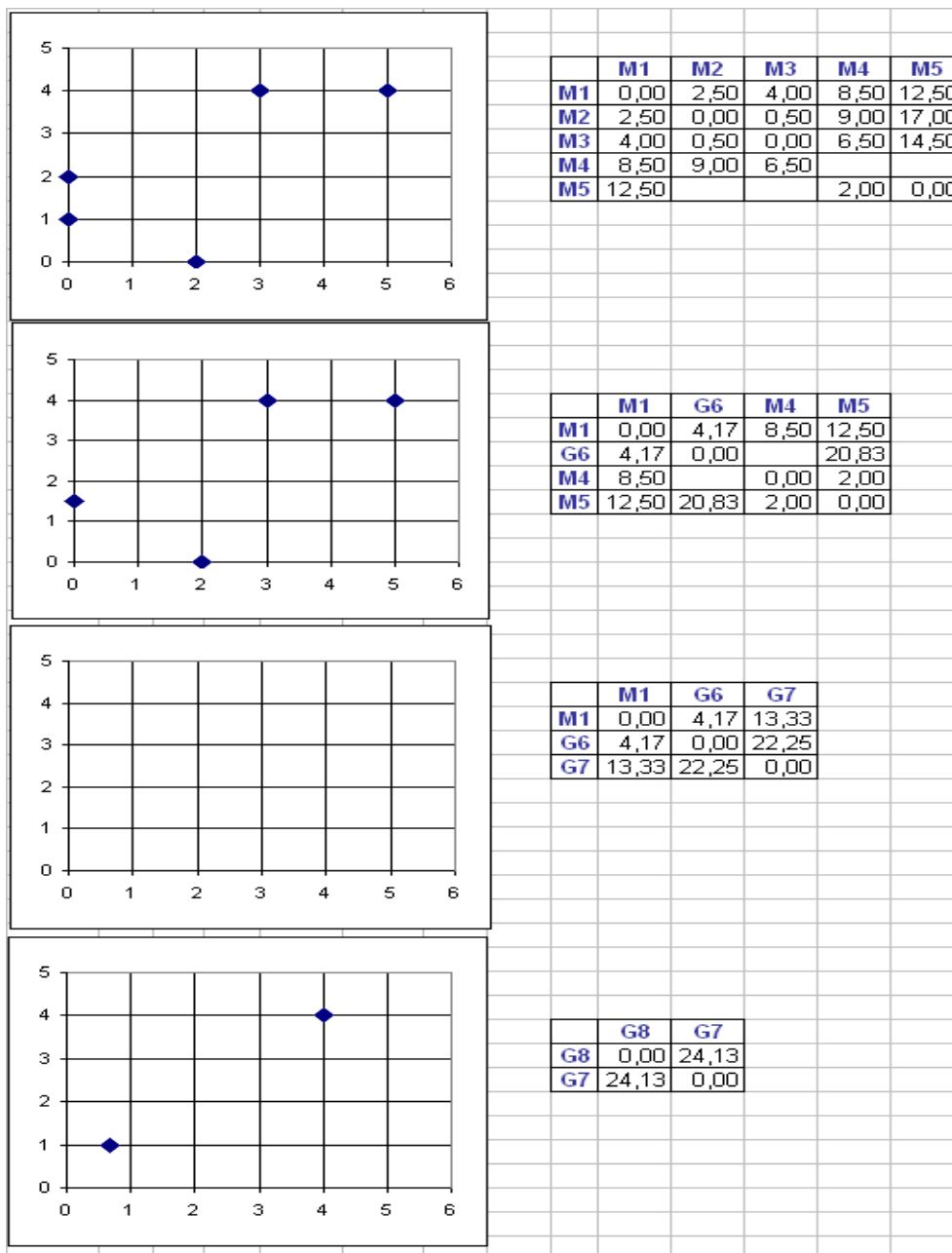


NOM :
 PRENOM :

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.



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

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

5. Tracer un dendrogramme résumant cette classification.

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 $R1$, $R2$, $R3$, $R4$, $R5$ et $R6$, et 3 sites d'étangs, notés $E1$, $E2$ et $E3$, 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 $R1$ et qu'il y a 7 espèces présentes à la fois au site $R1$ et au site $R2$.

	$R1$	$R2$	$R3$	$R4$	$R5$	$R6$	$E1$	$E2$	$E3$
$R1$	11	7	4	6	6	7	4	4	3
$R2$	7	15	8	8	9	6	3	3	2
$R3$	4	8	13	7	7	4	2	3	2
$R4$	6	8	7	15	7	6	6	8	6
$R5$	6	9	7	7	12	4	3	5	4
$R6$	7	6	4	6	4	10	6	5	5
$E1$	4	3	2	6	3	6	13	10	9
$E2$	4	3	3	8	5	5	10	15	11
$E3$	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 :

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

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

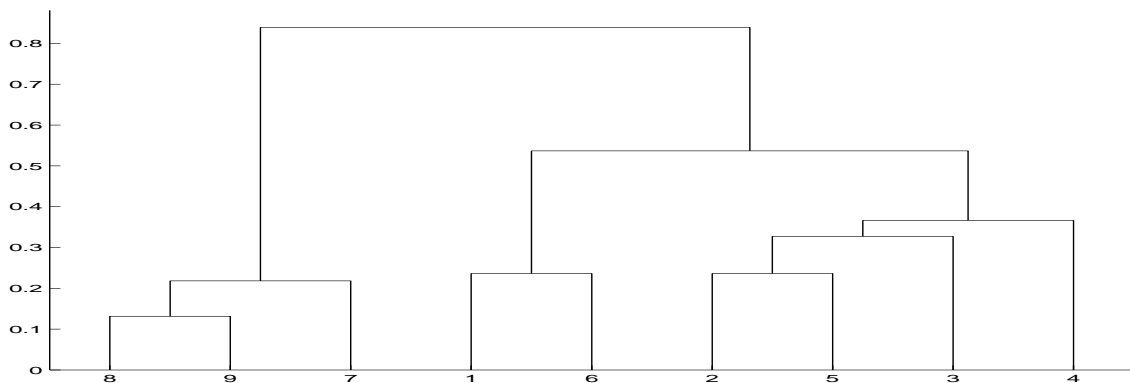


FIG. 1 – Classification des 9 sites

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

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

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