

Partie II = statistique

Chapitre I: Généralité et notions de base

Definition: La statistique est une méthode scientifique qui consiste à réunir des données chiffrées sur des ensembles nombreux, organisé puis à analyser, à commenter, à critiquer ces données et à tirer des conclusions et prendre la décision qu'il faut.

Vocabulaire général

Population: C'est l'ensemble des objets ou des personnes retenues dans notre étude statistique on appelle aussi **univers**

Individus: les éléments de la population seront appelés individus.

L'échantillon: le plus souvent il est difficile d'observer toutes les données ou bien examiner l'ensemble total, on examine une partie, tout sous ensemble de la population est un échantillon.

Le caractère c'est le critère retenu pour mener l'étude statistique, chaque individu présente un caractère un caractère peut être:

qualitatif: par exp: cause d'un accident ou

quantitatif: par exp: une durée de vie.

- Un caractère quantitatif est **continu**, s'il peut prendre toutes les valeurs d'un intervalle (par ex p = une longueur).
- Il est **discontinu** ou **discret** s'il ne peut prendre que des valeurs isolées (par ex p = le nombre d'enfants dans la famille).

Remarque: Quand le caractère est continu, ou quand il est discret avec beaucoup de valeurs possibles, on effectue des regroupements en classes statistiques.

Définition d'une série statistique

Une série statistique est l'ensemble des valeurs numériques résultants de l'observation: la série statistique est du même type que le caractère qui la constitue.

ex: si le caractère est quantitatif \Rightarrow la série est quantitative

- le nombre total des valeurs numériques de la population ou l'échantillon est appelé l'**effectif total** et on note n .

L'effectif partiel (fréquence absolue).

c'est le nombre de fois où l'on a rencontré cette valeur du caractère dans la série statistique, la distribution est

noté n_i et on a: $\sum_{i=1}^t n_i = n$ / n_i = l'effectif partiel

la fréquence (relative)

si n_i est l'effectif partiel de la modalité (ou valeurs)

x_i , le rapport $= f_i = \frac{n_i}{n}$ s'appelle la fréquence.

En multipliant la fréquence f_i par 100 on obtient le pourcentage $= P_i = 100 f_i$ de la valeur x_i .

Remarque : $0 \leq n_i \leq n \iff 0 \leq \frac{n_i}{n} \leq \frac{n}{n} \Rightarrow 0 \leq f_i \leq 1$.

$$\sum_{i=1}^m n_i = n, \quad \sum_{i=1}^m f_i = 1, \quad \sum_{i=1}^m P_i = 100.$$

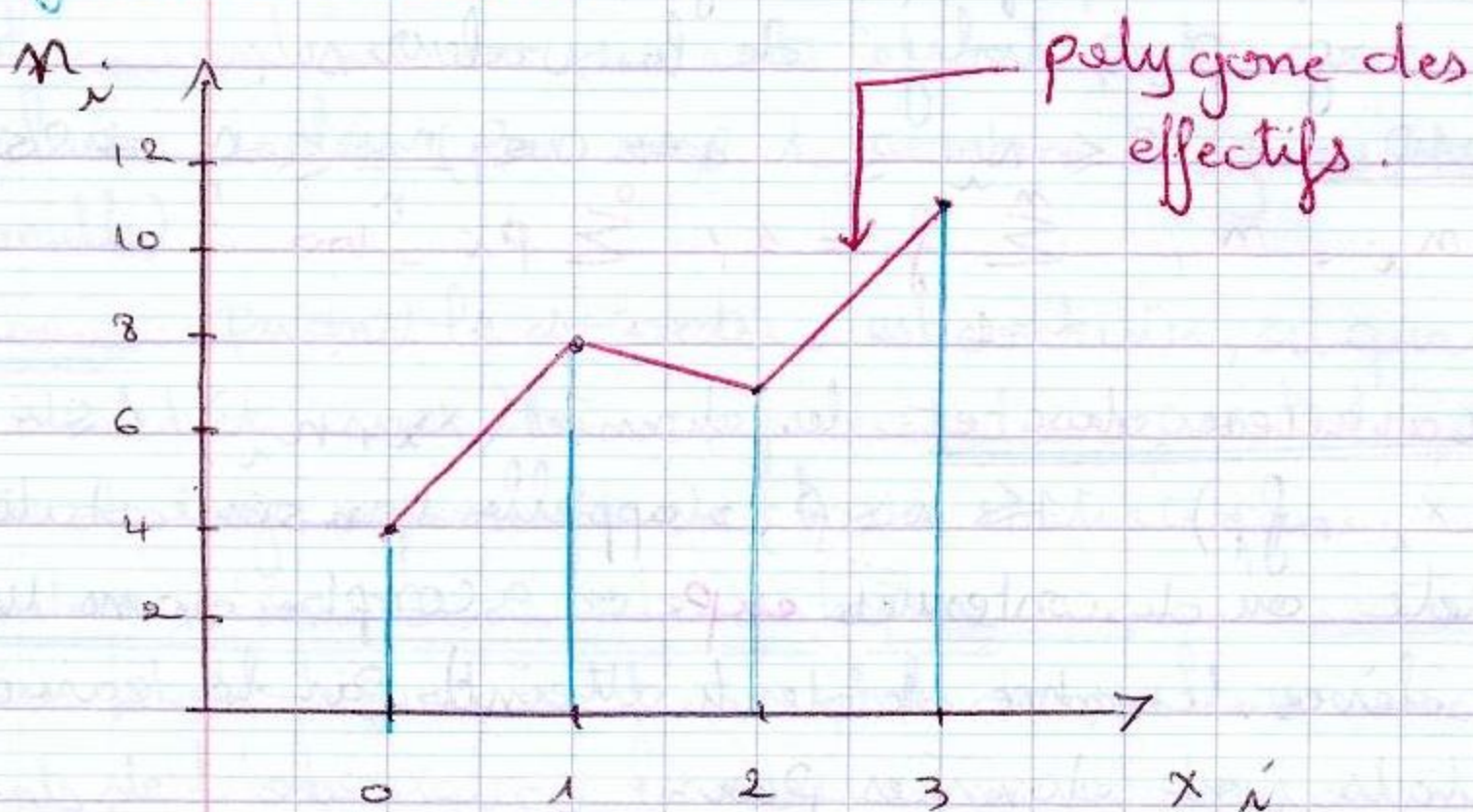
1) Un caractère discret = les données $(x_i, n_i) / 1 \leq i \leq p$ ou $(x_i, f_i) / 1 \leq i \leq p$, s'appelle une série statistique discrète ou discontinues *exp.* on a compté dans une classe de 30 élèves, le nombre de dents atteintes par la carie, les résultats sont données par :

1, 3, 3, 0, 1, 2, 3, 3, 0, 2, 2, 3, 3, 2, 3, 2, 1, 0, 2, 1, 1, 2, 3, 1, 3, 3, 3, 0, 1. Nous interprétons les résultats suivants en disant que :

" 2^{ème} " " " " 3 " " " " "
 " 3^{ème} " " " " 3 " " " "

x_i	n_i	f_i	P_i
0	4	0,1333	% 13,33
1	8	0,2666	% 26,66
2	7	0,2333	% 23,33
3	11	0,3666	% 36,66
Σ	30	1	100

La représentation graphique de cette série donnée un diagramme en bâtons



2) Un caractère continu

Si les données de la statistique ont été groupées en classes $[a_1, a_2[\dots [a_p, a_{p+1}[$ d'effectif n_1, \dots, n_p la famille $([a_i, a_{i+1}[, n_i) \ 1 \leq i \leq p$ s'appelle une série statistique groupée ou continue.

Definition pour la classe est $[a_i, a_{i+1}[$:

- Les n^{bres} a_i, a_{i+1} s'appellent les **limites**, ou **les bornes** de la classe.
- la demi somme $\frac{a_i + a_{i+1}}{2} = C_i$ s'appelle le **centre** de la classe.
- la différence $a_{i+1} - a_i = e_i$ s'appelle **endue** ou

L'amplitude de la classe.

- L'effectif de la classe est la somme des effectifs partiels des valeurs appartenant à la classe.

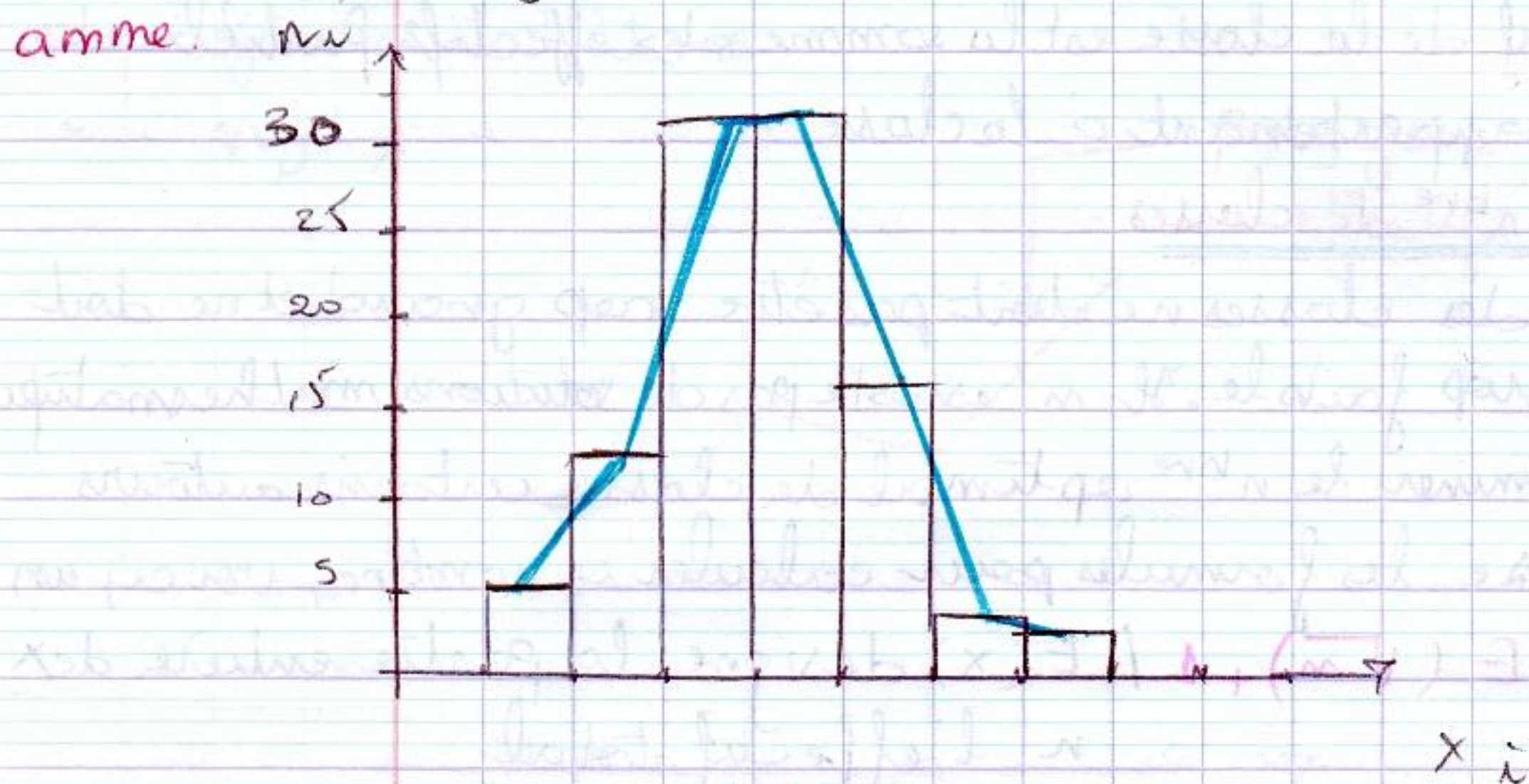
Choix du nombre de classes :

Le nombre des classes ne doit pas être trop grand et ne doit pas être trop faible. Il n'existe pas de solutions mathématique pour déterminer le nombre optimal de classes certains auteurs ont proposé les formules pour calculer ce nombre, voici, un exp. $C = E(\sqrt{n}) + 1$ / $E(x)$ désigne la partie entière de x .
 n l'effectif total

exp : On a étudié le poids dans un groupe de 100 adultes normaux de sexe féminin.

classe	C_i	n_i	f_i	P_i
39,5 - 44,5	4 1	5	0,05	5
44,5 - 49,5	4 7	12	0,12	12
49,5 - 54,5	5 2	31	0,31	31
54,5 - 59,5	5 7	31	0,31	31
59,5 - 64,5	6 2	16	0,16	16
64,5 - 69,5	6 7	3	0,03	3
69,5 - 74,5	7 2	2	0,02	2
Σ	/	100	1	100%

La représentation graphique de la série se fait par **histogramme**.



3/ Un caractère qualitatif :

lorsqu'il s'agit d'un caractère qualitatif tel que : le sexe, la couleur des yeux ou des cheveux, on classe les résultats en catégories ou modalités et on indique le nombre n_i de cas qui revient dans chaque catégorie.

Pour la représentation graphique on peut réaliser une série qualitative on plusieurs méthodes :

Un diagramme à secteurs.

Un diagramme en branches. (**Diagramme d'orgue**)

* Effectif et Fréquences cumulées - croissantes :

Soit $(x_i, n_i) \quad 1 \leq i \leq p$ une série statistique

discrète prenant les valeurs x_1, x_2, \dots, x_p classées par ordre croissant.

• On appelle effectif cumulé croissant de la valeur x_i la somme des effectifs partiels des valeurs

$$N_i = n_1 + n_2 + \dots + n_i = \sum_{k=1}^i n_k = n_i \rightarrow$$

ainsi l'effectif cumulé croissant d'une variable x_i est le nombre d'individus de la population sur lesquels la série X prend une valeur inférieure ou égale x_i ($\leq x_i$).

On appelle fréquence cumulée croissante de la valeur x_i la somme des fréquences des valeurs x_1, x_2, \dots, x_i .

$$F_i = f_1 + f_2 + \dots + f_i = \sum_{k=1}^i f_k = f_i \rightarrow$$

Effectif et fréquences cumulés décroissants

On appelle effectif cumulé décroissant de la valeur x_i le nombre

$$n_i \rightarrow = n - [n_1 + n_2 + \dots + n_{i-1}] = n - \sum_{k=1}^{i-1} n_k$$

ainsi l'effectif cumulé décroissant d'une valeur x_i est le nombre d'individus de la population

sur les quels la série prend une valeur supérieur ou égale x_i , ($\geq x_i$).
 * On appelle fréquence cumulée décroissante de la valeur x_i le nombre

$$f_{i \downarrow} = 1 - [f_1 + f_2 + \dots + f_i]$$

$$= 1 - \sum_{k=1}^{i-1} f_k$$

Exemple --

le nombre d'intervention par jour des pompiers dans une caserne est distribué comme indiqué le tableau (pendant une année)

x_i : le nombre de sortie par jour.

x_i	m_i	f_i	m_i^{\uparrow}	f_i^{\uparrow}	$m_{i \downarrow}$	$f_{i \downarrow}$
0	84	0,230	84	0,230	365	1
1	105	0,287	189	0,517	281	0,77
2	72	0,197	261	0,714	176	0,483
3	59	0,161	320	0,875	104	0,286
4	28	0,076	348	0,951	45	0,125
5	15	0,041	363	0,992	17	0,049
6	2	0,005	365	1	2	0,005
N	365	1				

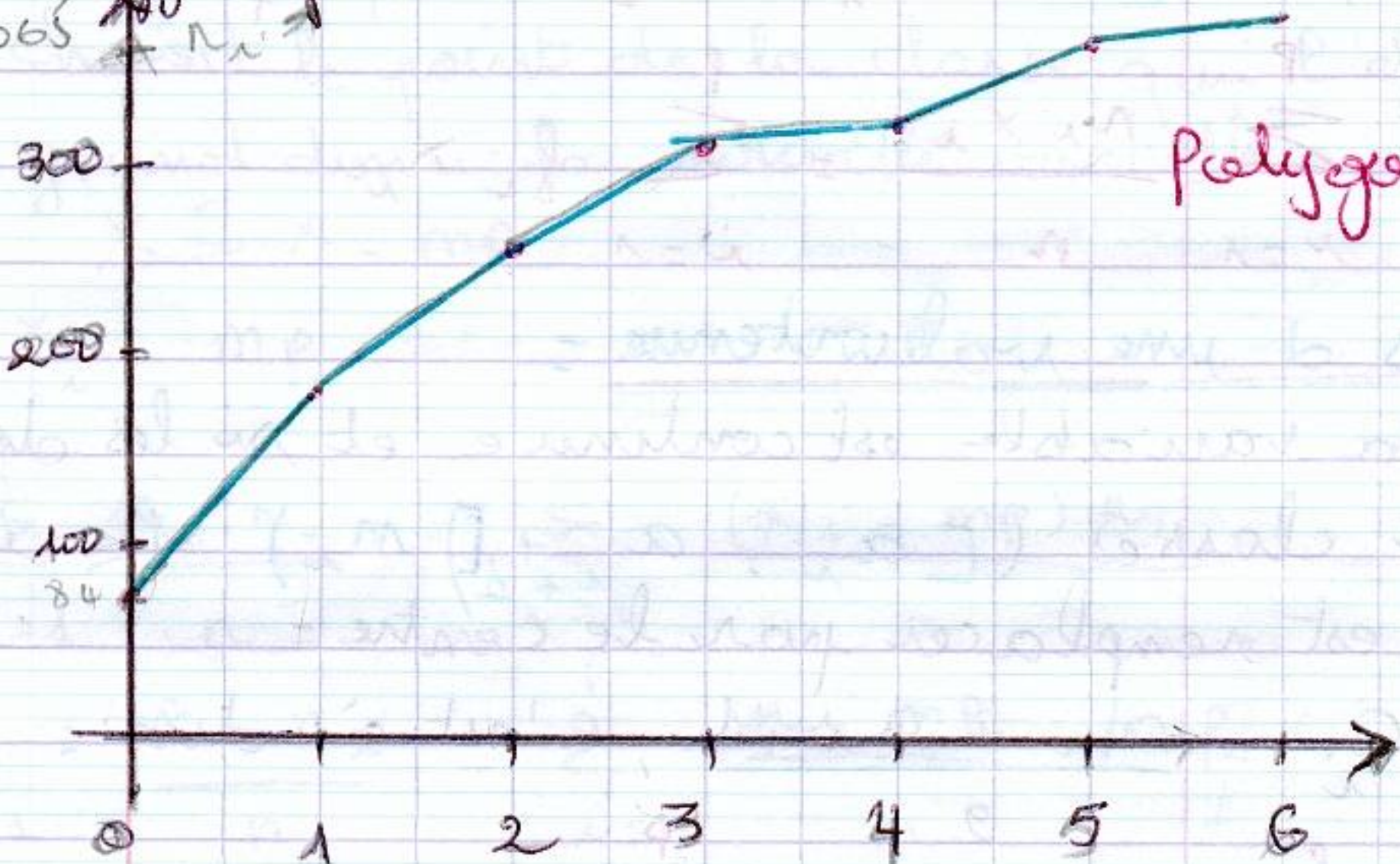
La représentation graphique de la série (x_i, n_i)

$1 \leq i \leq p$.

(x_i, n_i) $1 \leq i \leq p$ donne un diagramme appelé courbe des effectifs cumulés croissante (décroissante)

esp

365



Polygone

• Les paramètres de position :

1/ la moyenne arithmétique

la moyenne arithmétique d'une série de valeurs d'une variable statistique est égale à la somme de ces valeurs divisée par leur nombre :

• Cas d'une N discrète :

Soit (X_i, n_i) $1 \leq i \leq p$ ou (x_i, f_i) $1 \leq i \leq p$ une série statistique discrète. On appelle moyenne

arithmétique de cette série le nombre.

$$\bar{X} = \frac{m_1 x_1 + m_2 x_2 + m_3 x_3 + \dots + m_p x_p}{m_1 + m_2 + \dots + m_p}$$

$$= \sum_{i=1}^p \frac{m_i x_i}{n} = \sum_{i=1}^p f_i x_i$$

* Cas d'une v. continue =

si la variable est continue et si les données sont classées $([a_i, a_{i+1}[, m_i) \quad 1 \leq i \leq p$.
 x_i est remplacé par le centre =

$$C_i = \frac{a_i + a_{i+1}}{2}, \text{ c'est à dire :}$$

$$\bar{X} = \sum_{i=1}^p \frac{m_i \cdot C_i}{n} = \sum_{i=1}^p f_i C_i$$

Remarque: Si on a pas de pondération (répétition) répétition \bar{X} est donnée par:

$$\bar{X} = \frac{x_1 + x_2 + \dots + x_p}{1 + 1 + \dots + 1} = \sum_{i=1}^p \frac{x_i}{p}$$

La moyenne prévisive :

Lorsque les données sont nombreuses on peut simplifier le calcul de la moyenne arithmétique en adoptant le procédé dit **La moyenne prévisive mp**. On choisit le point de la classe qui l'effectif de plus grand donc la nouvelle variable.

$$X' = X - mp$$

$$X'_i = X_i - mp$$

$$X' = \sum_{i=1}^p \frac{X_i n_i}{n} = \sum_{i=1}^p \frac{(X_i - mp) n_i}{n}$$

$$= \sum_{i=1}^p \frac{X_i n_i}{n} - \sum_{i=1}^p \frac{n_i mp}{n} \quad | \quad mp = c$$

$$= \sum_{i=1}^p \frac{X_i n_i}{n} - mp \sum_{i=1}^p \frac{n_i}{n}$$

$$\bar{X}' = \bar{X} - mp \quad \Rightarrow \quad \bar{X} = \bar{X}' + mp$$

exp: Revenons à la distribution de poids. / $mp = 52$

Classes	C_i	n_i	$n_i C_i$	$c_i = C_i - mp$	$C_i' n_i$
] ,]	42	5	210	-10	-50
] ,]	47	12	564	-5	-60
] ,]	52	31	1612	0	0
] ,]	57	31	1767	5	155
] ,]	62	16	992	10	160
] ,]	67	3	201	15	45
] ,]	72	2	144	20	40
Σ		100	5490		290

$$\bar{X} = \sum \frac{n_i c_i}{n} = \frac{5490}{100} = 54,9 \text{ Kg}$$

$$\bar{X}' = \sum \frac{n_i c_i'}{n} = \frac{290}{100} = 2,9 \text{ Kg}$$

$$\bar{X}' + m_p = 2,9 + 52 = 54,9 = \bar{X}$$

2/ le mode (la valeur dominante) $M_0 =$
le mode ou la valeur modale est la valeur
que la variable prend plus fréquemment

* Variable discrète: On appelle mode d'une
série discrète la valeur x_i correspondante à
 n_i maximale.

* variable continue: Pour déterminer le mode
d'une v.c. en 1^{er} lieu on détermine la classe
modale

On appelle classe modale d'une série continue,
toute classe d'effectif maximal.
le mode est donnée par la formule

$$M_0 = l_0 + (l_2 - l_1) \left[\frac{n_i - n_{i-1}}{2n_i - (n_{i-1} + n_{i+1})} \right]$$

$$M_0 \in [l_1, l_2] \leftarrow \text{classe modale.}$$

3/ la médiane : (M_e) :

La médiane est la valeur x_i de la variable statistique qui correspond à l'effectif cumulé $\frac{n}{2}$ ou à la fréquence cumulée 0,5 (50%).

la médiane est la valeur x_i qui divise une série statistique en 2 sous groupes de même effectif.

Cas d'une variable discrète :

Soit X une variable discrète définie dans une population D et prenant les valeurs x_1, x_2, \dots, x_n classés par ordre croissant.

* Si n est impair : la médiane M_e est la $\left(\frac{n+1}{2}\right)^{\text{ème}}$ valeur : $M_e = X_{\left(\frac{n+1}{2}\right)}$

* Si n est pair : M est $\frac{1}{2} \left[\left(\frac{n}{2}\right)^{\text{ème}} + \left(\frac{n}{2} + 1\right)^{\text{ème}} \right]$
 $M_e = \frac{1}{2} \left[X_{\frac{n}{2}} + X_{\frac{n}{2} + 1} \right]$

ex :

On a les 2 séries suivantes

1/ 2 2 3 3 3 4 5 6 6 7 8 8 9

$$n = 13 \Rightarrow M_e = X_{\frac{13+1}{2}} = X_7$$

$$M_e = 5$$

2/ 2 2 3 3 3 (4) ↓ (5) 6 6 7 8 8.

$$n = 12 \Rightarrow M_e = \frac{1}{2} \left[X_{\frac{12}{2}} + X_{\frac{12}{2} + 1} \right] = \frac{1}{2} [X_6 + X_7]$$

$$M_e = \frac{1}{2} [4 + 5] = 4,5$$

exemple 2: si les données sont classées dans un tableau

X_i	n_i	n_i^{\uparrow}
0	4	4
1	8	12
2	7	19
3	11	30
Σ	30	

$$\frac{n}{2} = \frac{30}{2} = 15$$

$$15 \in [12, 19]$$

$$M_e = 2$$

Cas d'une variable continue

Pour une série continue ou groupée

$$M_e = l_1 + \frac{(l_2 - l_1)}{n} \left[\frac{n}{2} - F^{\uparrow} \right]$$

• l_1, l_2 : sont les bornes de la classe contenant la médiane.

• n : est l'effectif de cette classe médiane.

• F^{\uparrow} : c'est l'effectif cumulé de toutes les classes précédant la classe médiane.

excp.

* Considerons la série suivantes :

classes	n_i	$\sum_{j=1}^i n_j$ ↗
[10, 15[10	10
[15, 20[11	21
[20, 25[17	38
[25, 30[06	44
[30, 35[12	56
[35, 40[04	60
Σ	60	

$$n = 60$$

$$\Rightarrow \frac{n}{2} = 30$$

$$\frac{n}{2} = 30 \in [21, 38]$$

$$M_e \in [20, 25[$$

$$M_e = 20 + \left(\frac{25 - 20}{17} \right) [30 - 21] = 22,64$$

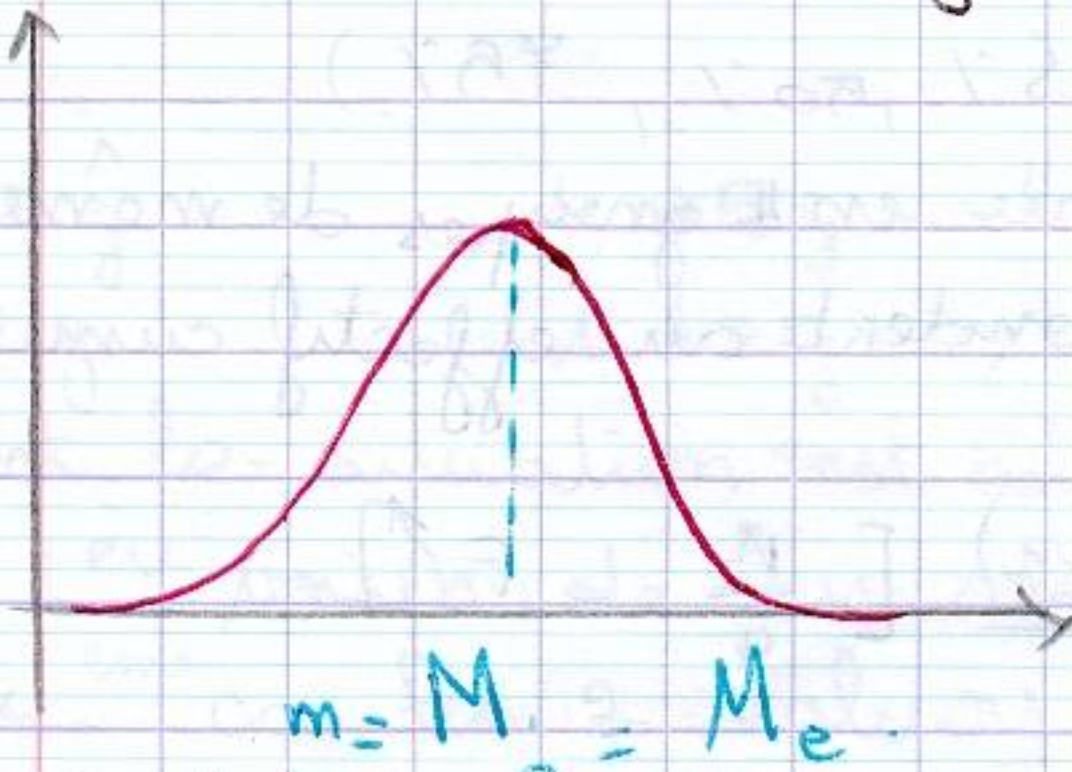
Détermination graphique de la médiane :

La détermination graphique de la médiane se fait, on trace le polygone des effectifs cumulés (des fréquences cumulées) de rechercher l'abscisse du point de ce polygone d'ordonnée égale $\frac{n}{2}$ ou $(0,5)$. Une autre méthode : on trace les 2 polygones n_i ↗, n_i ↘

• M_e est l'abscisse du point d'intersection

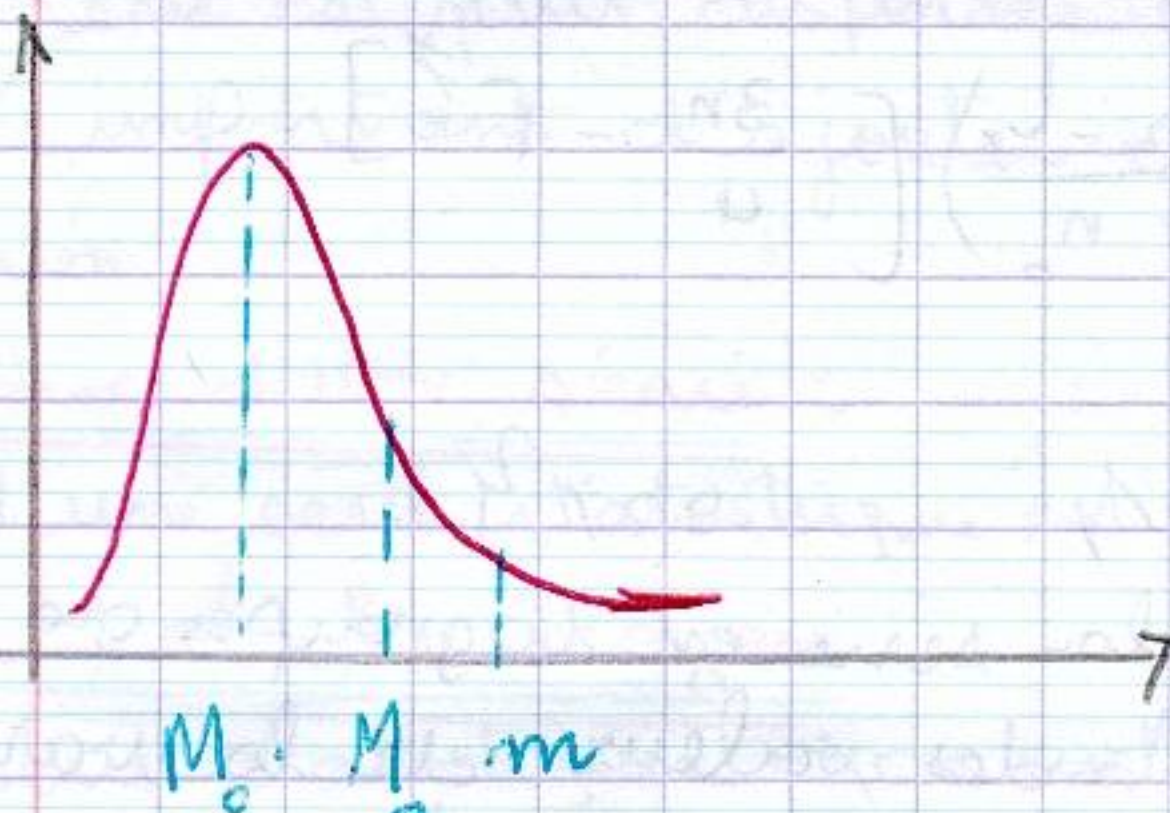
Remarque =

* Si mode = médiane = moyenne.



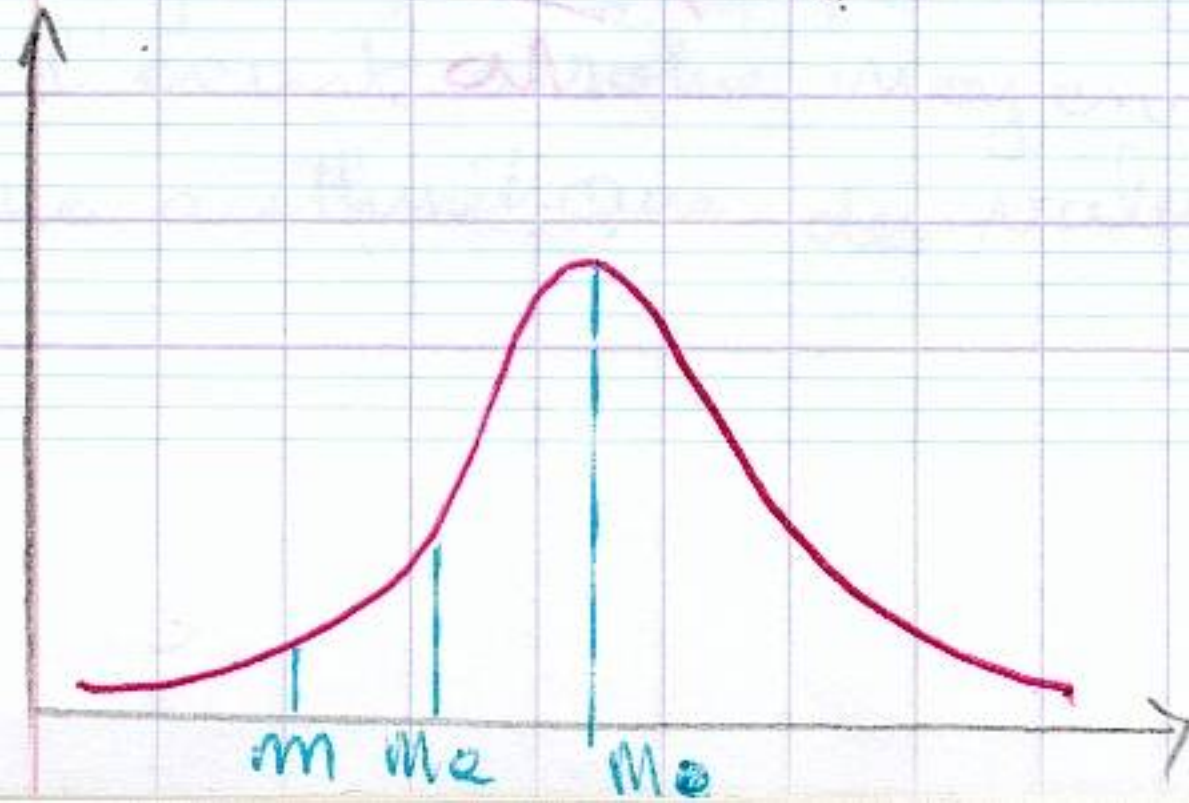
La série est dite symétrique.

* Si Mode < Médiane < Moyenne.



on dit que la série est biaisée à droite ou positivement.

* Si moyenne < médiane < Mode.



On dit que la série est biaisée à gauche ou négativement

• Les quartiles : (25%, 50%, 75%).

Ce qui divise la série en 4 groupes de même effectif. La 1^{ère} quartile correspondent au effectif cumulé $\frac{n}{4}$

• le cas continue :

$$Q_1 = l_1 + \frac{(l_2 - l_1)}{n_0} \left[\frac{n}{4} - F^{\uparrow} \right]$$

• le cas continue

• la 2^{ème} quartile est la médiane. $Q_2 = M_e$

• la 3^{ème} quartile correspondent aux $\frac{3n}{4}$

• le cas continue :

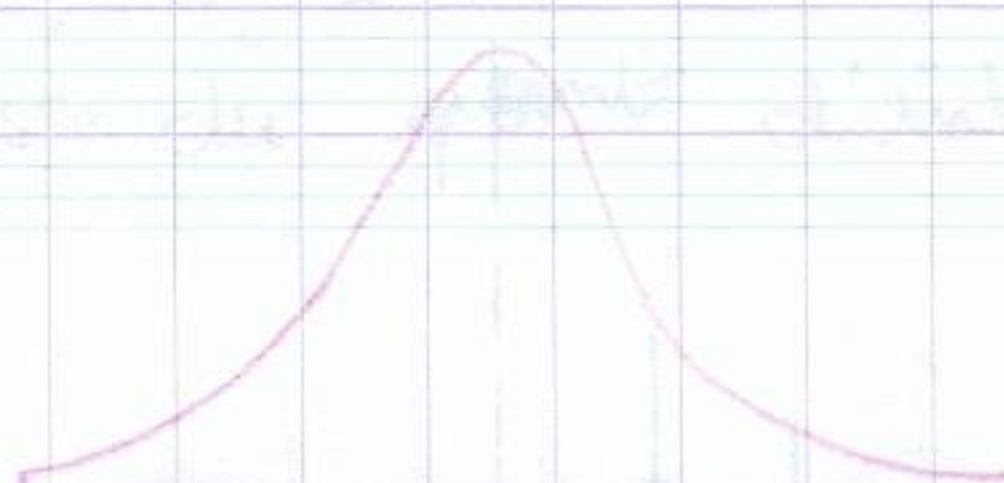
$$Q_3 = l_1 + \frac{(l_2 - l_1)}{n_0} \left[\frac{3n}{4} - F^{\uparrow} \right]$$

• les déciles :

$K = \{ 10\%, 20\%, \dots, 90\% \}$

Ce qui divise la série en 10 groupes de même effectif. Il s'agit des valeurs de la variable qui correspondent au effectifs cumules

$\frac{n}{10}, \frac{2n}{10}, \dots, \frac{9n}{10}$



2] Les paramètres de dispersion

Considérons les 2 séries suivantes qui donnent les glycémiés -

0,95 g 1,05 g $m = 1$, $M_e = 1g$

0,50 g 1,50 g $m = 1$, $M_e = 1g$

et pourtant la situation est tout à fait différente dans la 1^{ère} voisines de leur moyenne.

dans le 2^{ème} cas les 2 valeurs très éloignées de leur moyenne, c'est ce qu'en statistique on nomme

dispersion.

La 2^{ème} série est plus dispersée que la 1^{ère}

Il est donc important d'ajouter autres paramètres de dispersion

1/ L'étendue d'une série :

L'étendue d'une série statistique est : $E = X_{\max} - X_{\min}$

2/ L'écart absolu moyen :

Soit une série statistique prenant les valeurs x_1, \dots, x_p . On considère son écart à la moyenne.

$$e_1 = x_1 - \bar{x}, \quad e_2 = x_2 - \bar{x}, \quad \dots, \quad e_p = x_p - \bar{x}$$

On définit l'écart absolu moyen comme étant la moyenne arithmétique des valeurs absolues des écarts.

$$\bar{e} = \sum_{i=1}^p \frac{e_i n_i}{n} = \sum_{i=1}^p \frac{|x_i - \bar{x}| n_i}{n}$$

Dans le cas continu x_i est remplacé par C_i centre de classe $[a_i, a_{i+1}]$. $i \leq i' \leq p$.

$$\bar{e} = \sum_{i=1}^p \frac{|C_i - \bar{x}| n_i}{n}$$

3/ L'écart type et la variance

La caractéristique de dispersion la plus usuelle est en effet l'écart type

• la variance: Soit la série statistique prenant les valeurs x_1, x_2, \dots, x_p / $e_1 = x_1 - \bar{x}, \dots, e_p = x_p - \bar{x}$.
On définit la variance par la moyenne arithmétique des carrés des écarts à la moyenne.

$$V = \sum_{i=1}^p \frac{n_i e_i^2}{n} = \sum_{i=1}^p \frac{n_i (x_i - \bar{x})^2}{n}$$

L'écart type d'une série

c'est la moyenne quadratique des écarts à la moyenne autrement dit:

$$\sigma = \sqrt{V} = \sqrt{\sum_{i=1}^p \frac{n_i (x_i - \bar{x})^2}{n}}$$

pour le cas continue on utilise la même formule
remplaçant x_i par C_i centre de classe

Méthode de calcul:

$$V = \frac{\sum_{i=1}^p n_i (x_i - \bar{x})^2}{n}$$

$$= \sum_{i=1}^p \frac{1}{n} [n_i (x_i^2 - 2\bar{x}x_i + \bar{x}^2)]$$

$$= \sum_{i=1}^p \frac{1}{n} [n_i \cdot x_i^2 - 2n_i \bar{x}x_i + n_i \bar{x}^2]$$

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

$$= \frac{1}{n} \sum n_i x_i^2 - \frac{2\bar{x}}{n} \sum n_i x_i + \bar{x}^2 \sum \frac{n_i}{n}$$

$$= \sum \frac{n_i x_i^2}{n} - 2\bar{x}\bar{x} + \bar{x}^2$$

$$V(x) = \sum_{i=1}^p \frac{n_i x_i^2}{n} - \bar{x}^2$$

Signification de l'écart type :

On se fait une idée simple de la signification de l'écart type lorsqu'on compare 2 séries statistiques de même nature.

Celle qui a l'écart type le plus grand est la plus dispersée.

excp: Considérons les 2 séries:

1/ 95 97 100 103 105

2/ 50 75 100 125 150

$$\bar{X}_1 = \sum_{i=1}^5 \frac{X_{i,1}}{5} = 100, \quad \bar{X}_2 = \sum_{i=1}^5 \frac{X_{i,2}}{5} = 100$$

$$n(1) = n(2) = 5 \Rightarrow Me_1 = \frac{X_{n+1}}{2} = X_3 = 100$$

$$Me_2 = 100$$

$$V_1(X_1) = \sum_{i=1}^5 \frac{X_{i,1}^2}{n} - (\bar{X}_1)^2 = 1316$$

$$\Rightarrow \sigma_1 = \sqrt{V_1} = 3,68$$

$$V_2(X_2) = \sum_{i=1}^5 \frac{(X_{i,2})^2}{n} - (\bar{X}_2)^2 = 12,50$$

$$\sigma_2 = \sqrt{V_2} = 3,53$$

$\sigma_2 > \sigma_1$ donc la 2^{ème} série est plus dispersée que la 1^{ère}.

L'interval interquartile (L'ecart interquartile) :

La dispersion d'une serie statistique est par fois mesuree par l'interval interquartile ou la semi interquartile.

* L'interval interquartile est $[Q_1, Q_3]$ est l'interquartile est la difference entre le 1^{er} quartile et le 3^{eme} quartile $Q_2 = Q_3 - Q_1$.

* Le semi interquartile : $Q_2 = \frac{Q_3 - Q_1}{2}$

* Le coefficient interquartile est $\frac{Q_3 - Q_1}{Q_2} = \frac{Q_2}{M_e}$.

On definit de la meme maniere l'interval interdecile $[D_9, D_1]$ de l'interdecile est : $D_2 = D_9 - D_1$

le coefficient de variation

Le coefficient de variation permet de comparer les dispersion de divers series statistique, le coefficient de variation est le rapport entre la moyenne et l'ecart type. $C_V = \frac{S}{\bar{X}}$