

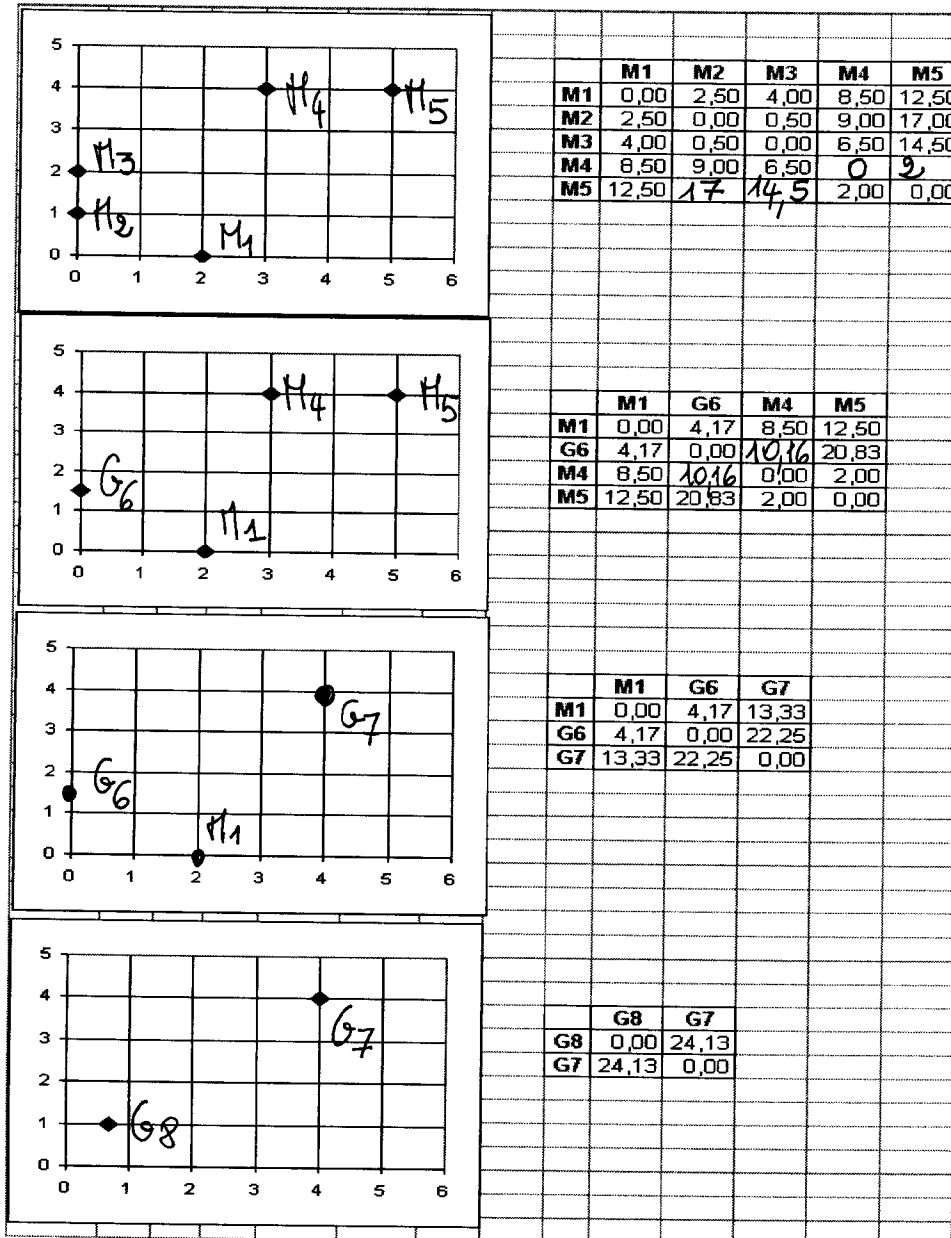
NOM :  
PRENOM :

*Corrigé*

Date :  
Groupe :

**Mathématiques pour la Biologie (semestre 2) : Feuille-réponses du TD 7**  
**Classification hiérarchique ascendante**

**Exercice 1.** : La succession des quatre dessins suivants correspond aux étapes successives d'une classification hiérarchique ascendante des cinq points  $M_1(2,0)$ ,  $M_2(0,1)$ ,  $M_3(0,2)$ ,  $M_4(3,4)$  et  $M_5(5,4)$  progressivement regroupées en classes de deux ou trois points dont les centres de gravité sont notés  $G_6$ ,  $G_7$  et  $G_8$ . On suppose que les cinq points initiaux sont tous affectés du poids 1. La distance choisie pour cette classification, qui apparaît dans les quatre matrices de distance, est l'écart de Ward.



*en complète  
par  
symétrique*

*Le point G6 a  
pour coordonnées  
(0  
3/2) et poids 2  
 $d(G6, M4) = \frac{(1)(2)}{1+2} (3^2 + (\frac{5}{2})^2)$   
 $\approx 10,1666$   
 $= d(M4, G6)$*

1. Compléter le troisième dessin en y plaçant les trois points devant y figurer et indiquer sur les quatre dessins le nom des points.
2. Compléter les six distances manquantes dans les matrices de distances.

3. Préciser les coordonnées des points  $G_6, G_7$  et  $G_8$

$G_6$  est le milieu de segment  $\Pi_2(0), \Pi_3(2)$  donc  $G_6 = \left(\frac{0+2}{2}\right) = \left(\frac{2}{2}\right)$   
 $G_7$  est le milieu de segment  $\Pi_4(1), \Pi_5(2)$  donc  $G_7 = \left(\frac{1+2}{2}\right) = \left(\frac{3}{2}\right)$   
 $G_8$  est le centre de gravité de  $G_6$  (pid 2) et  $\Pi_1$  (pid 1) - donc  $G_8 = \left(\frac{2(0)+1(2)}{3}\right)$   
 d'où  $G_8 = \left(\frac{2}{3}\right)$

4. Calculer les coordonnées du centre de gravité  $G_9$  des cinq points.

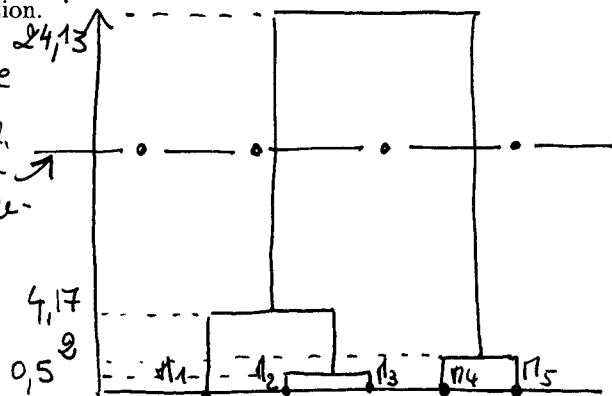
Les coordonnées de  $G_9$  sont la moyenne de celles des 5 points:

$$G_9 = \left(\frac{2+0+0+3+5}{5}\right) = \left(\frac{20}{5}\right) = (4)$$

C'est le centre de gravité de  $G_7$  (pid 2) et  $G_8$  (pid 3)

5. Tracer un dendrogramme résumant cette classification.

La meilleure partition est obtenue en découpant le dendrogramme à ce moment de ce regroupement car c'est au moment de ce regroupement qu'on a la perte d'inertie intermaximale (seul maximal).



**Exercice 2.** : (Sujet inspiré d'un article de John Hartshorne, paru dans le journal de la "British Ecological Society")

Un laboratoire d'écologie étudie les espèces micro-animales (larves, ...) présentes dans les rivières et les étangs. Il réalise, dans 6 sites de rivière, notés  $R_1, R_2, R_3, R_4, R_5$  et  $R_6$ , et 3 sites d'étangs, notés  $E_1, E_2$  et  $E_3$ , des prélèvements répétés qui lui permettent d'avancer une liste des espèces présentes dans chacun de ces sites et de repérer les espèces présentes dans plusieurs sites à la fois. La matrice suivante contient, pour chaque paire de sites  $A$  et  $B$ , le nombre d'espèces communes aux 2 sites. Ainsi on y lit par exemple que 11 espèces sont présentes au site  $R_1$  et qu'il y a 7 espèces présentes à la fois au site  $R_1$  et au site  $R_2$ .

	$R_1$	$R_2$	$R_3$	$R_4$	$R_5$	$R_6$	$E_1$	$E_2$	$E_3$
$R_1$	11	7	4	6	6	7	4	4	3
$R_2$	7	15	8	8	9	6	3	3	2
$R_3$	4	8	13	7	7	4	2	3	2
$R_4$	6	8	7	15	7	6	6	8	6
$R_5$	6	9	7	7	12	4	3	5	4
$R_6$	7	6	4	6	4	10	6	5	5
$E_1$	4	3	2	6	3	6	13	10	9
$E_2$	4	3	3	8	5	5	10	15	11
$E_3$	3	2	2	6	4	5	9	11	12

On se propose de regrouper les 9 sites en trois ou quatre classes composées de sites où ce sont pratiquement les mêmes espèces qui sont présentes. Pour réaliser cette classification, on propose de mesurer la distance entre deux sites  $A$  et  $B$  par la formule

$$d(A, B) = \frac{n_A + n_B - 2n_{AB}}{n_A + n_B}$$

où  $n_A$  (resp.  $n_B$ ) désigne le nombre d'espèces présentes au site  $A$  (resp. au site  $B$ ) et  $n_{AB}$  le nombre d'espèces en commun entre les sites  $A$  et  $B$ . On obtient la matrice des distances suivante :

$$d(R_2, R_4) = \frac{13+15-2(7)}{13+15} = \frac{14}{28} = 0,5$$

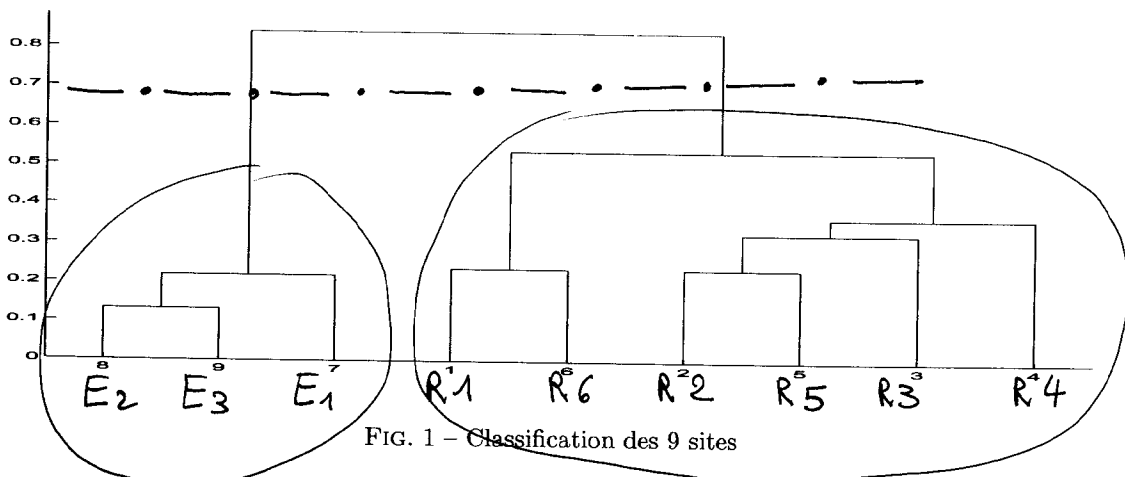
$$d(E_1, E_2) = \frac{13+15-2(10)}{13+15} = \frac{8}{28} = 0,285$$

	R1	R2	R3	R4	R5	R6	E1	E2	E3
R1	0	0,462	0,666	0,538	0,478	0,334	0,666	0,692	0,74
R2	0,462	0	0,428	0,466	0,334	0,52	0,786	0,8	0,852
R3	0,666	0,428	0	0,5	0,44	0,652	0,846	0,786	0,84
R4	0,538	0,466	0,5	0	0,482	0,52	0,572	0,466	0,556
R5	0,478	0,334	0,44	0,482	0	0,636	0,76	0,63	0,666
R6	0,334	0,52	0,652	0,52	0,636	0	0,478	0,6	0,546
E1	0,666	0,786	0,846	0,572	0,76	0,478	0	0,285	0,28
E2	0,692	0,8	0,786	0,466	0,63	0,6	0,285	0	0,186
E3	0,74	0,852	0,84	0,556	0,666	0,546	0,28	0,186	0

1. Compléter les coefficients manquants de cette matrice.
2. Préciser quels sont les deux sites les plus proches ainsi que les deux sites les plus éloignés.

On constate que les 2 sites les plus proches sont E2 et E3 car  $d(E_2, E_3) = 0,186$  et les plus éloignés sont E3 et R2 car  $d(E_3, R_2) = 0,852$ .

3. La classification conduit au dendrogramme représenté ci-dessous.



Décrire la composition des classes de la partition qui vous semble la plus appropriée.

Le seul maximal du dendrogramme a lieu au dernier regroupement (mais on pourrait avoir considéré l'avant dernier). Dans ce cas on trouve les 2 classes  $(R_1, \dots, R_6)$  de rivières et  $(E_1, \dots, E_3)$  des étangs.

4. Un autre choix de distance entre les sites aurait-il pu conduire à une partition différente?

Ouï. La décomposition dépend souvent de la distance choisie. Il faut donc s'attendre à trouver peut être une autre partition.

Pourquoi n'a-t-on pas choisi la distance euclidienne?

Impossible ici car les différents sites ne sont pas repérés par des coordonnées : ce ne sont pas les points d'un espace euclidien.