X_3	X_4	X_5
	1	2	9	12	20
$g_1=1,5$	0,5	0,5	7,5	10,5	18,5
$g_2=9$	8	7	0	3	11
$g_3=16$	15	14	7	4	4

④

$$C_1 = \{X_1, X_2\} \rightarrow g_1 = 1,5$$

$$C_2 = \{X_3, X_4\} \rightarrow g_2 = 10,5$$

$$C_3 = \{X_5\} \rightarrow g_3 = 20$$

3ème Itération:

	X_1	X_2	X_3	X_4	X_5
	1	2	9	12	20
$g_1 = 1,5$	0.5	0.5	7.5	10.5	18.5
$g_2 = 10,5$	8.5	8.5	1.5	1.5	9.5
$g_3 = 20$	19	18	11	8	0

$$C_1 = \{X_1, X_2\}$$

$$C_2 = \{X_3, X_4\}$$

$$C_3 = \{X_5\}$$

; les classes n'ont pas changé, donc on arrête l'algorithme.

$$\begin{aligned} I_{\text{inter}} &= 1/5 [2 \times (1,5 - 8,8)^2 + 2 \times (10,5 - 8,8)^2 + 1 \times (20 - 8,8)^2] \\ &= 47,56. \end{aligned}$$

$$R^2 = \frac{I_{\text{inter}}}{I_{\text{total}}} = \frac{47,56}{48,56} = 0,98$$

(5)

d) On peut utiliser l'entité inter-classes pour déterminer le meilleur parmi deux partitions qui ont m même nombre de classes. plus ce nombre est grande plus le pourcentage d'entité expliquée par la partition obtenue est grande. Ici, le meilleur partition entre "a" et "b" est celle de deuxième partie. (R^2 plus élevée que 1°).

⚠ Le partition du "c" a le R^2 le plus élevée, mais c'est mieux d'avoir un petit nb de classes;

2) Méthode CAH. en utilisant distance minimale:

1^{ère} itération: On considère chaque point X_i comme un classe et on calcule distance entre chaque deux points (classes).

	C_1 $\{1\}$	C_2 $\{2\}$	C_3 $\{9\}$	C_4 $\{12\}$	C_5 $\{20\}$
C_1 $\{1\}$	0	1	8	11	19
C_2 $\{2\}$	1	0	7	10	18
C_3 $\{9\}$	8	7	0	3	11
C_4 $\{12\}$	11	10	3	0	8
C_5 $\{20\}$	19	18	11	8	0

$$C_1 = \{X_1\}$$

$$C_2 = \{X_2\}$$

$$C_3 = \{X_3\}$$

$$C_4 = \{X_4\}$$

$$C_5 = \{X_5\}$$

Les classes les plus proches sont (C_1) et (C_2) , ont les regroupe als.

$$C_1 = \{X_1, X_2\}$$

$$C_2 = \{X_3\}$$

$$C_3 = \{X_4\}$$

$$C_4 = \{X_5\}$$

⑥

2^{ème} itération: On regroupe les classes les plus proches.

	C_1 $\{1,2\}$	C_2 $\{9\}$	C_3 $\{12\}$	C_4 $\{20\}$
C_1 $\{1,2\}$	0	7	10	18
C_2 $\{9\}$	7	0	3	11
C_3 $\{12\}$	10	3	0	8
C_4 $\{20\}$	18	11	8	0

(*) on prend $(9-2)=7$
et pas $(9-1)=8$
car on veut la distance minimale

Les classes les plus proches sont : C_2 et C_3 : ont les groupes

donc $C_1 = \{x_1, x_2\}$

$C_2 = \{x_4, x_3\}$

$C_3 = \{x_5\}$

On regroupe x_3 avec x_4

3^{ème} itération

	C_1 $\{1,2\}$	C_2 $\{9,12\}^*$	C_3 $\{20\}$
C_1 $\{1,2\}$	0	7	18
C_2 $\{9,12\}$	7	0	8
C_3 $\{20\}$	18	8	0

pour calculer
distances entre
 $\{9,12\}$ et $\{1,2\}$
on prend la
combinaison qui
va nous donner la
plus petite distance.
ex: $9-2=7$

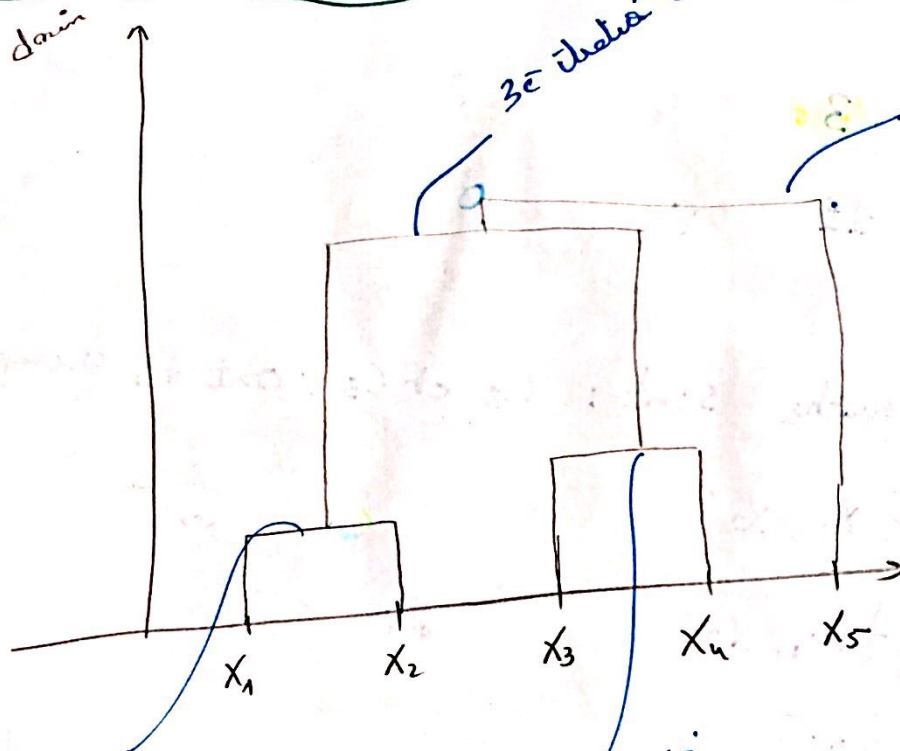
Les classes les plus proches sont : C_2 avec C_1
ont les regroupement : $C_1 = \{x_1, x_2, x_3, x_4\}$
 $C_2 = \{x_5\}$ (7)

4ème Itération:

	C1 {1, 2, 9, 12}	C2 {20}
C1 {1, 2, 9, 12}	0	8
C2 {20}	8	0

On groupe à la fin tous les classes ensemble.

Dendrogramme:



1^{ère} itération on a regroupé x_1 avec x_2

2^{ème} itération on a regroupé x_3 avec x_4

3^{ème} itération on a regroupé les x_1, x_2, x_3 et x_4

4^{ème} itération. A la fin il faut avoir tout regroupé.