

3 Statistiques descriptives

La statistique est la science qui traite des éléments suivants :

- Récolte de données ;
- Classification des données ;
- Représentation et analyse des données (**statistiques descriptives**) ;
- Interprétations, conclusions et prévisions pouvant être tirées de l'analyse des données (**inférence statistique** au chapitre 5).

Aujourd'hui la statistique a pris une grande place grâce aux nouvelles techniques et à la puissance des ordinateurs. Géographie, médecine, sciences humaines, sciences économiques, biologie, politique, etc... : aucun domaine n'est épargné.

3.1 Caractéristiques d'un ensemble d'individus

Vocabulaire statistique :

- **Population** : ensemble de toutes les personnes, de tous les objets ou de tous les faits sur lesquels porte l'étude.
- **Individu** : élément de la population.
- **Sondage** : étude réalisée sur un **échantillon** de la population, soit un sous-ensemble de la population (lorsque la population est trop grande pour être étudiée dans son ensemble).
- **Variable statistique** (v.s.) : caractéristique étudiée dans la population.
- **Modalités ou valeurs** : différents états ou valeurs prises par la variable statistique.

Notation : on emploie une lettre majuscule (par exemple X) pour désigner une v.s. et une lettre minuscule (x_i) pour désigner une de ses valeurs.

Si l'échantillon est choisi au hasard, que sa taille est suffisamment grande et qu'il peut être considéré comme représentatif, on peut généraliser certains résultats obtenus à l'ensemble de la population.

Une variable statistique est dite **quantitative** si les valeurs qu'elle peut prendre sont des nombres. Dans le cas contraire, elle est **qualitative**, et ses valeurs s'appellent des **modalités** (ou des **catégories**).

Si les valeurs que peut prendre une variable statistique quantitative sont isolées les unes des autres, on dit que la variable est **discrète**.

Si les valeurs constituent des intervalles de nombres, on dit qu'elle est **continue**.

Modèle 6. Dans chaque situation exposée ci-dessous :

- 1) Décrire la population étudiée;
- 2) Décrire l'échantillon;
- 3) Nommer la variable étudiée;
- 4) Donner le type de variable étudiée.
- 5) Décrire l'ensemble des modalités ou des valeurs de la variable.

Exemples :

- a) On demande à 200 passagers pris au hasard dans un aéroport de donner leur nationalité.
 - 1) Population : ...
 - 2) Echantillon : ...
 - 3) Variable statistique : ...
 - 4) Type de variable : ...
 - 5) Modalités ou valeurs : ...
- b) Durant le mois d'avril, on mesure la température maximale de la journée à La tour-de-Peilz.
 - 1) Population : ...
 - 2) Echantillon : ...
 - 3) Variable statistique : ...
 - 4) Type de variable : ...
 - 5) Modalités ou valeurs : ...
- c) On demande à 5'000 familles résidant en Suisse le nombre d'enfants de leur foyer.
 - 1) Population : ...
 - 2) Echantillon : ...
 - 3) Variable statistique : ...
 - 4) Type de variable : ...
 - 5) Modalités ou valeurs : ...

3.2 Variables qualitative et quantitative discrète

Notation : on désigne la taille de l'échantillon par n et les valeurs prises par la v.s. X dans l'échantillon par x_1, x_2, \dots, x_n .

3.2.1 Tableau de distribution d'une variable qualitative

Une fois les données récoltées, on les regroupe par modalité dans un tableau de distribution.

Modèle 7. On a demandé à tous les élèves d'une classe quelle était leur matière préférée parmi les matières suivantes : français, anglais, maths et musique.

musique	français	français	anglais	français
anglais	musique	maths	musique	musique
musique	musique	musique	français	français
français	anglais	musique	anglais	maths

Tableau de distribution :

Répartition des selon

Matière	Effectif d'élèves = n_i	Fréquence = f_i
Total		

Remarques :

1. La somme des effectifs est toujours égale au nombre N d'individus de la population :

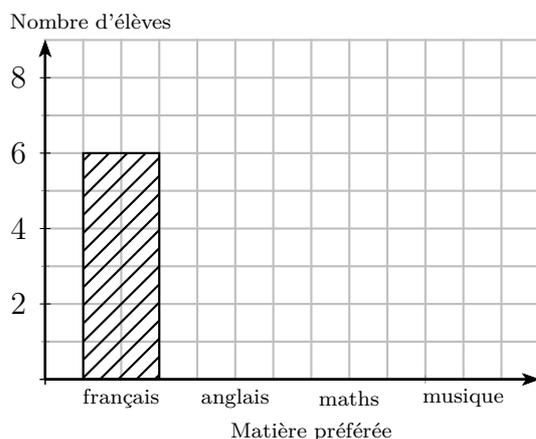
$$n_1 + n_2 + \dots + n_k = N$$

2. La somme des fréquences ($f_i = \frac{n_i}{N}$) est toujours égale à $1 = 100\%$.

3.2.2 Représentation graphique d'une variable qualitative

Modèle 8. Répartition des 20 élèves d'une classe selon leur matière préférée.

Diagramme à rectangles



3.3 Variables quantitative continue et discrète à grand nombre de valeurs

3.3.1 Tableau de distribution d'une variable quantitative continue

Modèle 11. On a demandé à ces élèves leur taille en centimètre.

172 157 162 156 167 179 173 173 178 160
 168 171 165 166 184 170 165 164 160 175

.....

Tableau de distribution :

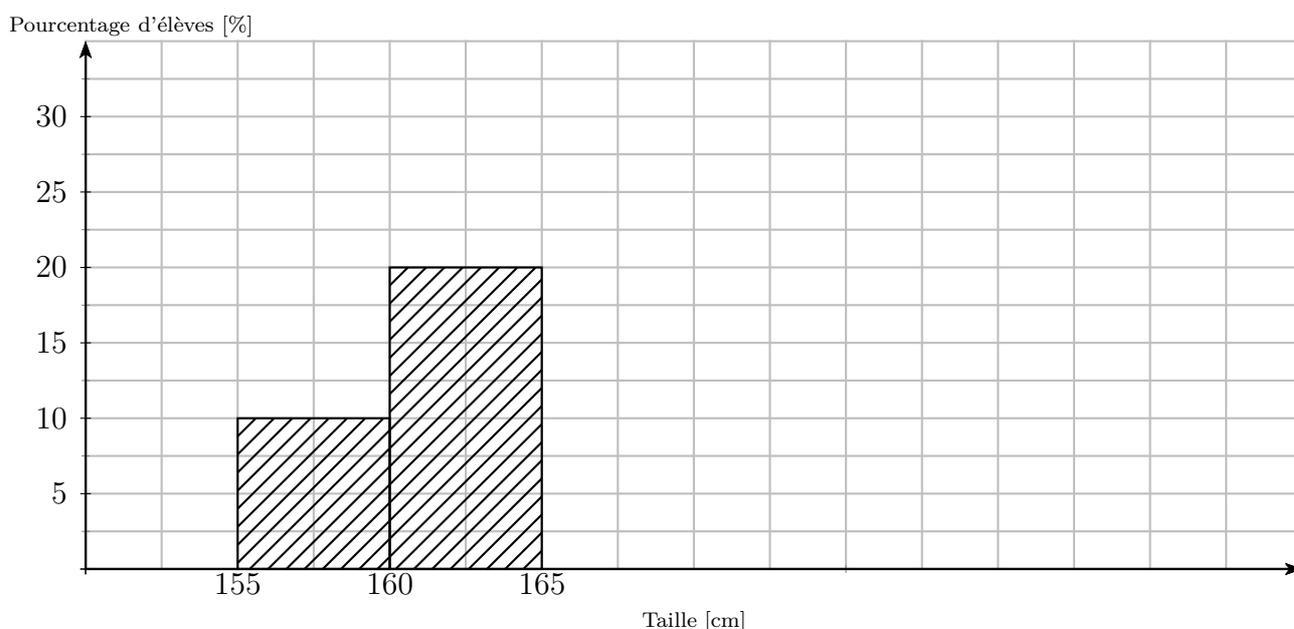
Répartition des selon

Taille	Effectif = n_i	Fréquence = f_i
Total		

3.3.2 Représentation graphique d'une variable quantitative continue

Modèle 12. Répartition des 20 élèves d'une classe selon leur taille.

Histogramme et polygone des fréquences



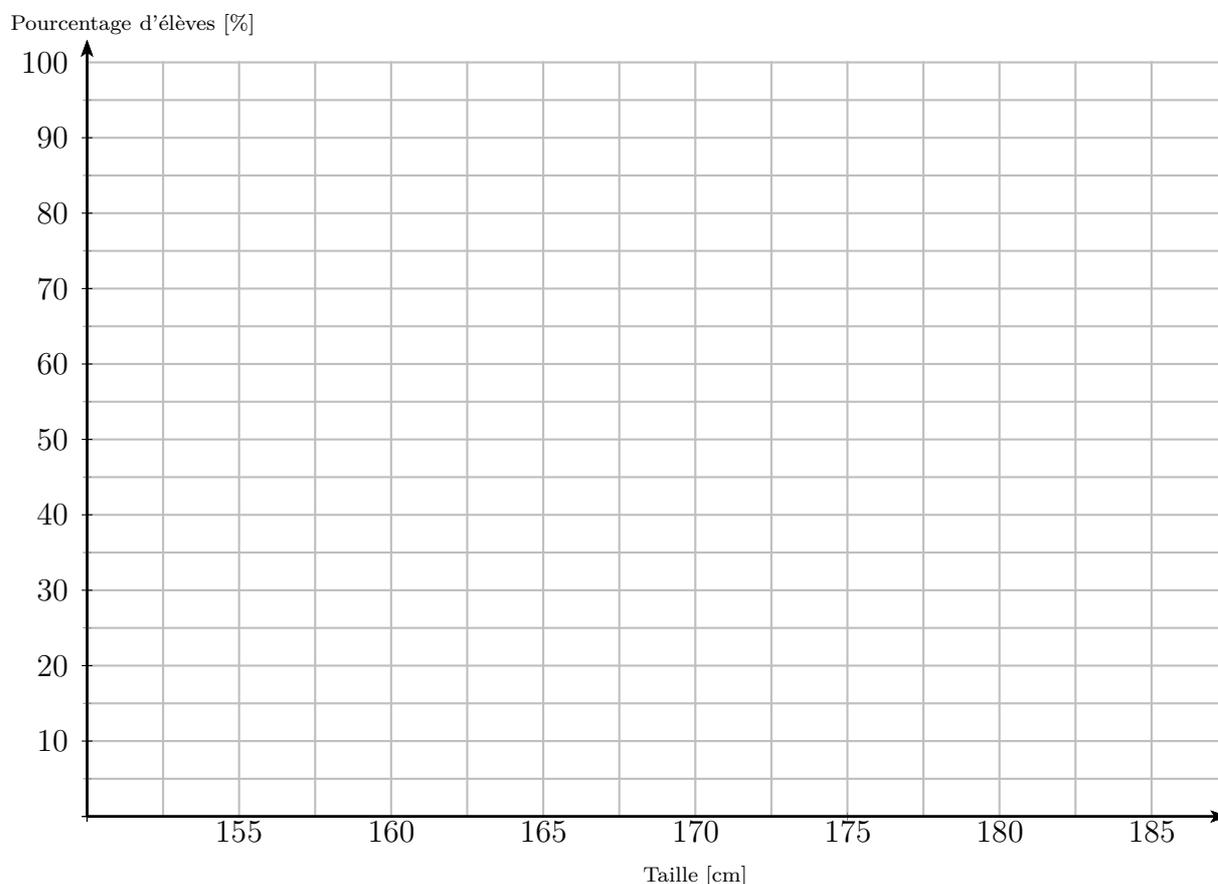
Le **polygone des fréquences** est la ligne polygonale obtenue en joignant les milieux respectifs des côtés supérieurs des rectangles de l'histogramme. On commence et on termine le polygone des fréquences en ajoutant une classe de fréquence nulle avant la première classe et une autre après la dernière classe.

Modèle 13. Répartition des 20 élèves d'une classe selon leur taille.

Polygone des fréquences cumulées (ou courbe de fréquences cumulées)

Taille	Effectif = n_i	Fréquence = f_i	Fréquence cumulée = F_i
[155 ; 160 [2	10 %	
[160 ; 165 [4	20 %	
[165 ; 170 [5	25 %	
[170 ; 175 [5	25 %	
[175 ; 180 [3	15 %	
[180 ; 185 [1	5 %	
Total	20	100 %	

polygone des fréquences cumulées



30 % des élèves mesurent

..... % des élèves mesurent moins de 1.75 mètre.

5 % des élèves mesurent au moins

Remarque

L'abscisse du point correspondant à une fréquence cumulée croissante de 50 % s'appelle la de la v.s. (ici \cong ).

3.4 Mesure de tendance centrale

But : résumer une série statistique par une seule valeur.

3.4.1 Moyenne

Calcul de la moyenne dans le cas discret

Modèle 14. On reprend les dernières notes d'anglais des 20 élèves d'une classe (modèle p.8).

5 4.5 3.5 5 6 3.5 4 2.5 4 4.5
4 4.5 4.5 4.5 3 4 4.5 5 3.5 4

Calcul de la moyenne arithmétique \bar{x} (se lit x barre) :

$$\bar{x} = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n}$$

$\bar{x} = \dots$

Calculatrice :

Calcul de la moyenne à partir du tableau de distribution (formule plus rapide) :

Répartition des 20 élèves d'une classe selon leur note d'anglais

Note (modalité) = c_i	Effectif = n_i	Fréquence = f_i
1 / 1.5 / 2	0	0 %
2.5	1	5 %
3	1	5 %
3.5	3	15 %
4	5	25 %
4.5	6	30 %
5	3	15 %
5.5	0	0 %
6	1	5 %
Total	20	100 %

$$\bar{x} = \frac{n_1 \cdot c_1 + n_2 \cdot c_2 + \dots + n_k \cdot c_k}{n} = f_1 \cdot c_1 + f_2 \cdot c_2 + \dots + f_k \cdot c_k$$

$\bar{x} = \dots$

Calculatrice :

Calcul de la moyenne dans le cas continu

Modèle 15. On reprend les tailles des 20 élèves d'une classe (modèle p.9).

172 157 162 156 167 179 173 173 178 160
 168 171 165 166 184 170 165 164 160 175

Calcul direct de la moyenne (à partir des données brutes) :

$$\bar{x} = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n}$$

$\bar{x} = \dots$

Calcul de la moyenne à partir du tableau de distribution (formule plus rapide) :

Répartition des 20 élèves d'une classe selon leur taille

Taille	Valeur centrale = c_i	Effectif = n_i	Fréquence = f_i	Fréquence cumulée = F_i
[155 ; 160 [2	10 %	10 %
[160 ; 165 [4	20 %	30 %
[165 ; 170 [5	25 %	55 %
[170 ; 175 [5	25 %	80 %
[175 ; 180 [3	15 %	95 %
[180 ; 185 [1	5 %	100 %
Total		20	100 %	

$$\bar{x} = \frac{n_1 \cdot c_1 + n_2 \cdot c_2 + \dots + n_k \cdot c_k}{n} = f_1 \cdot c_1 + f_2 \cdot c_2 + \dots + f_k \cdot c_k$$

$\bar{x} = \dots$

Remarque

La valeur de la moyenne n'est pas la même selon la méthode utilisée. Dans le deuxième cas, on ne dispose plus des données brutes, on doit donc estimer la valeur de \bar{x} avec l'information disponible.

Représentation graphique de la moyenne

On place un pivot (▲) sous l'axe horizontal, au point correspondant à la moyenne, celle-ci coïncide avec le centre d'équilibre du diagramme. Les modèles 10 et 12 sont à compléter.

3.4.2 Médiane

Calcul de la médiane dans le cas discret

La médiane partage une série de données triées en deux parties égales.

Si \tilde{x} est la médiane d'une série statistique, il y a donc 50 % des données qui sont plus petites ou égales à \tilde{x} et 50 % qui sont plus grandes ou égales à \tilde{x} .

Modèle 16. Dans l'exemple des notes d'anglais des 20 élèves d'une classe (modèle p.8).

Calcul de la médiane \tilde{x} (se lit x tilde) :

$$\tilde{x} = \begin{cases} x_{\frac{n+1}{2}} & \text{si } n \text{ est impair} \\ \frac{1}{2} (x_{\frac{n}{2}} + x_{\frac{n}{2}+1}) & \text{si } n \text{ est pair} \end{cases}$$

Il y a deux formules différentes, car si n est impair et les données sont triées, il y a une donnée au milieu de la série. Si n est pair, par contre, il y a deux données qui sont au milieu, et on utilise donc la moyenne arithmétique de ces deux valeurs.

Remarque importante :

Cette mesure de tendance centrale est plus robuste que la moyenne, elle est moins affectée par les valeurs extrêmes.

Exemple :

Dans une entreprise de 35 employés, supposons que le patron gagne 40'000 francs par mois, alors que les 34 employés gagnent 3'000 francs par mois.

Calcul du revenu mensuel moyen : $\bar{x} = \dots$

Cette moyenne ne reflète en rien la réalité des travailleurs de cette entreprise. La valeur extrême du salaire du patron a un impact trop grand sur la moyenne.

Calcul de la médiane : $\tilde{x} = \dots$

Il est correct de dire que le salaire moyen dans cette entreprise est de mais il vaudrait mieux dire que le salaire **médian** est de par mois, donc qu'au moins la moitié des employés gagnent

Calcul de la médiane dans le cas continu

Modèle 17. Dans l'exemple des tailles des 20 élèves d'une classe (modèle p.9).

Calcul direct de la médiane (à partir des données brutes) :

$$\tilde{x} = \frac{1}{2} (x_{\frac{n}{2}} + x_{\frac{n}{2}+1})$$

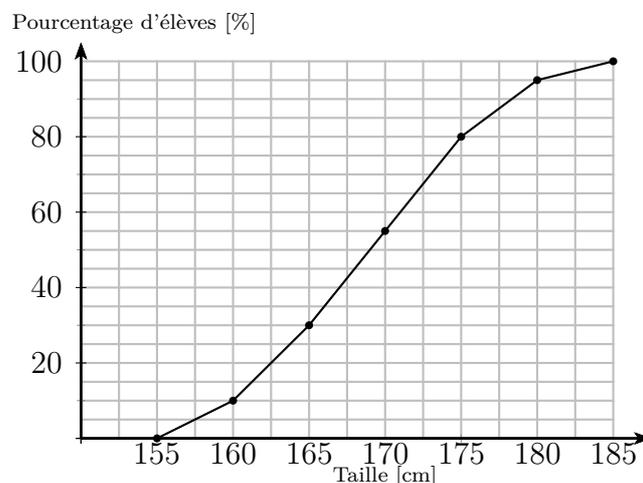
$\tilde{x} = \dots$

Si on ne dispose que du tableau de distribution, on doit estimer la médiane autrement.

Répartition des 20 élèves d'une classe selon leur taille

Polygone de fréquences cumulées

Taille	Centre	Effectif	Fré.	Fré.cum.
[155 ; 160 [157.5	2	10 %	10 %
[160 ; 165 [162.5	4	20 %	30 %
[165 ; 170 [167.5	5	25 %	55 %
[170 ; 175 [172.5	5	25 %	80 %
[175 ; 180 [177.5	3	15 %	95 %
[180 ; 185 [182.5	1	5 %	100 %
Total		20	100 %	



Classe médiane :

Calcul de la médiane par proportionnalité entre les fréquences et les valeurs :

Remarque

La valeur de la médiane n'est pas la même selon la méthode utilisée. Dans le deuxième cas, on ne dispose plus des données brutes, on doit donc estimer la valeur de \tilde{x} avec l'information disponible.

3.4.3 Mode et classe modale

Le mode est la valeur (ou la catégorie dans le cas qualitatif) qui revient le plus souvent dans une série statistique.

La classe modale est la classe qui regroupe le plus de données dans le cas d'une variable continue.

Remarques

1. Le mode ou la classe modale ne sont significatifs que si leur effectif est largement plus grand que celui des autres modalités ou des autres classes.
2. Le mode est la seule mesure de tendance centrale qui peut être utilisée pour une variable qualitative.

Modèle 18. Déterminer le mode pour chacune des parties :

a) Dans l'exemple des matières préférées, le mode est

Interprétation :

b) Dans l'exemple des notes d'anglais, le mode est

Interprétation :

c) Dans l'exemple des tailles, la classe modale est

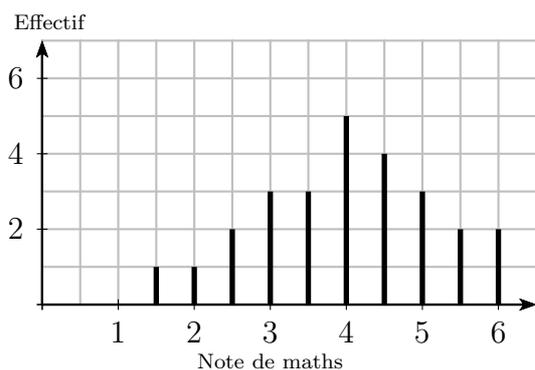
Interprétation :

3.5 Mesures de dispersion

Lorsqu'on résume une série statistique par une mesure de tendance centrale (souvent la moyenne), on ne donne aucune information sur la manière dont les données se répartissent autour de cette valeur : sont-elles toutes assez proches de la moyenne, ou trouve-t-on des valeurs très dispersées autour de celle-ci ? Cette question nécessite de donner une valeur supplémentaire, appelée mesure de dispersion.

Pour illustrer ces mesures de dispersion, nous allons nous baser sur les notes de maths de trois classes parallèles, données par des diagrammes en bâtons.

Classe A

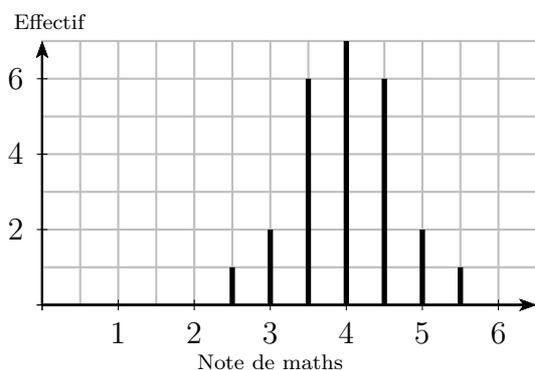


Moyenne :

Médiane :

Mode :

Classe B

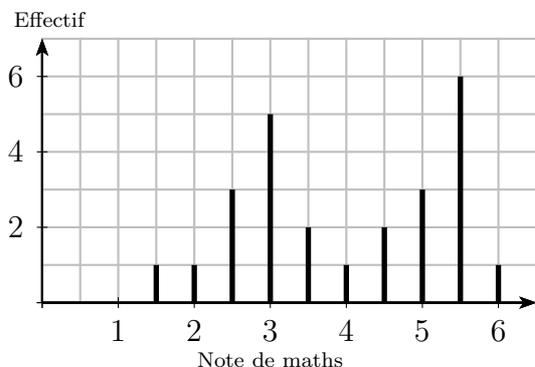


Moyenne :

Médiane :

Mode :

Classe C



Moyenne :

Médiane :

Mode :

Remarque

Ces mesures de tendance centrale ne suffisent pas à décrire les différences entre ces trois classes.

3.5.1 Etendue

Modèle 19. L'étendue est la "distance" entre la plus petite et la plus grande valeur.

Classe	A	B	C
Etendue			

Cette mesure permet de différencier les situations des classes ... et ... , mais pas des classes ... et ...

3.5.2 Variance et écart-type

La variance s^2 est l'écart quadratique moyen à la moyenne. Elle se calcule par la formule suivante :

$$s^2 = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n}$$

Dans le cas de données regroupées par modalités ou par classes, la formule devient :

$$s^2 = \frac{n_1(c_1 - \bar{x})^2 + n_2(c_2 - \bar{x})^2 + \dots + n_k(c_k - \bar{x})^2}{n} = f_1(c_1 - \bar{x})^2 + f_2(c_2 - \bar{x})^2 + \dots + f_n(c_k - \bar{x})^2$$

Comme la variance est calculée à partir de grandeurs au carré, on définit l'écart-type, noté s , comme la racine de la variance. On obtient ainsi une mesure de la dispersion dans la même unité que les mesures initiales.

Calculatrice :

Modèle 20. Calculer la variance et l'écart-type des classes A , B et C à l'aide des diagrammes en bâtons de la page 16.

Classe	A	B	C
Variance			
Ecart-type			

Grâce à ces nouvelles mesures, on peut maintenant affirmer que les notes de la classe ... sont les moins dispersées autour de la moyenne, et que celles de la classe ... sont les plus dispersées.

3.6 Mesure de position

3.6.1 Quantiles

La médiane est une valeur qui partage les données de l'échantillon en deux groupes de taille égale : 50% des données sont inférieures ou égales à la médiane, et 50% des données lui sont supérieures ou égales.

Cette idée se généralise pour n'importe quel pourcentage. Par exemple, quelle est la valeur qui sépare les 25% les plus petits des 75% les plus grands ?

Un quantile à $p\%$ est une valeur qui est supérieure ou égale aux $p\%$ des données les plus petites, et inférieure ou égale au reste des données. On le note $q_p\%$.

Cas particuliers

- Les quartiles (Q_1, Q_2, Q_3) sont les quantiles à 25%, 50% et 75%. Ils partagent les données en quatre parties égales. Le deuxième quartile (Q_2) est égal à la médiane.
- Les quintiles (V_1, V_2, V_3, V_4) sont les quantiles à 20%, 40%, 60% et 80%. Ils partagent les données en cinq parties égales.
- Les déciles (D_1, D_2, \dots, D_9) sont les quantiles à 10%, 20%, ..., 90%. Ils partagent les données en 10 parties égales.
- Les centiles (C_1, C_2, \dots, C_{99}) sont les quantiles à 1%, 2%, ..., 99%. Ils partagent les données en cent parties égales.

Les quantiles se déterminent en utilisant le même principe que pour la médiane.

Remarque

Pour que les quantiles aient du sens, il faut que l'échantillon soit suffisamment grand. On ne calculera jamais le premier décile d'une distribution composée d'une dizaine de valeurs !

3.6.2 Boxplot

Un boxplot (ou boîte à moustache) est une manière de représenter graphiquement la distribution d'une variable statistique en faisant apparaître la médiane, les quartiles et les deux valeurs extrêmes (la plus petite et la plus grande).

Modèle 21. On suppose que le nombre de périodes d'absences par année des élèves d'un gymnase se répartit de la manière suivante :

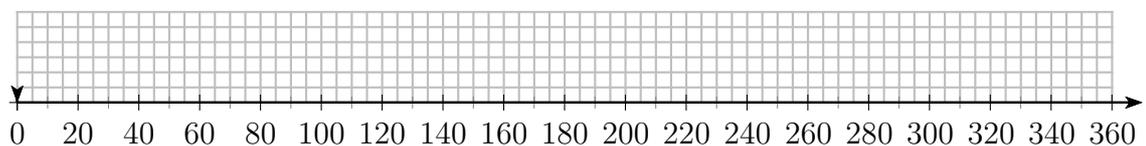
Périodes d'absence	$[0 ; 10 [$	$[10 ; 30 [$	$[30 ; 60 [$	$[60 ; 90 [$	$[90 ; 120 [$	$[120 ; 360 [$	Total
Fréquence	14 %	60 %	15 %	7 %	2 %	2 %	100 %

a) Calcul de la médiane (Q_2) :

b) Calcul du premier quartile (Q_1) :

c) Calcul du troisième quartile (Q_3) :

d) Box-Plot :



3.6.3 Cote Z

Modèle 22. Un gymnase souhaite engager un ancien étudiant pour donner des cours d'appui de mathématiques. Les quatre candidats ont suivi leur troisième année dans quatre gymnases différents, mais on souhaite tout de même déterminer le meilleur étudiant en fonction de ses résultats à l'examen de maturité.

Candidat	Nom de l'élève	Note moyenne de son gymnase	Ecart-type de son gymnase
Loïc	4.5	3.7	1.1
Muriel	5	4.1	0.6
Antonin	5.5	5.1	0.4
Eloïse	5	4.0	0.9

Si le gymnase ne se fie qu'à la note de l'élève, il devrait engager ...

S'il tient aussi compte de la moyenne du gymnase, il engagera plutôt ...

Enfin, en tenant compte de l'écart-type du gymnase, il choisira alors ...

Pour décrire la position d'une donnée par rapport à une distribution, on utilise la cote Z .

$$\text{Cote } Z \text{ de } x_i : z_{x_i} = \frac{x_i - \bar{x}}{s}$$

La cote Z mesure la distance d'une valeur à la moyenne, mesurée en nombre d'écart-type.

Cote Z de Loïc :

Cote Z de Muriel :

Cote Z d'Antonin :

Cote Z d'Eloïse :

Interprétation de la cote Z

Une cote Z positive signifie que la valeur est supérieure à la moyenne, alors qu'une cote Z négative indique qu'elle est en dessous de la moyenne.

Une cote Z de 3 ou plus, ou de -3 ou moins indique une valeur très rare. La cote Z permet donc d'identifier des situations exceptionnelles ou peu plausibles.

Modèle 23. Un cinéma accueille en moyenne 120 spectateurs les soirs de semaine, avec un écart-type de 14 spectateurs. Il décide de proposer une offre spéciale le mardi soir, avec des places à tarif réduit. Le mardi suivant, 172 spectateurs assistent à la projection. Peut-on déduire que l'offre spéciale a eu de l'effet ?

Un lundi soir, une exposition a lieu tout près du cinéma. Ce même soir, le cinéma vend 104 billets. Le gérant se plaint de l'effet négatif de l'exposition, qui lui aurait "volé" des clients. Est-ce justifié ?