

Statistiques descriptives

Chapitre I: I. vocabulaire

les statistiques consistent en diverses méthodes de classement des données
à l'aide des Tableaux, les histogrammes et les graphiques permettant
d'organiser un grand nombre de données. Elles sont dotées d'un vocabulaire
particulier.

1. Expériences statistiques:

des statistiques descriptives visent à étudier les caractéristiques
d'un ensemble d'observations comme les données obtenues lors d'une expérience.

l'expérience est étape préliminaire à toute statistique.

Il s'agit "d'extraire" avec les observations de données générales, la méthode
statistique est basée sur la compréhension de ces données.

l'objectif principal des statistiques est de donner une explication de ce qu'on
observe.

Ex 1: "la durée de vie d'une lampe" dans un magasin de bricolage.

un fabricant de lampes veut savoir la durée de vie de ses lampes.

se propose d'étudier

Pour ce faire, il va faire fabriquer la relation d'âge de la lampe en fonction

de la durée de vie de la lampe. Pour ce faire, la durée de vie de la lampe est la

Pour organiser les résultats obtenus

2. Population :

En statistique, on s'intéresse aux populations. Le terme vient du fait que la démographie étudie des populations humaines à l'échelle nationale ou d'état de la statistique notant au travers des recensements des populations.

Nous définissons la notion de population :

Définition 1 :

On appelle population l'ensemble sur lequel porte notre étude statistique. Cet ensemble est noté Ω

Expl 1 : On considère l'ensemble des étudiants de la section 1. On s'intéresse au nombre de frères et sœurs de chaque étudiant dans la classe. Ω est l'ensemble des étudiants.

3. Individu (unité statistique)

une population est composée d'individus.

les individus qui composent une population statistique sont appelés unités statistiques.

Définition 2 : On appelle individu tout élément de la population. Elle est notée ω (ω dans Ω)

Remarque 1:

Ensemble \mathcal{N} peut être un ensemble de personnes, de doses ...

L'unité statistique est un objet pour lequel nous sommes intéressés à acquies l'information

4. Caractères : (variable statistique)

La statistique descriptive étudie comment l'individu de la donnée une population donnée. nous nous intéressons aux caractéristiques des individus qui peuvent prendre des plusieurs valeurs.

Définition 4: On appelle caractère ou variable statistique désignée V.S

Toute application $X: \mathcal{N} \rightarrow C$, l'ensemble C est dit ensemble des valeurs des caractères X
C'est à dire est mesuré par chaque un des individus

Ex)

taille, l'âge, le poids, le nombre d'enfants

Remarque 2: soit \mathcal{N} un ensemble, On appelle et on note Cardinal de \mathcal{N}

$\text{Card}(\mathcal{N})$: le nombre d'éléments de \mathcal{N}

$\text{Card}(\mathcal{N})$ nombre d'éléments de $\mathcal{N} = N$

Question 6: L'ensemble des valeurs est représenté par des chiffres, et même, elle est notée en 2 sortes de caractères

↗ direct
↘ continue

Exemple 1: Le salaire d'employés d'une usine:

Modalités: 2000D, 3000D

Type: Discret

La longueur des versements:

Modalités: [10, 20] N/m

Type: continu

III - Etude d'une variable statistique discrète

Une variable statistique peut prendre un nombre fini (borné) de valeurs (nombre d'élèves, nombre de jours, nombre). Dans ce cas, la variable statistique est appelée variable discrète. Dans toute la suite de ce chapitre, nous noterons la suite de valeurs $X: \omega \mapsto \{x_1, x_2, \dots, x_m\}$ avec $\text{card}(X) = m$ est le nombre d'individus dans notre étude.

Nous allons utiliser souvent les x_i à l'avenir pour illustrer les notions de ce chapitre.

Exemple 2: une enquête réalisée dans un village porte sur le nombre d'élèves à l'école primaire. On note X : nombre d'élèves les parents ont donné par le tableau

x_i	0	1	2	3	4	5	6
n_i	8	32	66	41	32	9	2

Notations :

Ω : ensemble des ω

W : une famille

X : nbre d'effets payés

$$X: W \rightarrow \mathbb{N}$$

On dit que :

à la famille w , on associe $X(w)$ = le nbre d'effets de cette famille

III - 2 Effet partiel - effectif cumulé

On étudie un caractère statistique mesurable représenté par un site

l'expérimentation (2) donne la valeur de caractéristique avec $x \in \{1, 2, \dots, k\}$

III.4 - 1 - Effet partiel (fraction absolue)

Def 7: Pour chaque valeur x_i on pose par définition :

$$n_i = \text{card} \{ \omega \in \Omega, X(\omega) = x_i \}$$

on dit que :

on appelle effet partiel de x_i

et le noter n_i

III.4.2 - 2 - Effet cumulé

Def 8: Pour chaque valeur x_i on pose par définition $N_i = n_1 + n_2 + \dots + n_i$

L'effet cumulé d'une valeur est la somme d'effets d'une valeur et toutes les valeurs inférieures à elle

Effet de valeur x_i dépend

Exemple 9:

Dans l'exemple précédent 50 est le nombre de familles ayant un nombre d'enfant inférieur à 1. Nous le regardons dans le tableau suivant:

x_i	0	1	2	3	4	5	6
N_i	18	50	116	157	189	198	200

Interprétation: N_i est le nombre d'individus dont la valeur du caractère est inférieur ou égale à x_i . De ce fait, l'effectif total est donné par

$$N = \text{Card}(\Omega) = \sum_{i=1}^n x_i$$

Dans notre exemple précédent. Nous avons $N = 200$

III - 1 - 3 - Fréquence partielle - Fréquence Cumulée:

Définition 9:

• Pour chaque valeur x_i , on pose par définition $f_i = \frac{n_i}{N}$
 f_i s'appelle la fréquence partielle de x_i

• La fréquence d'une valeur est le rapport de l'effectif de cette valeur par l'effectif totale.

Remarque:

On peut remplacer f_i par $f_i \times 100$ qui représentent alors un pourcentage.

Interprétation:

f_i : c'est le pourcentage de w tel que $X(w) = x_i$

Exemple 10:

Quel est le pourcentage de famille dont le nombre d'enfants égal à 1?

Solution: $f_i = \frac{66}{200} = 0,33$

Proposition:

Soit f_i défini comme précédemment Alors

$$\sum_{i=1}^n f_i = 1$$

Dém:

$$\sum_{i=1}^n f_i = \sum_{i=1}^n \frac{n_i}{N} = \frac{1}{N} \sum_{i=1}^n n_i = \frac{N}{N} = 1$$

- III - 1 - 4 - Fréquence Cumulée :

Définition: Pour chaque valeur x_i , on pose par définition

$$F_i = f_1 + f_2 + \dots + f_i$$

la quantité F_i s'appelle la fréquence Cumulée

Interprétation: F_i est la pourcentage de ω tel que la valeur $x(\omega)$ est inférieur ou égale à x_i .

Exemple 11:

* Dans l'exemple précédent 0,785 représente 78,5% de famille dont le nombre d'enfant est inférieur ou égale à 3.

* Dans un Deuxième exemple, nous nous intéressons aux nombres d'erreurs d'assemblage sur un ensemble d'appareils

Nombre d'erreurs	Nombre d'appareils	Fréquence cumulée
0	101	0,26
1	140	0,61
2	92	0,84
3	42	0,94
4	18	0,99
5	3	1

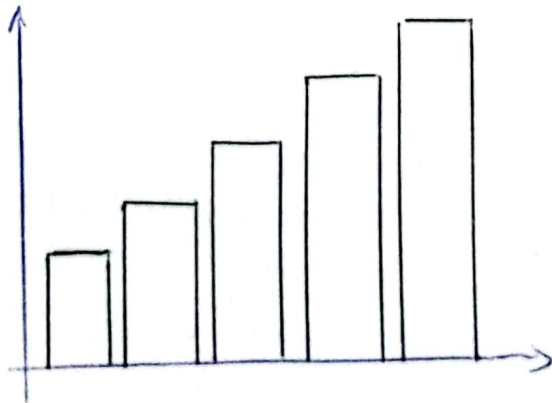
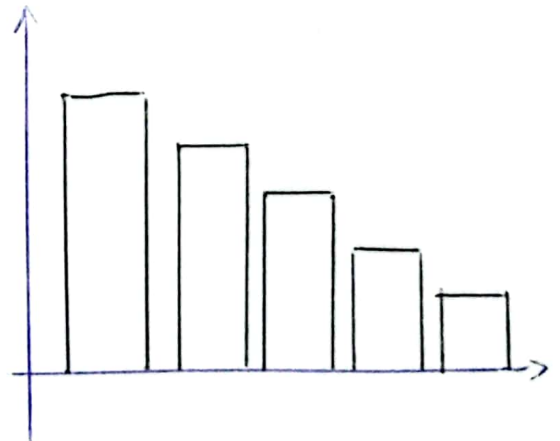
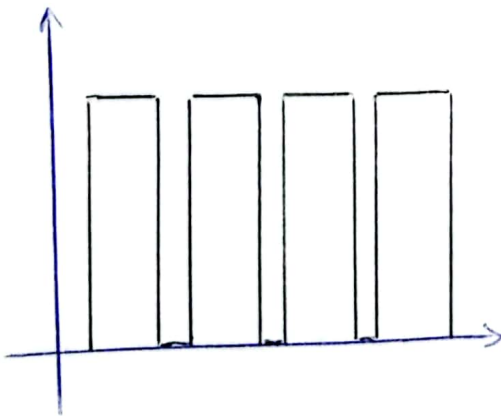
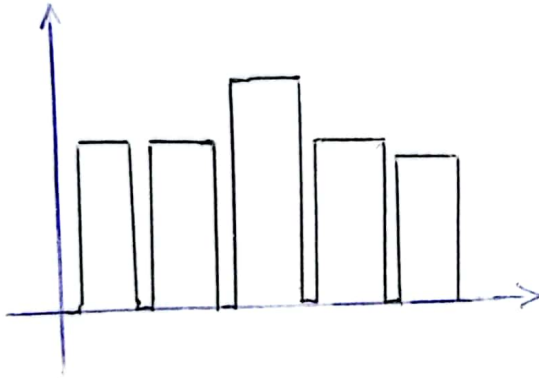
(8)

Représentation graphique des séries statistiques :

distingue la méthode de représentation d'une variable statistique en fonction de la nature de cette variable (qualitative ou quantitative).
Les représentations recommandées et le plus fréquent sont les tableaux et les diagrammes (Graphiques).

Le graphique est un support visuel qui permet :

- La synthèse : visualise d'un seul le principal caractéristique (mais on perd une quantité d'information)



- La découverte

Mettre en évidence tendance

- Le contrôle: on aperçoit les anomalies sur un graphique dans un tableau

- La recherche des régularités:

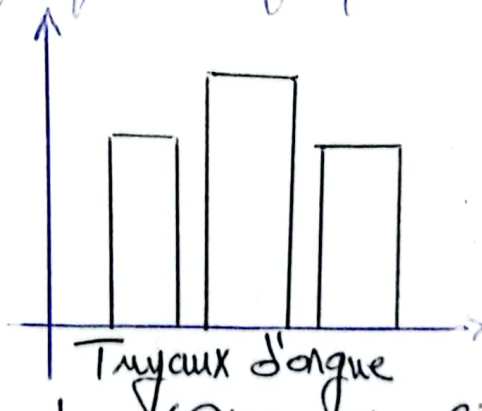
Régularités dans le mouvement, répétition du phénomène

- III.2 Distribution à caractère qualitatif:

À partir de l'observation d'une variable qualitative, les diagrammes permettent de représenter cette variable: les diagrammes en bandes.

(dit ~~(travaux d'orgue)~~ **travaux d'orgue**) et le diagramme à secteurs angulaire (dit **Camembert**)

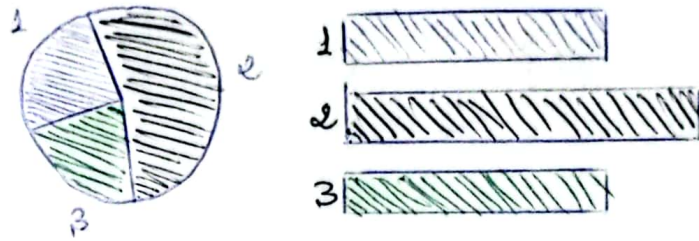
* Travaux d'orgue: Nous portons en abscisse les modalités, de façon arbitraire, nous portons en ordonnées (f) des rectangles dont la longueur de ses effectifs, ou de fréquence de chaque modalité.



* Diagramme par secteur (Diagramme Circulaire):

Les diagramme circulaires ou semi-circulaires consistent à partager un disque ou un demi-disque en tranche ou secteur

proportionnel au mode observé et dans la surface est proportionnel à l'effectif ou à la fréquence de la modalité.



* Diagramme par secteur :

Le degré d'un secteur est déterminé à l'aide de la règle de trois de la manière suivante :

$$N \longrightarrow 360^\circ$$

$$n_i \longrightarrow d_i \text{ (degré de la modalité)}$$

donc

$$d_i = \frac{n_i \times 360}{N}$$

2-2 - Distribution de Caractère quantitatif discret :

À partir de l'observation d'une variable quantitative discrète, deux diagrammes permettent de représenter cette variable le diagramme en bâtons et le diagramme Cumulatif.

Pour l'illustration, nous prenons l'exemple précédent de départ (nombre d'enfants par famille) nous rappelons le tableau statistique associé.

x_i	0	1	2	3	4	5	6
n_i	13	32	66	41	32	9	2

Diagramme à bâtons :

on veut représenter cette répartition sous la forme d'un diagramme

en bâtons. À chaque marque correspond un bâton. Les hauteurs sont proportionnelles aux effectifs représentés.

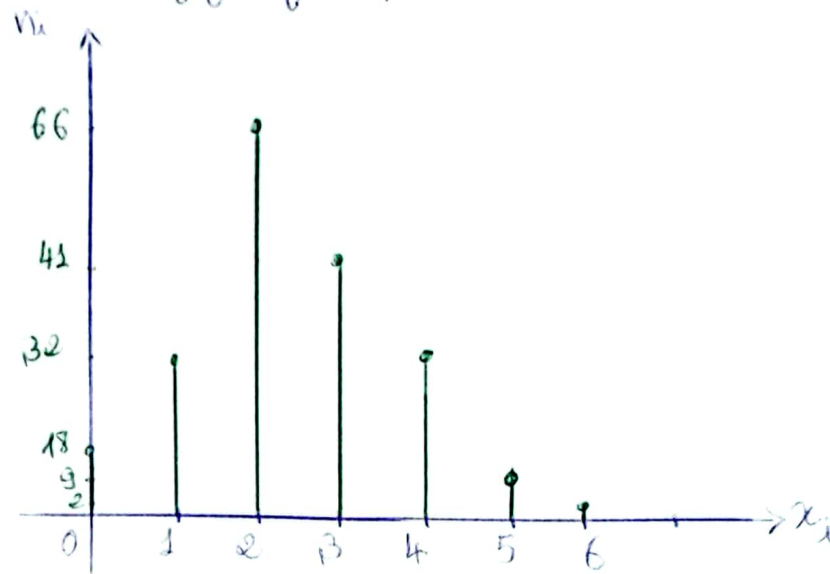


Diagramme à bâtons

2-3-4. Représentation sous forme de Courbe et fonction de répartition:

Nous avons déjà abordé la distribution cumulée d'une variable statistique. Nous allons dans cette partie exploiter ces valeurs cumulée pour introduire la notion de fonction de répartition. Cette fonction ne concerne que les variables quantitatives.

Soit la fonction $F_x: \mathbb{R} \rightarrow [0, 1]$, définie par: $F_x(x) = \text{pourcentage des individus dont la valeur de caractère } \leq x$

$F_x(x) =$ pourcentage des individus dont la valeur de caractère $\leq x$
 Cette fonction s'appelle la **fonction de répartition** d'un caractère x

Remarque:

Pour tout $i \in \{1, \dots, n\}$ on a $F_x(x_i) = F_i$

La courbe de F_x passe par les points: $(x_1, F_1), (x_2, F_2), \dots$ et (x_n, F_n)

En se basant sur notre exemple, la courbe F_x représentée ci-dessous sur

$$\mathbb{R} =]-\infty, 0[\cup [0, 2[\cup \dots \cup [6, +\infty[$$

① Dans ce cas nous avons:

- si $x < 0$, alors $F_n(x) = 0$

- si $x \in [0, 1[$, alors $F_n(x) = 0,09$ ($= \frac{18}{200}$)

- si $x \in [1, 2[$; $F_n(x) = \frac{36}{200} = 0,18 + 0,09 = 0,27$

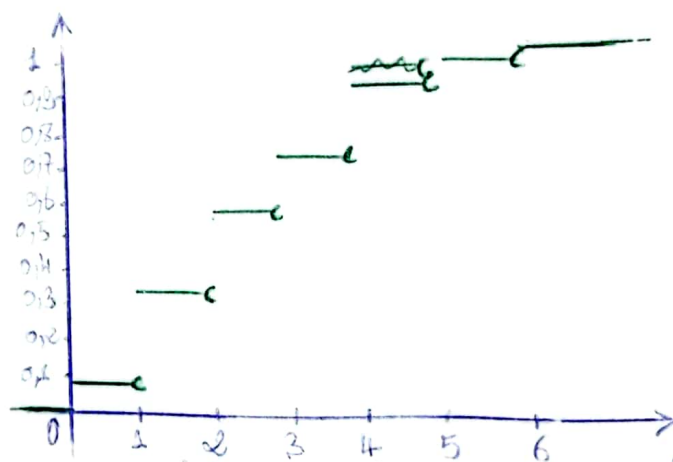
- si $x \in [2, 3[$; $F_n(x) = F_n(x) = \frac{66}{200} = 0,33 + 0,27 = 0,60$

- si $x \in [3, 4[$; $F_n(x) = 0,205 + 0,58 = 0,785$

- si $x \in [4, 5[$; $F_n(x) = 0,16 + 0,785 = 0,945$

- si $x \in [5, 6[$; $F_n(x) = 0,045 + 0,945 = 0,99$

- si $x \geq 6$; $F_n(x) = 0,99 + 0,01 = 1$



• Représentation d'une variable quantitative discrète par la fonction cumulative

Proposition:

La fonction de répartition satisfait pour $i \in \{1, \dots, n\}$

- l'égalité: $F_n(x_i) = F_i$

- l'expression:
$$F_n(x) = \begin{cases} 0 & \text{si } x < x_0 \\ F_1 & \text{si } x_0 \leq x < x_1 \\ F_i & \text{si } x_i \leq x < x_{i+1} \\ 1 & \text{si } x \geq x_n \end{cases}$$

- 2-4 Paramètre de Position: (Caractéristique de tendance Centrale)

Les indicateurs statistiques de tendance Centrale (dit aussi de Position)

Considère fréquemment sont la moyenne, la médiane et le mode.

• Le mode: Le mode d'une variable statistique est la valeur qui a le plus grand effectif partiel (par la plus grande fréquence partiel) il est noté M_o .

Exemple 11: Dans l'exemple précédent le mode est égale à la valeur qui correspond au plus grand effectif.

Remarque: on peut avoir plus d'un mode ou rien.

• La médiane: on appelle médiane la valeur M_e de la variable statistique qui vérifie la relation suivante $F_n(M_e^-) < 0,5 \leq F_n(M_e^+) = F_n(M_e)$

La médiane partage la série statistique en deux groupe de même effectif.

Exemple 12:

* $F_n(0) = 0 < 0,5$ ~~$F_n(0^+) = 0,03$~~

n'est pas satisfait donc la médiane est différente de 0.

* $F_n(2^-) = 0,25 < 0,5 \leq F_n(2^+) = 0,58$

donc $M_e = 2$

• La moyenne: on appelle moyenne de X la quantité

avec $N = \text{Card}(O)$.

$$\bar{x} = \frac{1}{N} \sum_{i=1}^n n_i x_i = \sum_{i=1}^n F_i x_i$$

On peut donc exprimer et calculer la moyenne dite "Arithmétique" avec des effectifs ou avec des fréquences.

Exemple 13: si $\bar{x} = 2,46$, alors nous avons en moyenne une famille à 2,46 enfant.

• La valeur de la moyenne est abstraite. Comme dans l'exemple précédent $\bar{x} = 2,46$ est un chiffre qui ne correspond pas à un fait concret.

• La moyenne arithmétique dont on vient de donner la formule est dite moyenne pondérée: c'est la signification de chaque valeur de la variable.

par la Coefficient, ici par l'effectif n_i qui correspond, dans ce cas chaque valeur x_i de la variable intervient sans le calcul de la moyenne, autant de fois qu'elle est observée. On parle de la moyenne arithmétique simple quand on effectue par la configuration. Par exemple si 5 étudiants en ~~pourage~~ ^{âge} respectivement 18, 19, 20, 21 et 22 ans, leur âge moyenne est donné par :

$$\frac{18+19+20+21+22}{5} = 20 \text{ ans (moyenne simple).}$$

Remarque : nous mentionnons qu'il existe d'autre moyenne que la moyenne arithmétique.

- 2 - 5 - Paramètre de dispersion: (variabilité)

Les indicateurs statistiques de dispersion usuel sont : l'étendue la variance et l'écart-type.

La différence entre la plus grande valeur et la plus petite valeur du caractère donner par la quantité $e = x_{\max} - x_{\min}$ s'appelle **l'étendue de la variable statistique x** .

Le calcul de l'étendue est très simple, il donne une première idée de la dispersion des observation.

* **La variance** : on appelle variance de cette série statistique x , le nombre

$$\text{Var}(X) = \sum_{i=1}^n f_i (\bar{x} - x_i)^2$$

on dit que la variance est la moyenne des carrés des écart à la moyenne de \bar{x} .

Le théorème suivant (**Théorème de König - ~~Huygens~~ Huygens**) donne une identité remarquable reliant la variance et la moyenne, parfois plus pratique dans le calcul de la variance

Théorème:

Soit (x_i, n_i) une série statistique de moyenne \bar{x} et de variance $\text{Var}(X)$

Alors
$$\boxed{\text{Var}(X) = \sum_{i=1}^n f_i x_i^2 - \bar{x}^2}$$

Dém:

$$\text{Var}(X) = \sum_{i=1}^n f_i (\bar{x} - x_i)^2 = \sum_{i=1}^n f_i (\bar{x}^2 - 2\bar{x}x_i + x_i^2)$$

$$\sum_{i=1}^n f_i = \sum_{i=1}^n f_i \bar{x}^2 - 2\bar{x} \sum_{i=1}^n x_i f_i + \sum_{i=1}^n f_i x_i^2$$

$$\sum_{i=1}^n f_i x_i = \bar{x}^2 - 2\bar{x}^2 + \sum_{i=1}^n f_i x_i^2 = \sum_{i=1}^n f_i x_i^2 - \bar{x}^2$$

Remarque:

Dans l'utilisation de la formule du théorème précédent (1) il faut veiller à remplacer \bar{x} par sa valeur approchée la plus précise aussi.

• L'écart-type:

La quantité $\sigma_X = \sqrt{\text{Var}(X)}$ s'appelle l'écart type de V.S. X
v.a. statistique

Remarque:

Les paramètres σ_X mesurent la distance moyenne entre \bar{x} et la valeur de X . Il sert à mesurer la dispersion d'une série statistique autour de sa moyenne.

- Plus il est petit, plus les caractères sont concentrés autour de la moyenne (on dit que la série est homogène).

- Plus il est grand, plus les caractères sont dispersés autour de la moyenne (on dit que la série est hétérogène).