

Statistiques descriptives

Lasse9 statistiques descriptives | GHAZOUAN Oumaima



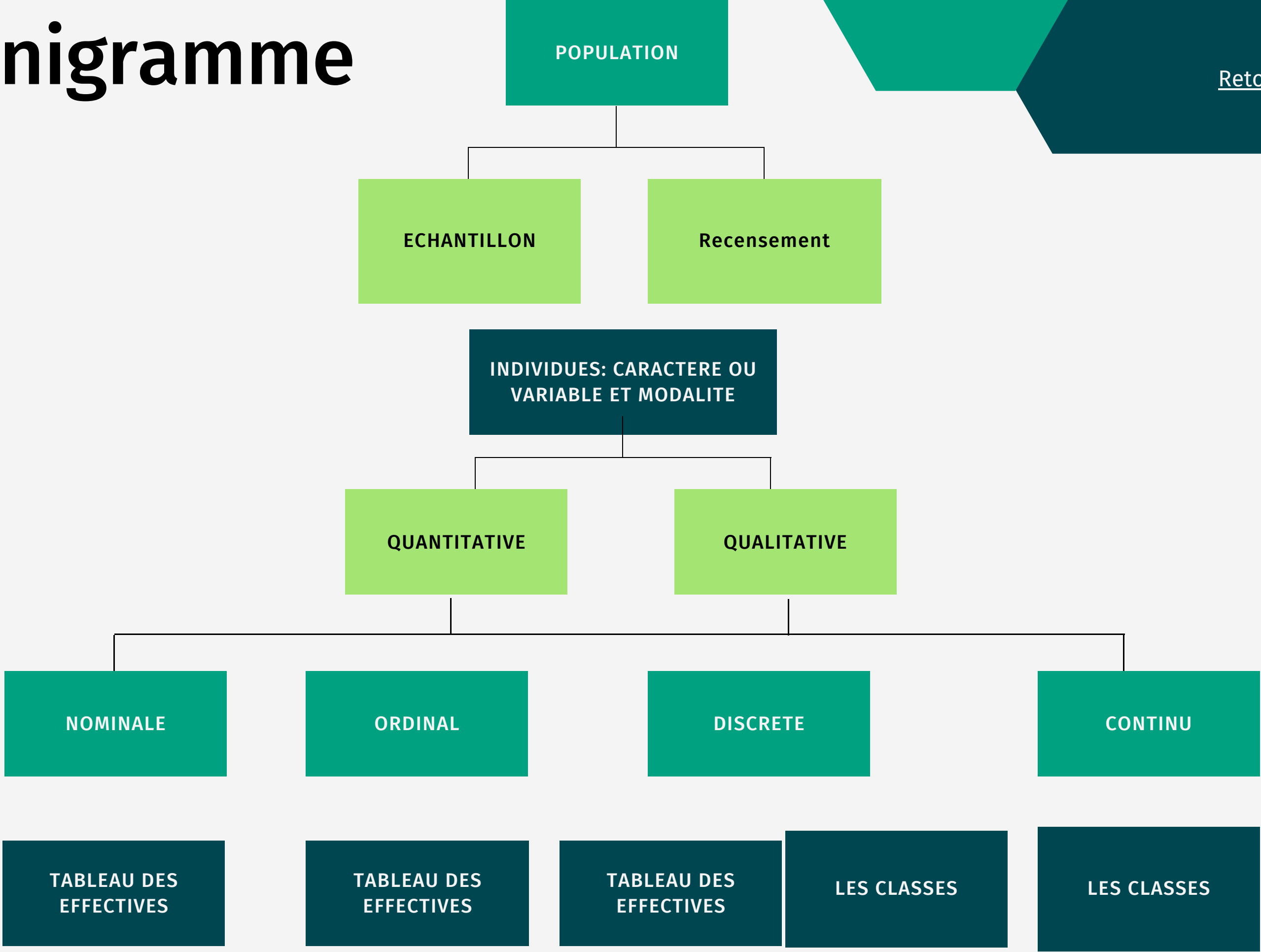
Programme

- Séries simples
- séries doubles
- séries chronologiques

Séries simples

The background features a dark teal color. On the left side, there are two overlapping hexagonal shapes. The top hexagon is a medium teal color, and the bottom hexagon is a light lime green color. The text 'Séries simples' is written in white, sans-serif font, positioned over the medium teal hexagon.

Organigramme



[Retourner à la page Programme](#)

Tableau récapitulatif Pour une variable quantitative continue:

Classes	n_i	a_i	c_i	f_i	f_{icc}	f_{icd}	n_{icc}	n_{icd}	d_i	n_{ic}
[0; 5[
[5;15[
[15;30[

- Xi : caractère
 Ni : effectif
 Nic : effectif cumulé
 Nicc: effectif cumulé croissant
 Nicd : effectif cumulé décroissant
 Fr : fréquence n_i/N
 N=somme des n_i
 Pi : pourcentage $(n/N)*100$
- Ci : centre de classe $(a+b)/2$
 Ai : amplitude $b-a$ ou x_i-x_{i-1}
 di : densité n_i / a_i
 Nic= $(N_i / a_i)*ppcm(a_i)$

Graphiques

- Caractère qualitatif

Diagramme en tuyaux d'orgues

Diagramme circulaire / demi-circulaire

- Caractère quantitatif discret

Diagramme en bâtons

Courbe cumulative des fréquences

- Caractère quantitatif continu

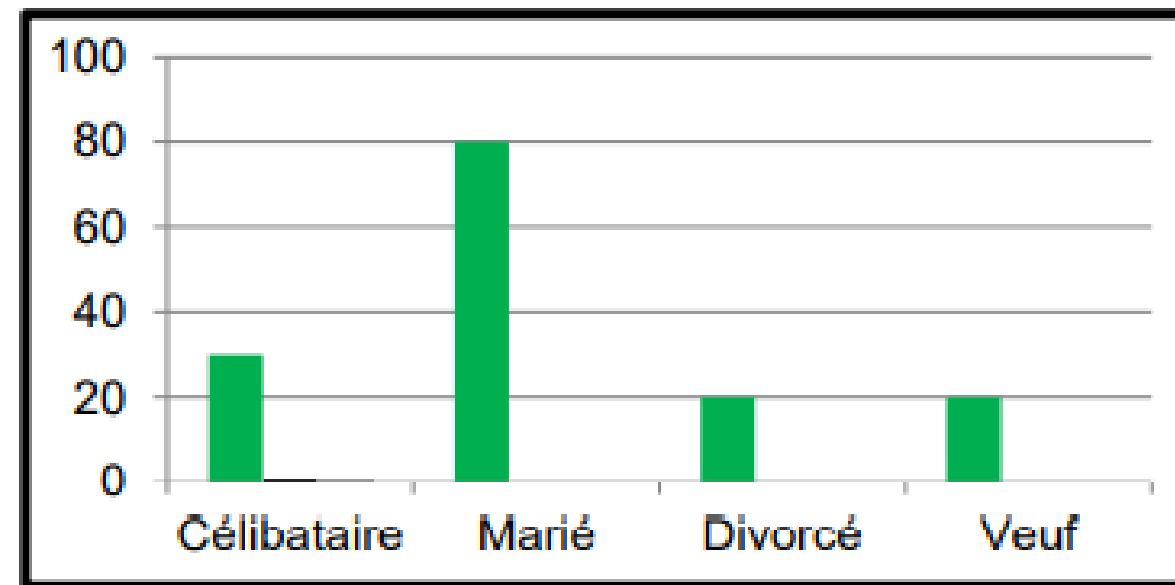
Histogramme

Polygone de fréquences

Courbe cumulative de fréquences

Pour une variable qualitative

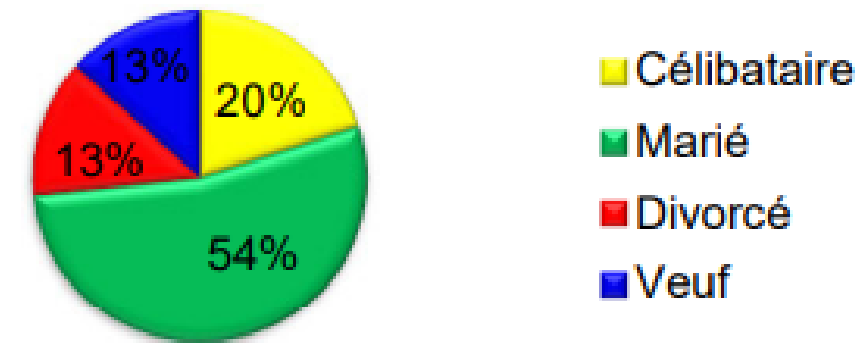
Diagramme des tuyaux d'orgue



La longueur des tuyaux = n_i

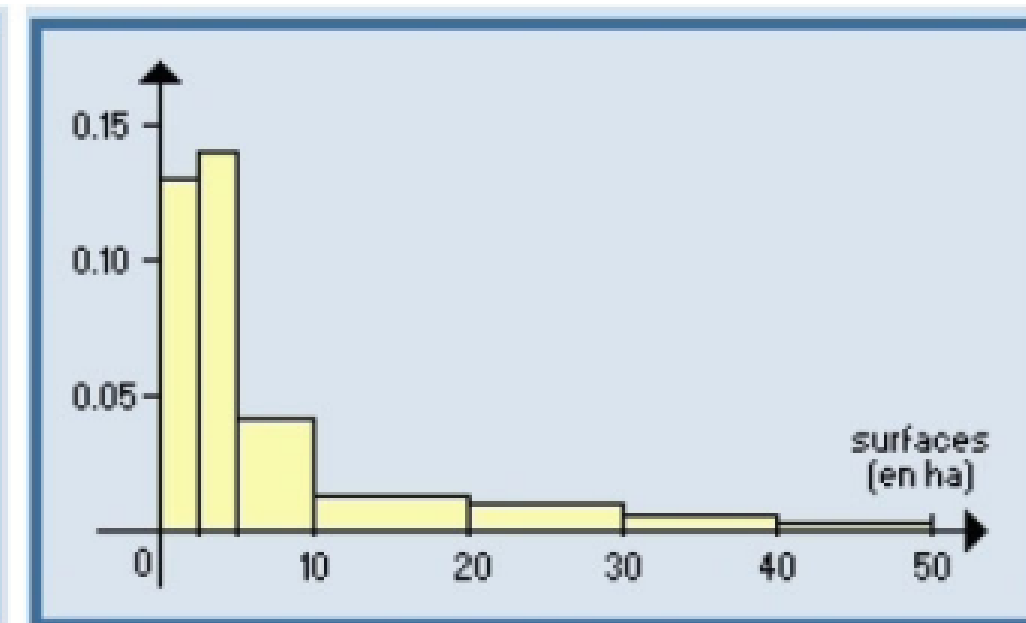
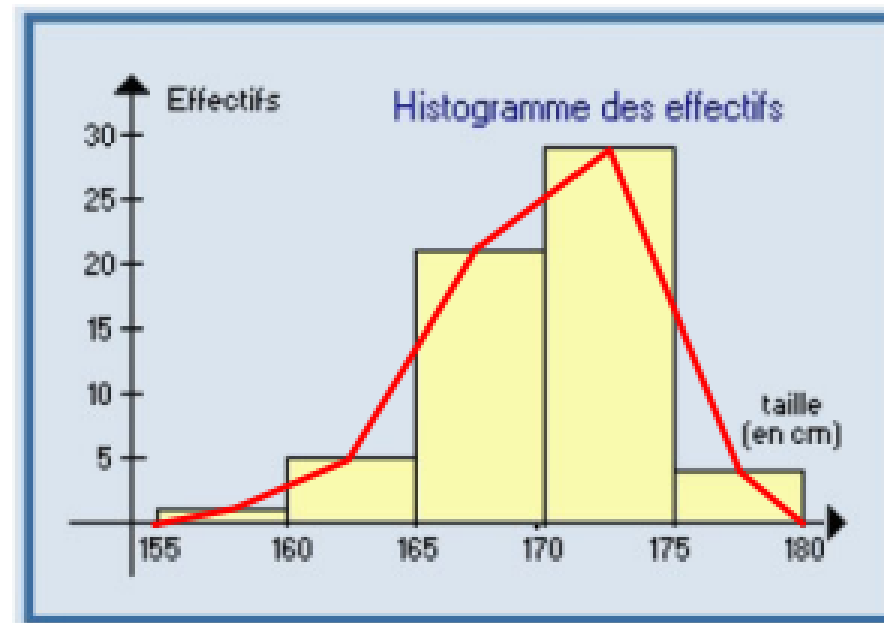
Le graphique à secteurs

Effectifs



$$\alpha_i = \frac{n_i}{n} \times 360 = f_i \times 360$$

Pour une variable quantitative



Fonction de répartition

La fonction de répartition est donnée par:

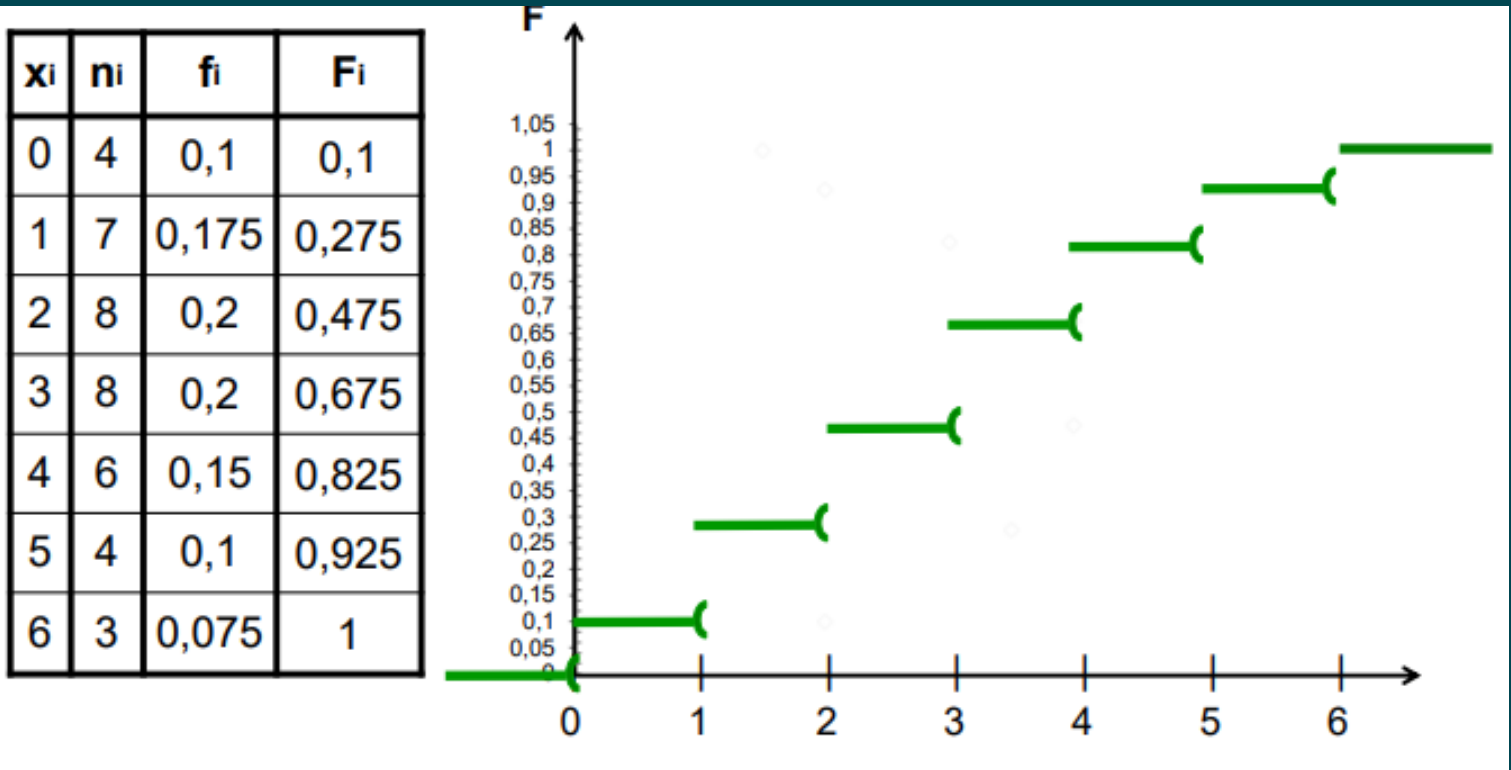
F_i = fréquence de nombre de famille qui ont moins de x_i enfants n_{icc}/n

$$F(x) = \begin{cases} 0 & \text{si } x < x_1 \\ F_i & \text{si } x_i \leq x < x_{i+1} \\ 1 & \text{si } x \geq x_k \end{cases}$$

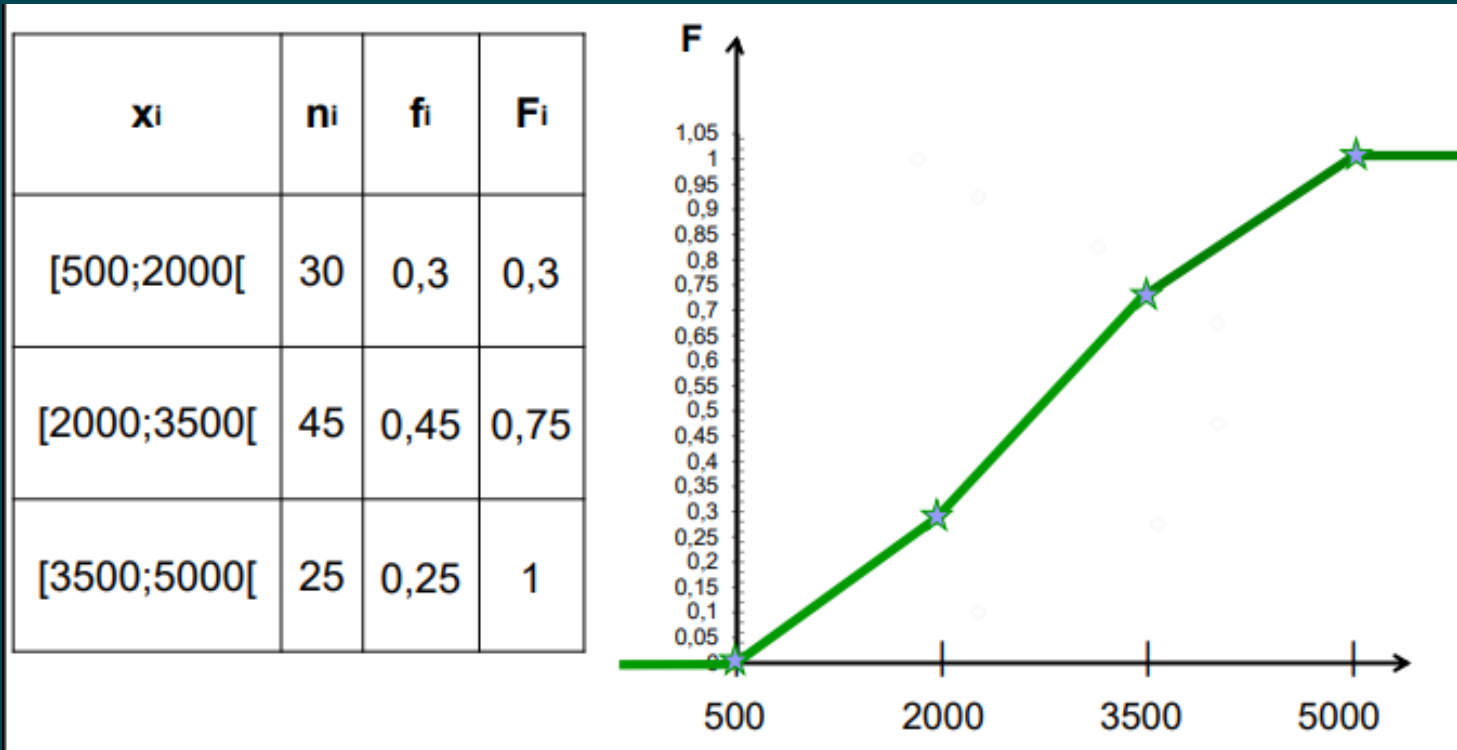
Exemple: le nombre d'enfants par famille,

x_i	n_i	N_i	F_i
0	4	4	0,1
1	7	11	0,275
2	8	19	0,475
3	8	27	0,675
4	6	33	0,825
5	4	37	0,925
6	3	40	1

$F(0)$ = fréquence de nombre de famille moins de 0 enfants = 0/40
 $F(1)$ = fréquence de nombre de famille moins de 1 enfants = 4/40
 $F(2)$ = fréquence de nombre de famille moins de 2 enfants = 11/40
 $F(3)$ = fréquence de nombre de famille moins de 3 enfants = 19/40
 $F(4)$ = fréquence de nombre de famille moins de 3 enfants = 27/40
 $F(5)$ = fréquence de nombre de famille moins de 3 enfants = 33/40
 $F(6)$ = fréquence de nombre de famille moins de 6 enfants = 37/40
 $F(7)$ = fréquence de nombre de famille moins de 7 enfants = 40/40



Représentation graphique de la fonction de répartition d'une variable discrète



Représentation graphique de la fonction de répartition d'une variable classée ou qualitative continue

Mode

Le mode de cette série statistique est la modalité de la variable correspondant à l'effectif le plus élevé

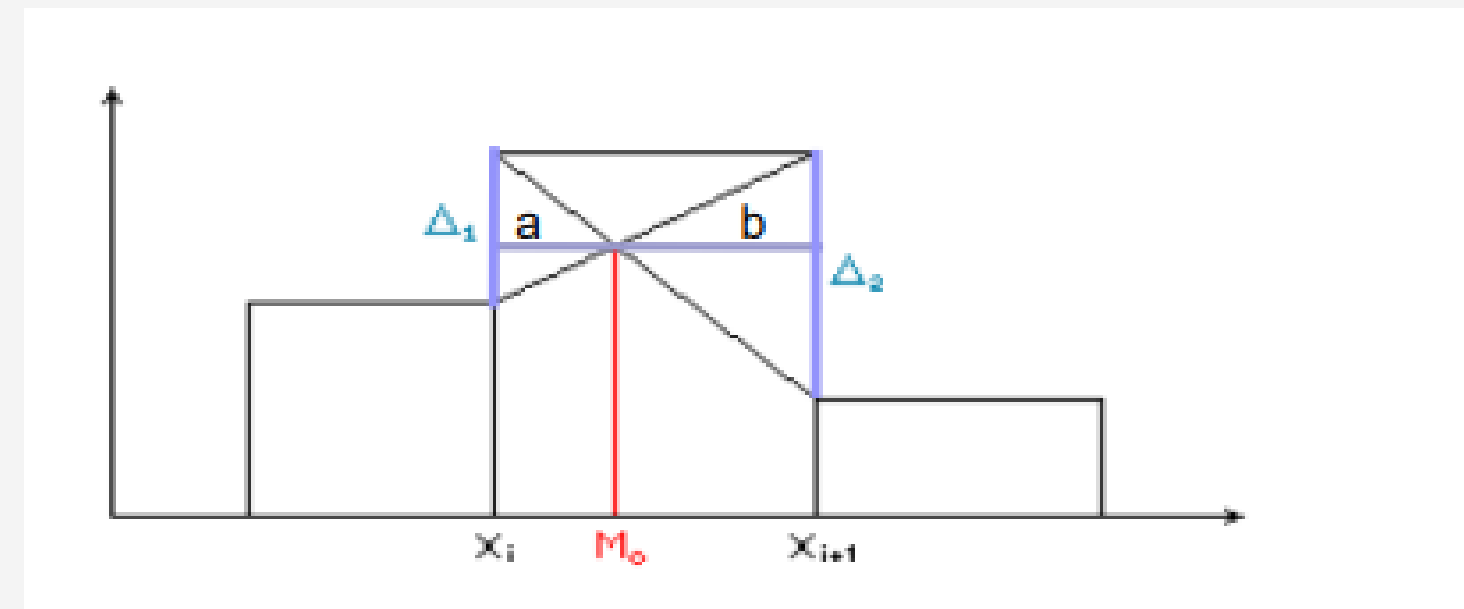
Les paramètres
de position

Pour les variables quantitatives classées

on parle d'abord de la classe modale:

- Si les classes sont d'égales amplitudes, la classe modale sera la classe où l'effectif est le plus élevé.
- Si les classes sont d'inégales amplitudes, la classe modale sera la classe où:

La densité ou La densité de fréquence est la plus élevée



$$M_o = \frac{x_i \Delta_2 + x_{i+1} \Delta_1}{\Delta_1 + \Delta_2} = x_i + \frac{\Delta_1 (x_{i+1} - x_i)}{\Delta_1 + \Delta_2}$$

Médiane

pour les variables quantitatives

Les paramètres
de position

Cas d'une variable non classée

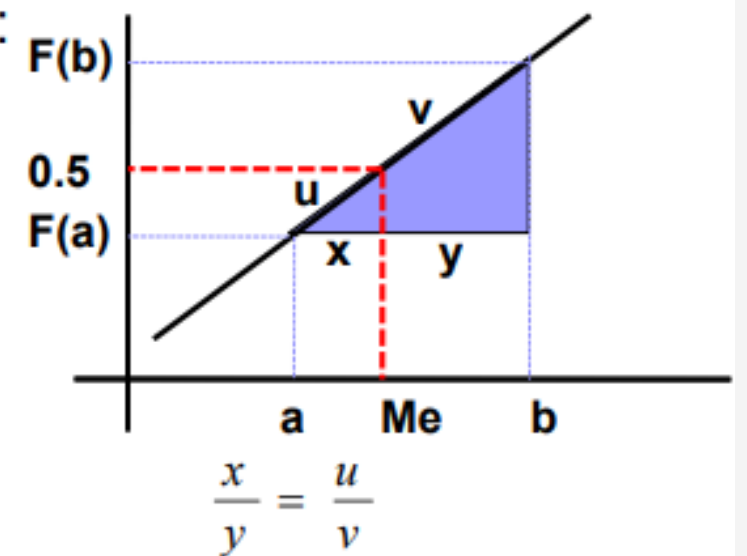
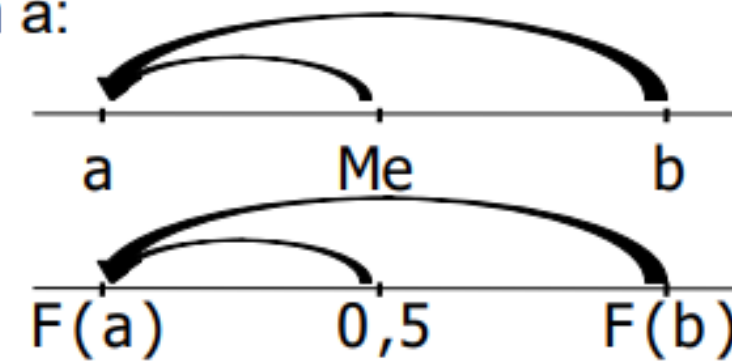
$x_1 < x_2 < \dots < x_n$

- Si n est impair et égal $2p+1$ la médiane sera: x_{p+1} .
- Si n est pair et égal $2p$, la médiane est: $x_p + x_{p+1} / 2$

Cas d'une variable classée

Soient a et b les bornes inférieures et supérieures de la classe contenant la médiane, $F(a)$ et $F(b)$ les valeurs des fréquences cumulées croissantes en a et b , alors:

On a:



$$\frac{Me - a}{b - a} = \frac{0,5 - F(a)}{F(b) - F(a)} \implies Me = a + (b - a) \times \frac{0,5 - F(a)}{F(b) - F(a)}$$

Moyenne

pour les variables classées on utilise ni

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i \quad \text{ou} \quad \bar{X} = \frac{1}{n} \sum_{i=1}^k n_i x_i = \sum_{i=1}^k \frac{n_i}{n} x_i = \sum_{i=1}^k f_i x_i$$

pour les variables non classées on utilise ci

$$\bar{X} = \frac{1}{n} \sum_{i=1}^k n_i c_i = \sum_{i=1}^k \frac{n_i}{n} c_i = \sum_{i=1}^k f_i c_i$$

Les paramètres
de position

Arithmétique	Géométrique	Harmonique	Quadratique
$\bar{X} = \sum_{i=1}^k f_i x_i$	$\bar{X}_g = \prod_{i=1}^k x_i^{f_i}$	$\bar{X}_h = \frac{1}{\sum_{i=1}^k f_i \frac{1}{x_i}}$	$\bar{X}_q = \sqrt{\sum_{i=1}^k f_i x_i^2}$

Comparaison des moyenne

$$\bar{X}_h \leq \bar{X}_g \leq \bar{X} \leq \bar{X}_q$$

La variance V

la dispersion des valeurs autour de la moyenne

L'écart-type σ

mesure la dispersion des valeurs d'un échantillon par rapport à la moyenne mais son interprétation dépend de l'échelle de la variable

Les paramètres
de dispersion

$$V(X) = \frac{1}{n} \sum_{i=1}^k n_i (x_i - \bar{x})^2 = \sum_{i=1}^k f_i (x_i - \bar{x})^2 = \frac{1}{n} \left(\sum_{i=1}^k n_i x_i^2 \right) - \bar{x}^2$$
$$\sigma(X) = \sqrt{V(X)}$$

Etendue $V_{\max} - V_{\min}$

Quartile $[Q1 ; Q3]$

des caractéristiques de position partageant la série statistique

Intervalles interquartiles $Q3 - Q1$

Décile $[D1 ; D9]$

Ecart interdécile $D9 - D1$

Exercice

Les paramètres
de dispersion

* Cas d'une Caractère Continue :

Ex : 97 à 100 DH

Donc : $q_1 = 25\%$ $q_3 = 75\%$
 $q_2 = 50\%$ $q_4 = 100\%$

x_i	m_i	m_{icc}
10-30	4	4
30-50	16	20
50-70	20	40
70-80	17	57
N	57	

* Rang = $\frac{N}{2}$

* q_1 : Rang = $\frac{n \times 25}{100} = \frac{57 \times 25}{100} = 14,25$

la classe de $q_1 = [30-50[$

$$q_1 = 30 + (50 - 30) \times \frac{14,25 - 4}{20 - 4}$$

$$q_1 = 32,03$$

* q_2 : Rang = $\frac{n \times 50}{100} = \frac{57 \times 50}{100} = 28,5$

la classe de $q_2 = [50-70[$

$$q_2 = 50 + (70 - 50) \times \frac{28,5 - 20}{40 - 20}$$

$$q_2 = 29,75$$

* Decile : $D_1, D_2, D_3, \dots, D_{10}$

$$\% = 10\%$$

$$D_1 = \frac{N \times 10}{100}$$

$$D_2 = \frac{N \times 20}{100}$$

Boîte à moustaches

Le diagramme en boîte à moustaches ou box-plot permet de représenter schématiquement les principales caractéristiques d'une distribution en utilisant les quartiles.

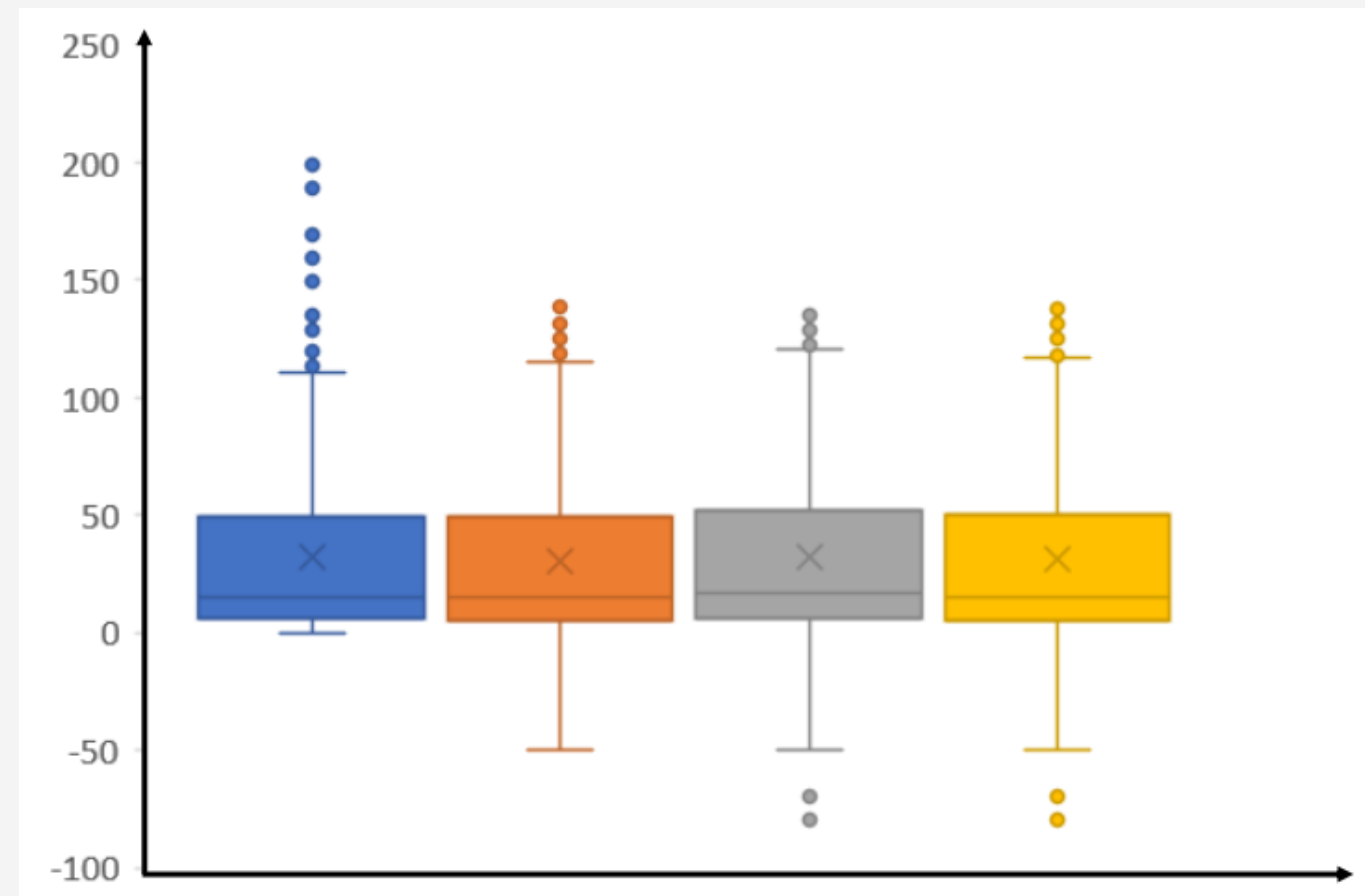
La partie centrale de la distribution est représentée par une boîte de largeur arbitraire et de

longueur la distance interquartile, la médiane est tracée à l'intérieur.

La boîte rectangle est complétée par des moustaches correspondant aux valeurs suivantes:

- Valeur supérieure : $\text{Min}(\text{la plus grande la valeur}; Q3 + 1,5(Q3 - Q1))$
- Valeur inférieure : $\text{Max}(\text{la plus petite valeur}; Q1 - 1,5(Q3 - Q1))$

Les valeurs extérieures « aux moustaches » sont représentées par des étoiles et peuvent être considérées comme aberrantes.



Coefficient de variation

comparer la distribution autour de la moyenne de deux variables statistiques de natures différentes:
Plus la valeur du coefficient de variation est élevée, plus la dispersion autour de la moyenne est grande.

$$C_v = \frac{\sigma_x}{\bar{X}}$$

Le moment d'ordre r et moment centré d'ordre r

Les paramètres
de dispersion

$$m_r = \frac{1}{n} \sum_{i=1}^k n_i x_i^r = \sum_{i=1}^k f_i x_i^r$$

$$\mu_r = \frac{1}{n} \sum_{i=1}^k n_i (x_i - \bar{x})^r = \sum_{i=1}^k f_i (x_i - \bar{x})^r$$

Le moment d'ordre r et moment centré d'ordre r

Les paramètres
de dispersion

$$m_r = \frac{1}{n} \sum_{i=1}^k n_i x_i^r = \sum_{i=1}^k f_i x_i^r$$

$$\mu_r = \frac{1}{n} \sum_{i=1}^k n_i (x_i - \bar{x})^r = \sum_{i=1}^k f_i (x_i - \bar{x})^r$$

Coefficient d'asymétrie de Fisher et Coefficient d'aplatissement de Fisher

Paramètres de forme

'une variable suit une loi normale ou de Gauss si :
Le coefficient d'asymétrie doit être inférieur à |1|
Le coefficient d'aplatissement ou encore de concentration doit être inférieur à |1,5|

Il est défini par:

$$\gamma_1 = \frac{\mu_3}{\sigma^3}$$

Si $\gamma_1 = 0$, la distribution est **symétrique** autour de la moyenne.

Si $\gamma_1 < 0$, la distribution est plus étalée **vers la gauche**.

Si $\gamma_1 > 0$, la distribution est plus étalée **vers la droite**.

Il est défini par:

$$\gamma_2 = \frac{\mu_4}{\sigma^4} - 3$$

Si $\gamma_2 = 0$, L'aplatissement est le même que celui de la **loi Normale** (de Gauss).

Si $\gamma_2 < 0$, la concentration des valeurs autour de la moyenne est faible: la distribution est **plus aplatie** que la loi Normale.

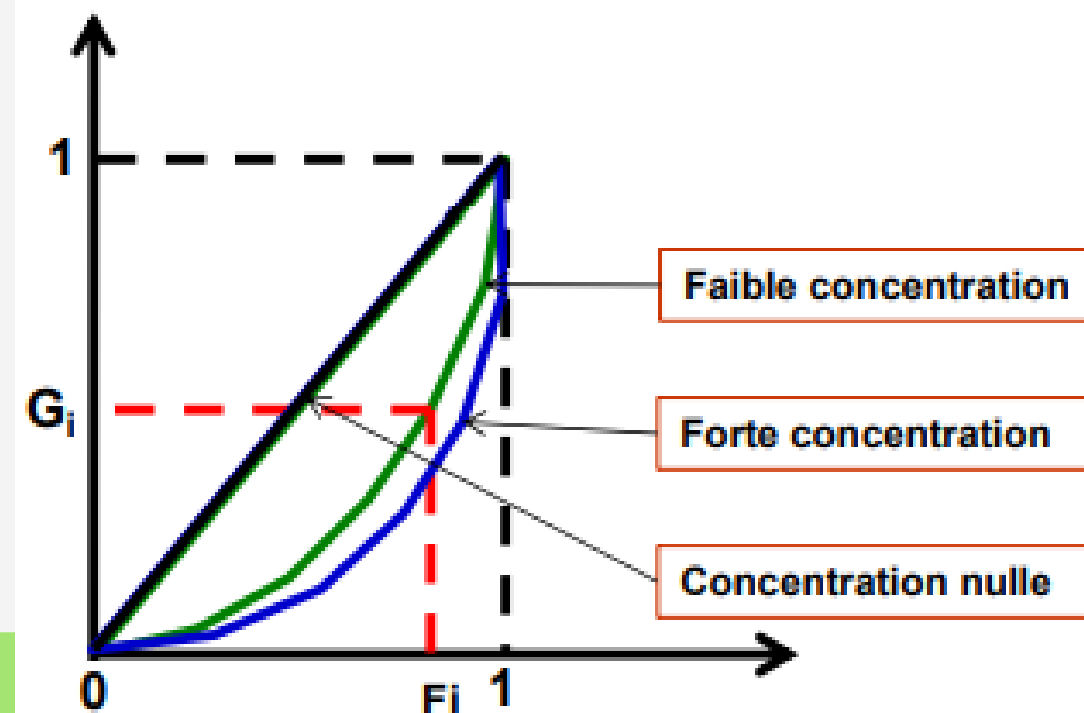
Si $\gamma_2 > 0$, la concentration des valeurs autour de la moyenne est forte: la distribution est **moins aplatie** que la loi Normale.

Indice de concentration de Gini

Paramètre de concentration

'une variable suit une loi normale ou de Gauss si :
Le coefficient d'asymétrie doit être inférieur à |1|
Le coefficient d'aplatissement ou encore de concentration doit être inférieur à |1,5|

- La $i^{\text{ème}}$ classe $[x_{i-1}, x_i[$ a, pour centre, c_i et, pour effectif, n_i .
- $s_i = n_i c_i$ la masse de caractère X dans la classe $[x_{i-1}, x_i[$.
 - $S = \sum_{i=1}^k s_i$ la masse globale de X
 - $g_i = \frac{s_i}{S}$ la fréquence de la masse de X possédée par les individus dans la classe $[x_{i-1}, x_i[$.
 - $G_i = \sum_{j=1}^i g_j$ La masse cumulée relative à la classe $[x_{i-1}, x_i[$.

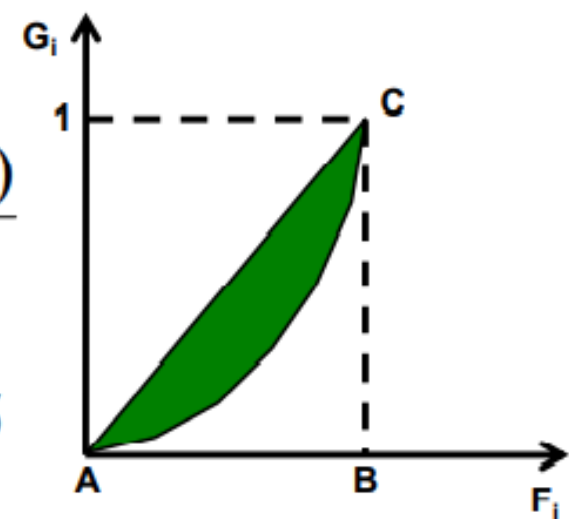


L'indice de Gini est égal à:

$$I_G = \frac{\text{aire de concentration (en vert)}}{\text{aire du triangle } ABC}$$

$$L'aire du triangle ABC = \frac{1 \times 1}{2} = 0,5$$

$$\text{Donc } I_G = 2 \times (\text{aire de concentration})$$





Séries doubles

Tableau de contingence

		Modalités de Y							
		Y	y ₁	...	y _j	...	y _q	Total	
Modalités de X	X								
	x ₁	n ₁₁					n _{1q}	n _{1.}	Distribution marginale de X
	
	x _i				n _{ij}			n _{i.}	
	
	x _p	n _{p1}					n _{pq}	n _{p.}	
Total		n _{.1}	...	n _{.j}	...	n _{.q}	n _{..}		Distribution marginale de Y

$$n_{i.} = \sum_{j=1}^q n_{ij}$$
$$n_{.j} = \sum_{i=1}^p n_{ij}$$
$$n_{..} = \sum_{i=1}^p \sum_{j=1}^q n_{ij}$$

Éléments d'un tableau de contingence

Les effectifs

Les effectifs partiels: n_{ij}

Les effectifs marginaux

Les fréquences

Les fréquences partielles

$$f_{ij} = \frac{n_{ij}}{n_{..}}$$

Les fréquences conditionnelles

- de X selon Y

$$f_{i|j} = \frac{n_{ij}}{n_{.j}}$$

- de Y selon X

$$f_{j|i} = \frac{n_{ij}}{n_{i.}}$$

Les fréquences marginales

$$f_{i.} = \frac{n_{i.}}{n_{..}} \quad \text{et} \quad f_{.j} = \frac{n_{.j}}{n_{..}}$$

Relations entre les fréquences marginales et conditionnelles

$$f_{i.} \times f_{j|i} = f_{.j} \times f_{i|j} = f_{ij}$$

Indépendance de deux variables

Deux variables X et Y sont totalement indépendantes si les fréquences $f_{i/j}$ conditionnelles ne dépendent plus de j.

$$f_{i/j} = \frac{n_{ij}}{n_{\cdot j}} = \frac{n_{i1}}{n_{\cdot 1}} = \frac{n_{i2}}{n_{\cdot 2}} = \dots = \frac{n_{iq}}{n_{\cdot q}} = \frac{\sum_{j=1}^q n_{ij}}{\sum_{j=1}^q n_{\cdot j}} = \frac{n_{i\cdot}}{n_{\cdot\cdot}} = f_{i\cdot}$$
$$\Rightarrow \frac{n_{ij}}{n_{\cdot j}} = \frac{n_{i\cdot}}{n_{\cdot\cdot}} \Rightarrow n_{ij} = \frac{n_{i\cdot} n_{\cdot j}}{n_{\cdot\cdot}} \Rightarrow \frac{n_{ij}}{n_{\cdot\cdot}} = \frac{n_{i\cdot}}{n_{\cdot\cdot}} \times \frac{n_{\cdot j}}{n_{\cdot\cdot}} \text{ donc } \boxed{f_{ij} = f_{i\cdot} \times f_{\cdot j}}$$

Exercice

Calculer les fréquences partielles et les fréquences marginales.
Montrer que les caractères X et Y sont indépendants.

X \ Y	y₁	y₂	Total
x₁	3	5	8
x₂	6	10	16
Total	9	15	24

La moyenne marginale, La variance marginale , L'écart-type, La covariance

$$\bar{x} = \frac{1}{n_{..}} \sum_{i=1}^p n_{i.} x_i \quad \text{avec} \quad n_{..} = \sum_{i=1}^p n_{i.} = \sum_{j=1}^q n_{.j} = \sum_{i=1}^p \sum_{j=1}^q n_{ij}$$

$$V(x) = \frac{1}{n_{..}} \sum_{i=1}^p n_{i.} (x_i - \bar{x})^2 \quad \text{ou} \quad V(x) = \frac{1}{n_{..}} \sum_{i=1}^p n_{i.} x_i^2 - \bar{x}^2$$

$$\sigma_x = \sqrt{V(x)}$$

$$\text{cov}(x, y) = \frac{1}{n_{..}} \sum_{i=1}^p \sum_{j=1}^q n_{ij} (x_i - \bar{x})(y_j - \bar{y})$$

Exercice

Pour 25 ménages, les âges de l'époux et de l'épouse, relevés sur le registres d'un état civil sont les suivants:

(22,17); (23,18); (24,17); (24,18); (24,20); (24,21); (25,18); (25,19); (25,20);
(26,18); (26,19); (26,21); (26,23); (27,19); (27,21); (28,21); (28,22); (30,22);
(30,23); (31,24); (31,25); (34,24); (35,24); (35,25); (36,25);

Sachant que chaque couple (x,y) représente respectivement l'âge de l'époux et l'âge de l'épouse au moment de mariage.

1. Ranger les données en classes de même amplitude 5, qui commencent par 20 pour X et par 15 pour Y.
2. Calculer l'âge moyenne des époux et des épouses.
3. Calculer la variance de l'âge d'épouse, et son écart-type.
4. Calculer la covariance des deux variables

Etude de liaison entre deux variables

Nulle(aucune influence), totale, relative

Notion de corrélation

même sens, sens inverse

Coefficient de corrélation

- Si r est proche de 1, il y a une forte corrélation positive entre X et Y (même sens de variation)
- Si r est proche de -1, il y a une forte corrélation négative entre X et Y (différence du sens de variation).
- Si $r=0$, X et Y sont non corrélées : il n'y a pas d'association linéaire entre X et Y .
- Si $r = \pm 1$, alors chacune de ces deux variables peut définir l'autre d'une façon exacte

$$r = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Ajustement par la méthode des moindres carrés

trouver les coefficients a et b de la droite de régression $y=ax+b$, qui minimisent la distance quadratique entre et qui revient à minimiser: $S(a,b)$

$$S(a,b) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (ax_i + b))^2$$

$$\hat{a} = \frac{\text{cov}(X, Y)}{V(X)} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x}\bar{y}}{\frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2}; \quad \hat{b} = \bar{y} - \hat{a}\bar{x}$$

Exercice

Soit X la note des mathématiques sur 20 points et Y la note de statistique sur 20 points pour 10 étudiants:

X	2	4	6	6	9	10	11	12	13	18
Y	3	6	6	7	9	10	10	11	14	14

1. Donner la droite des moindres carrés de Y en X ,
2. Donner la droite des moindres carrés de X en Y ,