

TD 1-2

Exercice 1

Considérer les données suivantes de l'attribut *âge* écrites par ordre croissant :

13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70.

- 1) Quelle est la *moyenne* de ces données ? Quelle est la *médiane* ?

Moyenne = 29,96

Médiane = 25

- 2) Quel est le *mode* des données ? Commenter la modalité des données (i.e., bimodales, trimodales, etc.).

Mode1 = 25

Mode2 = 35

Les données sont bimodales

- 3) Pouvez-vous déterminer (approximativement) le premier quartile ( $Q_1$ ) et le troisième quartile ( $Q_3$ ) des données ?

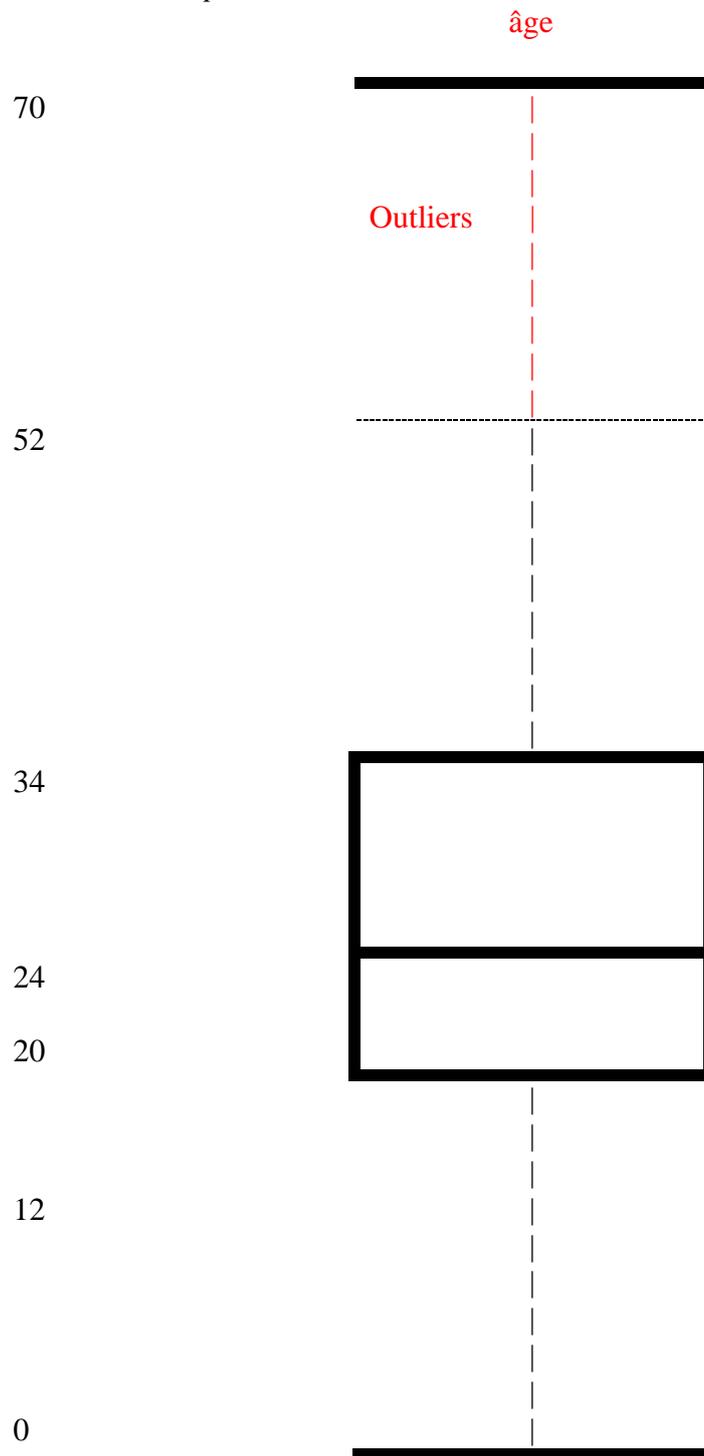
$Q_1 = 20$

$Q_3 = 35$

- 4) Donner le *résumé des cinq nombres* des données.

(13, 20, 25, 35, 70)

5) Dessiner un *boxplot* des données.



$$\text{IQR} = 35 - 20 = 15$$

6) quelle est la différence entre un *quantile-quantile plot* et un *quantile plot* ?

un quantile-quantile plot fait intervenir 2 sources de données alors que le quantile plot se fait pour une seule source de données

### **Exercice 2**

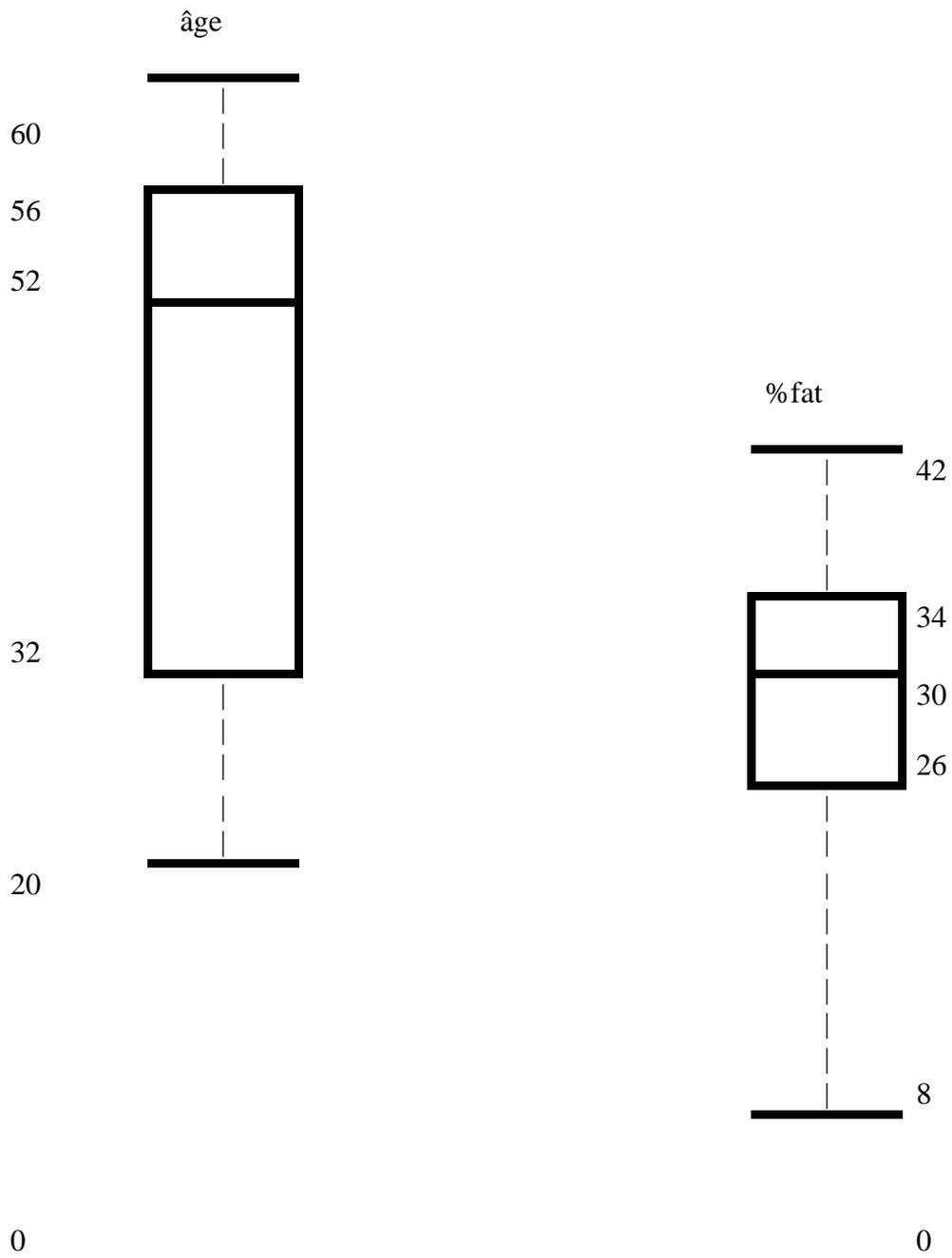
Supposer que dans un hôpital, on ait effectué des tests sur le taux des graisses de 18 adultes sélectionnés au hasard et dont les résultats accompagnés de l'âge sont comme suit :

<i>âge</i>	23	23	27	27	39	41	47	49	50
<i>%fat</i>	9.5	26.5	7.8	17.8	31.4	25.9	27.4	27.2	31.2
<i>âge</i>	52	54	54	56	57	58	58	60	61
<i>%fat</i>	34.6	42.5	28.8	33.4	30.2	34.1	32.9	41.2	35.7

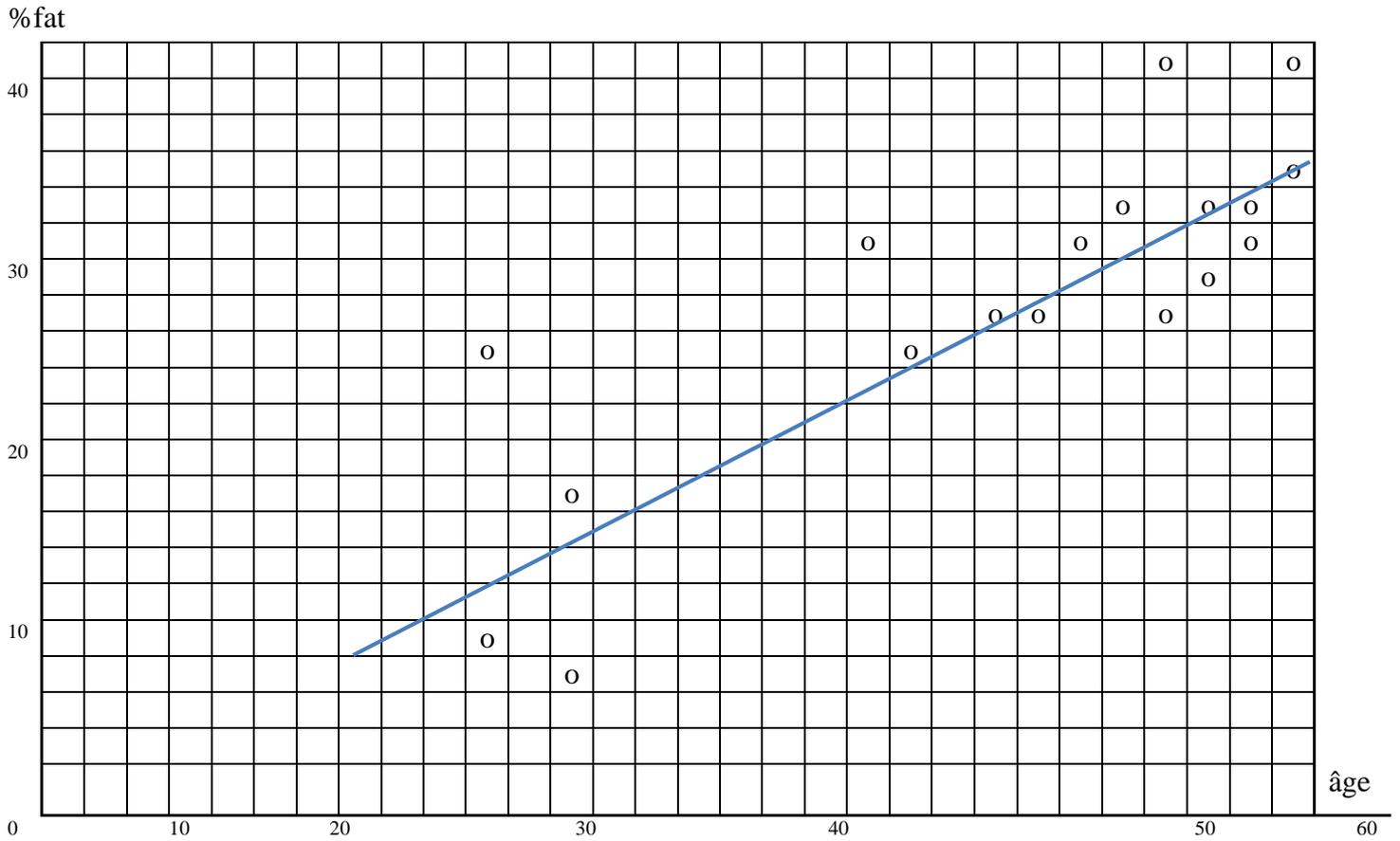
(a) Calculer la moyenne, la médiane et l'écart-type de l'*âge* et du taux de graisse *%fat*.

	moyenne	médiane	$\sigma$	Q1	Q3	IQR	1,5xIQR
<i>âge</i>	46.44	51	12.86	33	57.5	24.5	36.75
<i>%fat</i>	28.78	30.7	9.00	26.2	34.35	8.15	12.22

(b) Dessiner les boxplots pour *âge* and *%fat*.



(c) Dessiner un *scatter plot*.



(d) Normaliser les deux variables avec la technique du *z-score*.

$$v'_{\hat{age}} = \frac{v - \mu_{\hat{age}}}{\sigma_{\hat{age}}}$$

$$v'_{\%fat} = \frac{v - \mu_{\%fat}}{\sigma_{\%fat}}$$

âge	âge-normalisé
23	-1.82
23	-1.82
27	-1.51
27	-1.51
39	-0.57
41	-0.42
47	0.04

âge	âge-normalisé
52	0.43
54	0.58
54	0.58
56	0.74
57	0.82
58	0.89
58	0.89

49	0.19
50	0.27
<b>%fat</b>	<b>%fat-normalisé</b>
9.5	-2.14
26.5	-0.25
7.8	-2.33
17.8	-1.22
31.4	0.29
25.9	-0.32
27.4	-0.15
27.2	-0.17
31.2	0.26

60	1.05
61	1.13
<b>%fat</b>	<b>%fat-normalisé</b>
34.6	0.64
42.5	1.52
28.8	0.002
33.4	0.51
30.2	0.15
34.1	0.59
32.9	0.45
41.2	1.38
35.7	0.76

- (e) Calculer le *coefficient de corrélation* (Pearson's product moment coefficient).  
Ces deux variables sont-elles corrélées positivement, négativement ou pas du tout?

$$\tau_{A,B} = \frac{\sum AB - n\bar{A}\bar{B}}{(n-1)\sigma_A\sigma_B}$$

A= âge

B= %fat

$$\sum AB = 7986.3 + 17776.9 = 25763.2$$

$$n\bar{A}\bar{B} = 24057.7776$$

$$(n-1)\sigma_A\sigma_B = 1967.58$$

$$\tau_{A,B} = 0.86 > 0$$

Les deux variables sont positivement corrélées.

### Exercice 3

Considérer les notes suivantes d'un groupe de 10 étudiants pour les cours respectifs de 'datamining' et de 'méta-heuristiques' :

	Datamining				Méta-heuristiques			
	TP1/10	Test1/10	EMD1/20	Final1/20	TP2/10	Test2/10	EMD2/20	Final2/20
<b>1</b>	7,75	5	13,5	13,13	7,25	6,50	16	14,88
<b>2</b>	6	6	11	11,50	7,75	5,50	11,5	12,38
<b>3</b>	6	4	7,5	8,75	7,87	6,25	10	12,06
<b>4</b>	6,17	4	5,5	7,84	7,75	4,00	11,5	11,44
<b>5</b>	8	5	11	12,00	6,62	6,50	13	13,63
<b>6</b>	7,75	2	10,5	10,13	7,50	7,50	12,5	13,31
<b>7</b>	6	4	5,5	7,75	7,37	3,50	12,5	11,75
<b>8</b>	6,75	5	9	10,38	7,62	7,50	5,5	10,19
<b>9</b>	5,67	6	6	8,84	8,37	4,00	8,5	10,06
<b>10</b>	6,5	6	9	10,75	7,37	5,50	17	15,44

La note 'Final' est obtenue comme suit :

$$\text{Final} = (\text{TP} + \text{Test} + \text{EMD})/2$$

- 1) A-t-on besoin de normaliser les notes de TP, Test et EMD pour calculer la note 'Final' ? Pourquoi ? Si oui, procéder à la normalisation des données.

Nous n'avons pas besoin de normaliser les notes car la note finale appartient à l'intervall [0-20] tout comme EMD et TP+Test.

- 2) Que représente concrètement 'Final' par rapport aux notes de TP, Test et EMD ?  
Final représente la moyenne pondérée entre EMD de coefficient 2 et TP+Test de coefficient 1.

- 3) Pour chacune des notes Final1 et Final2, Calculer :
- La moyenne
  - La médiane
  - Le mode
- Que peut-on en conclure?

	Final1	Final2
moyenne	10.10	12.51
médiane	10.25	12.22
mode	10	Pas de mode

La distribution de Final1 est symétrique car moyenne=médiane=mode. En conséquence, Final1 suit la loi normale.

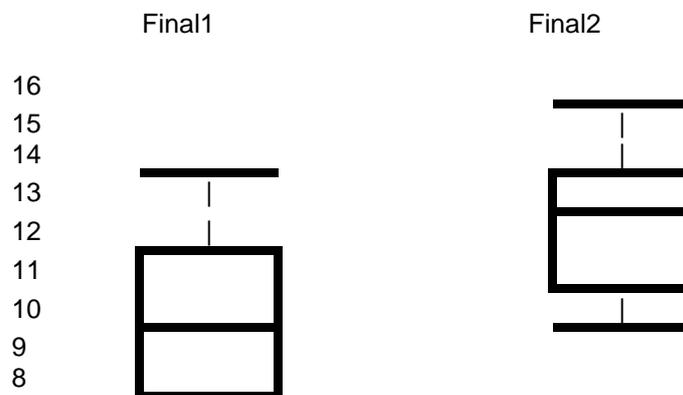
- d. Calculer approximativement le premier et le troisième quartiles

	Final1	Final2
Q1	8.29	11.59
Q3	11.75	14.25

- e. Donner le résumé des cinq nombres pour ces données.

	Final1	Final2
minimum	7.75	10.06
Q1	8.29	11.59
médiane	10.25	12.22
Q3	11.75	14.25
maximum	13.13	15.44

- f. Dessiner une *boxplot* pour ces données.



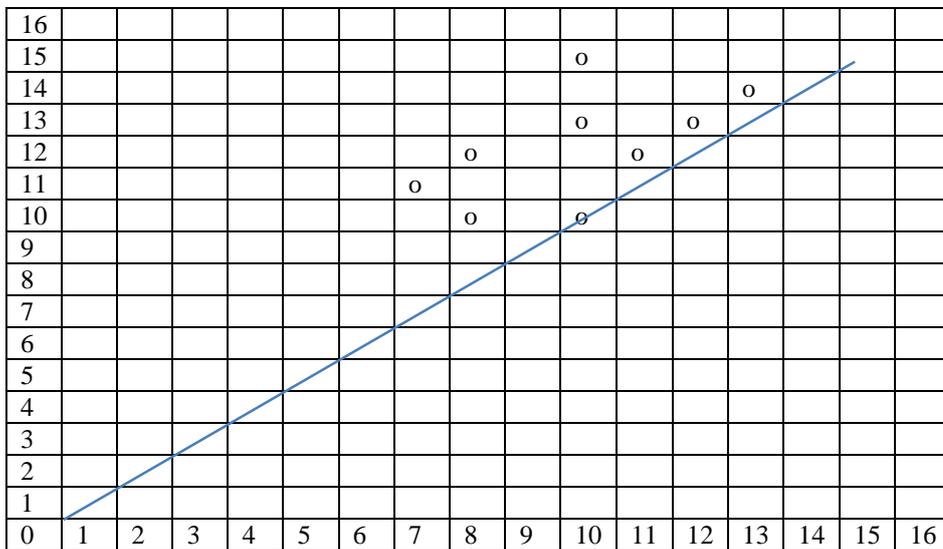
7



0

4) Dessiner un *q-q plot* pour les deux variables Final1 et Final2. Que peut-on en conclure ?

Final2



Final1

Le q-q graphe montre bien que les notes de métaheuristiques sont toutes supérieures à celles des notes de datamining.

5) Calculer le *coefficient de corrélation* (Pearson's product moment coefficient) pour les deux variables. Que peut-on en conclure ?

$$r_{A,B} = \frac{\sum AB - n\bar{A}\bar{B}}{(n-1)\sigma_A\sigma_B} \quad A = \text{Final1} \quad B = \text{Final2}$$

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^n (x_i - \mu)^2 = \frac{1}{N} \sum_{i=1}^n x_i^2 - \mu^2$$

	Final1	Final2
moyenne	10.10	12.51
Ecart-type ( $\sigma$ )	1.75	1.74

$$\tau_{A,B} = \frac{\sum AB - n\bar{A}\bar{B}}{(n-1)\sigma_A\sigma_B} = 0,71$$

Les deux notes sont positivement corrélées.

#### **Exercice 4**

L'analyse de données et le data mining sont des technologies transversales très utiles ces dernières années à cause de la présence des données dans quasiment toute entreprise.

- 1) Expliquer clairement la différence entre l'analyse de données et le datamining.
- 2) Quelles sont les fonctionnalités du datamining ? Citer au moins un algorithme de datamining pour chacune des fonctionnalités.
- 3) Enumérer les différentes étapes d'un système de data mining orienté intelligence économique ou business intelligence.
- 4) Considérer le domaine de la finance et plus particulièrement l'analyse et la gestion financière.
  - a. Imaginer au moins trois sources de données pour une telle application. A-t-on besoin d'analyse de données ou de datamining pour traiter les données ? pourquoi ?
  - b. Imaginer trois types de fonctionnalités pour le traitement des données de l'application et donner une technique pour chacune d'elles.
  - c. Considérer la requête suivante de datamining exprimée dans le langage DMQL (Data Mining Query Language) :

```

use database AllElectronics-db
use hierarchy location-hierarchy for T.branch, age-hierarchy for C.age
mine classification as promising-customers
in relevance to C.age, C.income, I.type, I.place-made, T.branch
from customer C, item I, transaction T
where I.item-ID = T.item-ID and C.cust-ID = T.cust-ID
and C.income >= 40,000 and I.price >= 100
group by T.cust-ID
having sum(I.price) >= 1,000
display as rules
  
```

Quels sont les résultats que la requête doit rendre ? expliquer.

#### **Exercice 5**

Supposer qu'un groupe de 18 étudiants ont obtenu les notes suivantes dans deux examens de cours différents :

Note1	2	2	4	4	4	7	8	8	10
Note2	4	5	6	7	7	9	10	19	11
Note1	10	12	12	12	13	14	14	15	17
Note2	8	11	15	15	16	16	16	19	19

- 1) Calculer la moyenne, la médiane et l'écart-type pour les deux notes.

- 2) Dessiner les boîtes à moustache pour *les deux notes*. *Que pouvez-vous conclure ?*
- 3) Dessiner un *histogramme* pour représenter ces données.
- 4) Calculer le *coefficient de corrélation* de Pearson. Ces deux variables sont-elles corrélées ?
- 5) Comment peut-on réduire ces données ? Donner le résultat de la réduction.

### **Exercice 6**

Considérer les attributs *longueur du sépale en cm* et *longueur de la pétale en cm* du dataset Iris. La table ci-dessous exhibe 12 instances du dataset pour les deux attributs apparaissant respectivement en deuxième et troisième colonnes.

Instance	Sépale	Pétale
1	4.9	1.4
2	5.0	1.4
3	5.4	1.7
4	4.6	1.4
5	5.5	4.0
6	5.1	3.0
7	5.7	4.5
8	5.0	3.3
9	4.9	4.5
10	5.7	5.0
11	5.8	5.1
12	5.6	4.9

- 1) Dessiner les boîtes à moustaches pour chacun des attributs *Sépale* et *Pétale*. *Que pouvez-vous conclure ?*
- 2) Appliquer l'algorithme *Chimerge* ci-dessous pour discrétiser l'attribut *Sépale*.

### **Algorithme ChiMerge**

1. trier les valeurs de l'attribut par ordre croissant.
2. considérer chaque valeur dans un intervalle distinct.
3. calculer la valeur de  $\chi^2$  pour tous les intervalles adjacents.
4. fusionner les paires d'intervalles qui ont la plus petite valeur de  $\chi^2$ .
5. arrêter le processus quand le nombre d'intervalles est égal à 4 sinon aller à (3).

La formule du  $\chi^2$  est donnée comme suit:

$$\chi^2 = \sum_{i=1}^{i=m} \frac{(R_i - E)^2}{E}$$

où:

$m$  est le nombre d'intervalles à comparer (2 dans ce cas),

$R_i$  est le nombre de valeurs de l'intervalle  $i$ ,

$E$  est la fréquence moyenne calculée comme:  $E = n / \text{MaxIntervalles}$ ,

$n$  est le nombre total de valeurs,

$\text{MaxIntervalles}$  est le nombre maximum d'intervalles.

Considérer les attributs *longueur du sépale en cm* et *longueur de la pétale en cm* du dataset Iris. La table ci-dessous exhibe 12 instances du dataset pour les deux attributs apparaissant respectivement en deuxième et troisième colonnes.

Instance	Sépale	Pétale
1	4.9	1.4
2	5.0	1.4
3	5.4	1.7
4	4.6	1.4
5	5.5	4.0
6	5.0	3.0
7	5.7	4.5
8	5.0	3.3
9	4.9	4.5
10	5.7	5.0
11	5.8	5.1
12	5.6	4.9

1) Dessiner les boîtes à moustaches pour chacun des attributs Sépale et Pétale. Que pouvez-vous conclure ?

### Exercice 6

A)

	Moyenne	Médiane	Mode
Sépale <b>0.75pts</b>	5.0	5.0	5.1 et 5.4
Pétale <b>0.75pts</b>	1.4	1.4	1.4

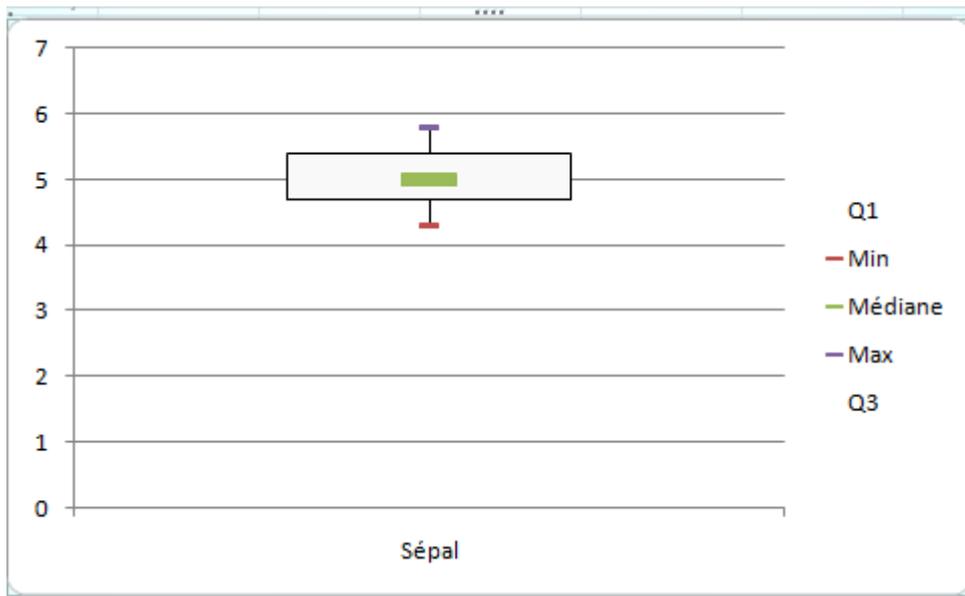
- Les données de sépale sont asymétriques. **0.25pts**
- Les données de pétale sont symétriques. **0.25pts**

B)

	Min	Q1	Médiane	Q3	Max
Sépale <b>0.5pts</b>	4.3	4.7	5.0	5.4	5.8
Pétale <b>0.5pts</b>	1.1	1.4	1.4	1.5	1.7

Sépale :

Box plot : **0.25pts**



Distribution homogènes des données .

La médiane se confond avec la moyenne.

Outliers  $> 1.5 \cdot \text{IQR}$

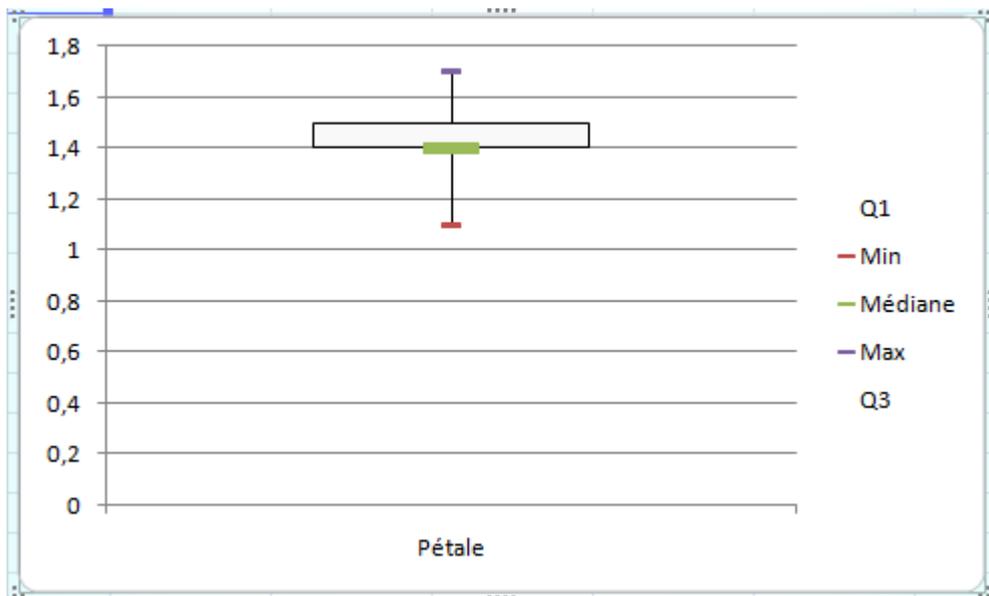
$\text{IQR} = 0.7$       **0.25pts**

Outliers  $> 1.5 \cdot 0.7 = 1.05$

→ 0 outliers

Pétale :

Box plot : **0.25pts**



La plupart (60%) des données sont comprises entre 1.4 et 1.5

La médiane se confond avec Q1.

Plus de 35% sont égale à 1.4.

Ecart entre les données est faible :  $\text{min-max} = 1.7 - 1.1 = 0.6$  **0.25pts**

$\text{IQR} = 1.5 - 1.4 = 0.1$

Outliers  $> 1.5 * 0.1 = 0.15$

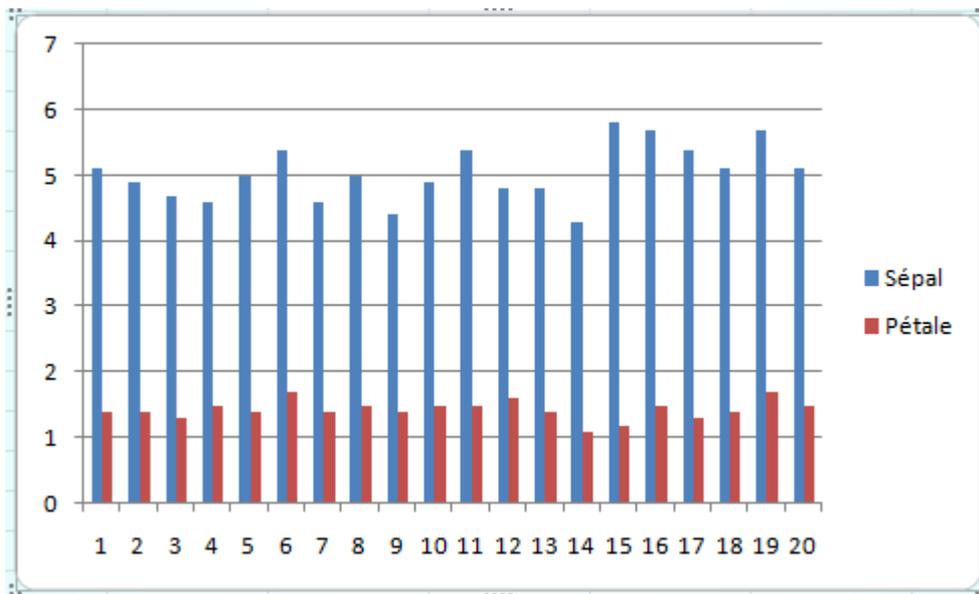
$> 1.65$  (6, 19)

$< 1.25$  (15, 14)

→ 4 outliers

C)

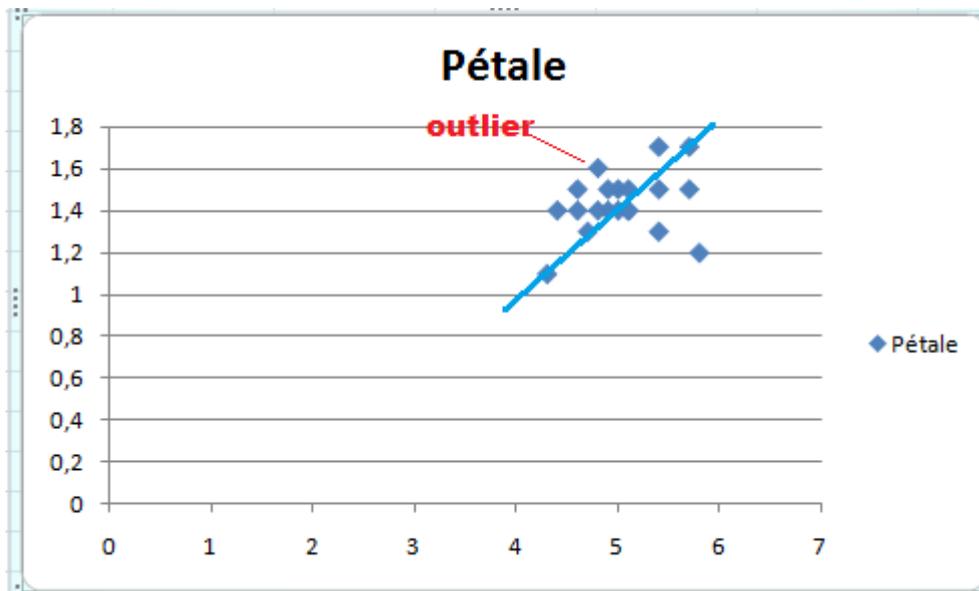
Histogrammes. **0.5pts**



On retrouve la symétrie pour pétale. **0.25pts**

Et l'asymétrie pour sépal. **0.25pts**

D)q-q plot **0.5pts**



Scatter plot. **0.5pts**

pétale : presque constante .

sépal : variable.

Le gap entre les deux est presque constant.

**0.5pts**

Corrélé positivement.

**E) coefficient de Pearson**

$$r_{AB} = \frac{\sum (A - \bar{A})(B - \bar{B})}{(n-1)\sigma_A\sigma_B}$$

$$B^2_{\text{sépal}} = 1/(n-1) * \sum_1^n (x_i - \bar{x})^2 = 1/19$$

$$(0.49+0.36+0.32+0.09+0.02+0.48+0.98+0.65+0.08)=3.49/19=0.18$$

**B<sub>sépal</sub>=0.42 0.25pts**

$$B^2_{\text{pétale}} = 1/19 (0.09+0.04+0.02+0.06+0.04+0.18)=0.43/19=0.02$$

**B<sub>pétale</sub>=0.15 0.25pts**

$$X = (0.3 * (-0.1)) + (-0.4 * 0.1) + (0.4 * 0.3) + (-0.1 * 0.1) +$$

$$(0.4 * 0.1) + (-0.2 * 0.2) + (-0.7 * -0.3) + (0.8 * -0.2) + (0.7 * 0.1) +$$

$$(0.4 * -0.1) + (0.7 * 0.3) + (0.1 * 0.1)$$

$$r_{\text{ sépal pétale}} = \frac{X}{19*0.42*0.15} = 0.40/1.097 = 0.36 \quad \mathbf{0.25pts}$$

Les deux attributs sont corrélés positivement. **0.25pts**

$\chi^2$  :

<b>1 pts</b>	>Moy pétale	<= moy pétale	Somme (ligne)
>Moy sépal	5 (4)	4 (5)	9
<= moy sépal	4 (5)	7 (6)	11
Somme colonne	9	11	20

$$\chi^2 = (5-4)^2 / 4 + (4-5)^2 / 5 + (4-5)^2 / 5 + (7-6)^2 / 4 = 1/4 + 1/5 +$$

$$1/5 + 1/4 = 2/4 + 2/5 = 0.5 + 0.4 = 0.9 \quad \mathbf{0.25pts}$$

Les deux groupes sont corrélés positivement. **0.25pts**

**G)** Si l'on tient compte du  $\chi^2$ , la réduction se fera en considérant les différentes classes (groupes).

Chaque groupe sera présenté par la moyenne de ses éléments ainsi, on obtient :

(Sépal > moy , pétale>moy) = 5.5 (sépal) ,1.6 (pétale)

(Sépal <= moy , pétale>moy) =4.8 ,1.5

(Sépal > moy , pétale<=moy) = 5.5 ,1.3

(Sépal ><=moy , pétale<=moy) =.4.7 ,1.3

Exemple :

**1pts**

$$(5.4+5.4+5.7+5.7+5.1)/5=5.5$$

2) Appliquer l'algorithme *Chimerge* ci-dessous pour discrétiser l'attribut Sépale.

#### **Algorithme ChiMerge**

1. trier les valeurs de l'attribut par ordre croissant.
2. considérer chaque valeur dans un intervalle distinct.
3. calculer la valeur de  $\chi^2$  pour tous les intervalles adjacents.
4. fusionner les paires d'intervalles qui ont la plus petite valeur de  $\chi^2$ .
5. arrêter le processus quand le nombre d'intervalles est égal à 4 sinon aller à (3).

La formule du  $\chi^2$  est donnée comme suit:

$$\chi^2 = \sum_{i=1}^{i=m} \frac{(R_i - E)^2}{E} \quad (5)$$

où:

$m$  est le nombre d'intervalles à comparer (2 dans ce cas),

$R_i$  est le nombre de valeurs de l'intervalle  $i$ ,

$E$  est la fréquence moyenne calculée comme:  $E = n/MaxIntervalles$ ,

$n$  est le nombre total de valeurs,

$MaxIntervalles$  est le nombre maximum d'intervalles.

**2) Trier les valeurs de l'attribut Sépale :**

4.6
4.9
4.9
5.0
5.0
5.0
5.4
5.5
5.6
5.7
5.7
5.8

Intervalles initiaux :

[4.6 (1)]
[4.9 (2)]
[5.0 (3)]
[5.4 (1)]
[5.5 (1)]
[5.6 (1)]
[5.7(2)]
[5.8 (1)]

$$E = 12/4 = 3$$

1<sup>ère</sup> itération

Intervalle	$\chi^2$
[4.6 (1), 4.9 (2)]	$4/3+1/3=5/3$
[4.9 (2), 5.0 (3)]	$1/3+0/3 = 1/3$
[5.0 (3), 5.4 (1)]	$0/3+4/3 = 4/3$
[5.4 (1), 5.5 (1)]	$4/3+4/3=8/3$
[5.5 (1), 5.6 (1)]	$8/3$
[5.6 (1), 5.7(2)]	$5/3$
[5.7(2), 5.8 (1)]	$5/3$

Résultat :

Intervalle
[4.6 (1)]
[4.9 (2), 5.0 (3)]
[5.4 (1)]
[5.5 (1)]
[5.6 (1)]
[5.7(2)]
[5.8 (1)]

2<sup>ème</sup> itération

Intervalle	$\chi^2$
[[4.6 (1)], [4.9 (2), 5.0 (3)]]	$4/3+4/3=8/3$
[4.9 (2), 5.0 (3)], [5.4 (1)]	$4/3+4/3 = 8/3$
[5.4 (1)], [5.5 (1)]	$4/3+4/3=8/3$
[5.5 (1)], [5.6 (1)]	$8/3$
[5.6 (1)], [5.7(2)]	$5/3$
[5.7(2), 5.8 (1)]	$5/3$

Résultat :

Intervalle
------------

[4.6 (1)]
[4.9 (2), 5.0 (3)]
[5.4 (1)]
[5.5 (1)]
[5.6 (1), 5.7(2)]
[5.8 (1)]

3<sup>ème</sup> itération

Intervalle	$\chi^2$
[[4.6 (1)], [4.9 (2), 5.0 (3)]]	$4/3+4/3=8/3$
[4.9 (2), 5.0 (3)], [5.4 (1)]	$4/3+4/3 =8/3$
[5.4 (1)], [5.5 (1)]	$4/3+4/3=8/3$
[5.5 (1)], [5.6 (1), 5.7(2)]	$4/3$
[5.6 (1), 5.7(2)], [5.8 (1)]	$4/3$

Résultat :

Intervalle
[4.6 (1)]
[4.9 (2), 5.0 (3)]
[5.4 (1)]
[5.5 (1), 5.6 (1), 5.7(2)]
[5.8 (1)]

4<sup>ème</sup> itération

Intervalle	$\chi^2$
[4.6 (1)], [4.9 (2), 5.0 (3)]]	$4/3+4/3=8/3$
[4.9 (2), 5.0 (3)], [5.4 (1)]	$4/3+4/3 =8/3$
[5.4 (1)], [5.5 (1), 5.6 (1), 5.7(2)]	$4/3+1/3=5/3$
[5.5 (1), 5.6 (1), 5.7(2)], [5.8 (1)]	$5/3$

Résultat :

Intervalle
[4.6 (1)]
[4.9 (2), 5.0 (3)]
[5.4 (1)], [5.5 (1), 5.6 (1), 5.7(2)]
[5.8 (1)]

### **Exercice 7**

Considérer les attributs *longueur du sépal en cm* et *longueur de la pétale en cm* du dataset Iris. La table ci-dessous exhibe les 20 premières entrées du dataset pour les deux attributs apparaissant respectivement en deuxième et troisième colonnes.

Instance	Sépale	Pétale
1	5.1	1.4
2	4.9	1.4
3	4.7	1.3
4	4.6	1.5
5	5.0	1.4
6	5.4	1.7
7	4.6	1.4
8	5.0	1.5
9	4.4	1.4
10	4.9	1.5
11	5.4	1.5
12	4.8	1.6
13	4.8	1.4
14	4.3	1.1
15	5.8	1.2
16	5.7	1.5
17	5.4	1.3
18	5.1	1.4
19	5.7	1.7
20	5.1	1.5

- 1) Calculer la moyenne, la médiane et le mode pour chacun des deux attributs. Que pouvez-vous conclure ?
  - 2) Dessiner les boîtes à moustaches pour les deux attributs. Que pouvez-vous conclure ?
  - 3) Dessiner un histogramme pour représenter ces données. Qu'observez-vous ?
  - 4) Donner le diagramme de dispersion des deux attributs (q-q plot et scatter plot). Que pouvez-vous conclure ?
  - 5) Calculer le coefficient de corrélation pour les deux variables. Que pouvez-vous conclure ?
  - 6) Les 20 instances se divisent en 2 catégories :
    - ✓ Celles qui dépassent la moyenne pour la longueur du sépal.
    - ✓ Celles qui dépassent la moyenne pour la longueur de la pétale.
- Calculer le  $\chi^2$  pour ces deux sous-groupes d'instances. Que pouvez-vous conclure ?
- 7) Comment peut-on réduire ces données ? Donner le résultat de la réduction.

### Exercice 8

Considérer les 10 premières instances du dataset Heart que vous avez étudié pour votre mini-projet.

- 1) Calculer le coefficient de corrélation entre l'attribut *sex* (2<sup>e</sup> colonne) et l'attribut *fasting blood sugar or FBS* (colonne 6). Que peut-on en déduire ?

Age	sex	chest pain	blood pressure	cholestorol	FBS > 120 mg/dl	ECG (0,1,2)	heart rate	Exercise angina	oldpeak	ST	vessels (0-3)	thal	
70	1	4	130	322	0	2	109	0	2.4	2	3	3	present
67	0	3	115	564	0	2	160	0	1.6	2	0	7	absent
57	1	2	124	261	0	0	141	0	0.3	1	0	7	present
64	1	4	128	263	0	0	105	1	0.2	2	1	7	absent
74	0	2	120	269	0	2	121	1	0.2	1	1	3	absent
65	1	4	120	177	0	0	140	0	0.4	1	0	7	absent
56	1	3	130	256	1	2	142	1	0.6	2	1	6	present
59	1	4	110	239	0	2	142	1	1.2	2	1	7	present
60	1	4	140	293	0	2	170	0	1.2	2	2	7	present
63	0	4	150	407	0	2	154	0	4	2	3	7	present

	Female (0)	Male (1)	Sum (row)
FBS > 120 (0)	3 (2,7)	6 (6,3)	9
FBS > 120 (1)	0 (0,3)	1 (0,7)	1
Sum (col)	3	7	10

$$E_{00} = (3 \times 9) / 10 = 2,7$$

$$E_{01} = (3 \times 1) / 10 = 0,3$$

$$E_{10} = (7 \times 9) / 10 = 6,3$$

$$\chi^2 = \frac{(3-2,7)^2}{2,7} + \frac{(6-6,3)^2}{6,3} + \frac{(0-0,3)^2}{0,3} + \frac{(1-0,7)^2}{0,7} = 0,033 + 0,014 + 0,30 + 0,128 = 0,475$$

$$E_{11} = (7 \times 1) / 10 = 0,7$$

Faible Corrélation.

- 2) Calculer le coefficient de corrélation entre l'attribut âge (1<sup>ère</sup> colonne) et l'attribut *serum cholestoral* (5<sup>ème</sup> colonne).

Coefficient de Pearson :

$$r_{A,B} = \frac{\sum(A - \bar{A})(B - \bar{B})}{(n - 1)\sigma_A\sigma_B}$$

$$\bar{A} = 63$$

$$\bar{B} = 305$$

Calcul des  $\sigma_A$  et  $\sigma_B$  :

Age	$A - \bar{A}$	$(A - \bar{A})^2$	cholesterol	$B - \bar{B}$	$(B - \bar{B})^2$	$(A - \bar{A})(B - \bar{B})$
70	7	49	322	17	289	119
67	4	16	564	259	67081	1036
57	6	36	261	44	1936	264
64	1	1	263	42	1764	42
74	11	121	269	36	1296	396
65	2	4	177	128	16384	256
56	7	49	256	49	2401	343
59	4	16	239	66	4356	264
60	3	9	293	12	144	36
63	0	0	407	102	10404	0
<b>Total</b>		<b>301</b>			<b>106055</b>	<b>2756</b>

$$\sigma_A = \sqrt{\frac{301}{10}} = \sqrt{30} = 5$$

$$\sigma_B = \sqrt{\frac{106055}{10}} = \sqrt{10605} = 103$$

$$r_{A,B} = \frac{\sum(A - \bar{A})(B - \bar{B})}{(n - 1)\sigma_A\sigma_B} = \frac{2756}{9 \times 5 \times 103} = \frac{2756}{4635} = 0.595$$

Faible corrélation

- 3) Dessiner le scatter plot de l'attribut *âge* et de l'attribut *serum Cholestorol*.  
Que peut-on en conclure ?
- 4) Dessiner la boîte à moustache de l'attribut *serum cholestorol*. Que peut-on en conclure ?

Calcul des 5 nombres à savoir la médiane, les quantiles, le minimum et le maximum.

Trions au préalable les valeurs de l'attribut :

177, 239, 256, 261, 263, 269, 293, 322, 407, 564

Maximum = 564

Q3 = 322

Médiane = 266

Q1 = 256

Minimum = 177

$$\text{IQR} = 1.5 \times (322 - 256) = 99$$

Outliers :

$$322 + 99 = 421 \text{ donc : } \mathbf{564}$$

$$256 - 99 = 157$$