



S4 2019-2020

# Statistiques

MINFO0401



Nathan TONNELLE

## Table des matières

Statistique descriptive.....	3
Vocabulaire .....	3
Population .....	3
Echantillon.....	3
Caractère et modalité .....	3
Représentation graphique .....	4
Variables qualitatives .....	4
Variable quantitative.....	4
La fonction cumulative .....	5
Variables continues.....	5
Définition (représentation graphique différentielle) .....	5
Exemple.....	5
Courbe cumulative .....	6
Description numérique d'une variable .....	6
Paramètres de position .....	6
Paramètre de dispersion .....	7
Coefficient de variation.....	8
Les moments .....	8
Caractéristiques de forme.....	9
Distributions à deux dimensions.....	10
Notations.....	10
Distribution des fréquences du couple (X,Y).....	10
Notation .....	10
Remarque.....	10
Distribution marginale .....	11
Moyennes et variances marginales .....	11
Distributions conditionnelles .....	11
Définition.....	11
Remarque.....	11
Indépendance des variables X et Y .....	12
Conséquence .....	12
Moyennes et variances conditionnelles .....	12
Résultat .....	12
Notion de corrélation.....	13
Coefficient de corrélation linéaire .....	13
Rapports de corrélation .....	13

## MINFO 0401

Définition.....	13
Propriété .....	14
Courbes de régression .....	14
Définition.....	14
Commentaire .....	14
Liaison entre deux variables .....	14

# Statistique descriptive

## Vocabulaire

### Population

Une population est l'ensemble des individus ou objets sur lesquels portent une étude statistique, on le note  $P$

#### Exemple

- 1) Logement d'une ville
- 2) Personnel d'une entreprise
- 3) Animaux d'un parc naturel

### Echantillon

Un échantillon est une partie de la population e étudier sur laquelle porte l'étude statistique.

Une étude statistique portant sur un échantillon est appelée sondage.

On appelle étude associative ou recensement si elle porte sur l'ensemble de la population.

### Caractère et modalité

Une étude statistique porte sur 1 ou plusieurs caractères communs à tous les individus de la population à étudier. Un caractère est aussi appelé variable.

#### Exemple

- 1) Surface de logement
- 2) Age, ancienneté, revenu
- 3) Régime alimentaire, espèces

#### Modalité

Les modalités d'un caractère sont les différentes valeurs que peut prendre se caractère sur les individu de la population étudiée

#### Exemple

- 1)  $\mathbb{R}^+$
- 2)  $[16,10)$
- 3) Herbivore, carnivore, omnivore

On distingue 2 types de caractères :

- Caractère qualitatif :

Les modalités sont des attributs qualitatif (ex: régime alimentaire)

- Caractère quantitatif

Les modalités sont des quantités numériques (ex : age)

Les caractères quantitatifs sont de 2 types

- Variable discrète :

Les modalités de la variable appartiennent à un ensemble discret tel  $\mathbb{N}$ ,  $\mathbb{Z}$ ,  $\mathbb{N}^2$

- Variable continue :

Les modalités de la variable prennent des valeurs dans un ensemble continu tel que  $\mathbb{R}$

#### Remarque :

Pour étudier une variable continue on constitue des classes de valeurs possibles, ces classes sont des intervalles d'amplitude égale ou inégale est constitué alors de nouvelle modalité ou des caractères.

#### Attention :

Le découpage en classes peut influencer sur les résultats et les interprétations que l'on peut faire.

S'il est trop important, il risque de faire apparaître des irrégularités artificielles car les effectifs des classes seront trop faibles.

S'il est trop grossier, il conduira à une perte d'information.

#### Effectif et fréquence

L'effectif d'une modalité ou d'une classe de modalité est le nombre d'individu de la population correspondant à cette modalité ou à cette classe de modalités. On ne note  $n_i$  pour la  $i^{\text{ème}}$  modalité.

La fréquence de la  $i^{\text{ème}}$  modalité (ou classe) est donnée par le rapport de son effectif sur l'effectif total de la population noté  $n$ .

On a alors  $f_i = n_i/n$

*Remarque*

Si la variable X possède K modalités  $x_1, \dots, x_k$  alors

$$\sum_{i=1}^k n_i = n \quad \sum_{i=1}^k f_i = 1$$

**Représentation graphique**

Il existe différentes façons de représenter graphiquement des variables

**Variables qualitatives***Diagramme circulaire*

C'est un disque dans lequel chaque modalité est représenté par un secteur angulaire proportionnel à sa fréquence.

Pour  $1 \leq i \leq f$   $\alpha_i = 360 \times f_i$

*Diagramme en tuyaux d'orgue*

C'est un diagramme formé de rectangles tous de même largeur et dont les hauteurs sont proportionnelles aux fréquences des modalités

**Variable quantitative**

Il y a 2 sortes de représentation graphique des variables quantitatives :

- Diagramme différentiel
- Diagramme intégrale

*Variables discrètes*

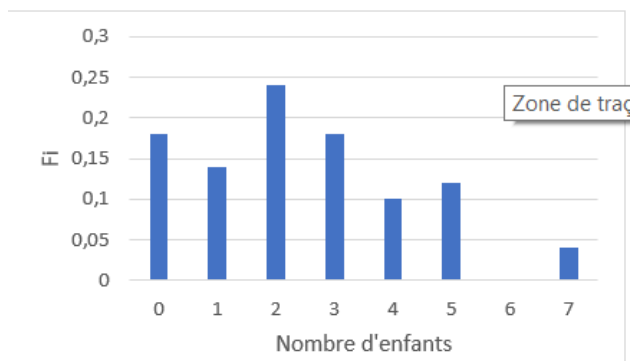
Le diagramme différentiel utilisé dans le cas d'une variable discrète est le diagramme à bâton.

Où chaque bâton est de longueur ou de hauteur proportionnel à la fréquence de la modalité correspondante

**Exemple**

Nombre d'enfants dans un échantillon de 50 familles.

Nombre d'enfant	$n_i$	$f_i$
0	9	0.18
1	7	0.14
2	12	0.24
3	9	0.18
4	5	0.10
5	6	0.12
6	0	0
7	2	0.04



## La fonction cumulative

La fonction cumulative d'une variable  $X$  prise en un point  $x$ , noté  $F(x)$  est définie comme la proportion de la population pour laquelle la variable  $X$  prend des valeurs  $\geq x$

Si les modalités de  $X$  sont  $x_1, x_2, \dots, x_k$  alors pour  $x_i \leq x < x_{i+1}$ ,  $F(x) = \sum_{j=1}^i f_j$

Exemple :

Si  $x < 0$   $F(x) = 0$

$0 \leq x < 1$   $F(x) = 0.18$

$1 \leq x < 2$   $F(x) = 0.32$

$2 \leq x < 3$   $F(x) = 0.56$

$3 \leq x < 4$   $F(x) = 0.74$

$4 \leq x < 5$   $F(x) = 0.84$

$5 \leq x < 7$   $F(x) = 0.86$

$x \geq 7$   $F(x) = 1$

## Variables continues

Si  $e_i$  et  $e_{i+1}$  sont les extrémités de la classe  $n^o i$ , noté  $[e_i, e_{i+1}]$ , on notera  $c_i$  soit milieu et  $a_i$  son amplitude,  $1 \leq i \leq k$

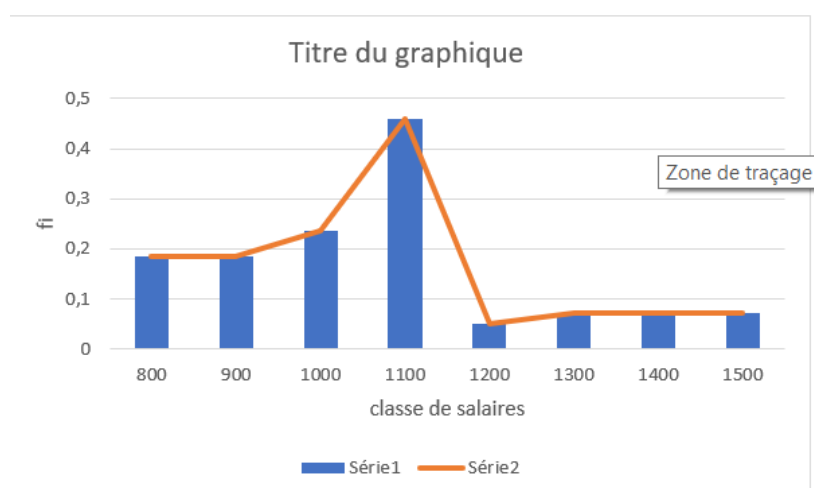
### Définition (représentation graphique différentielle)

Un histogramme est une représentation graphique où chaque classe est représenté par un rectangle de base proportionnelle à son amplitude et de surface proportionnelle à sa fréquence. Ainsi la hauteur de la classe  $n^o i$  est  $h_i = f_i / a_i$ . Un polygone statistique est un polygone reliant le milieu des bases supérieurs des rectangles de l'histogramme.

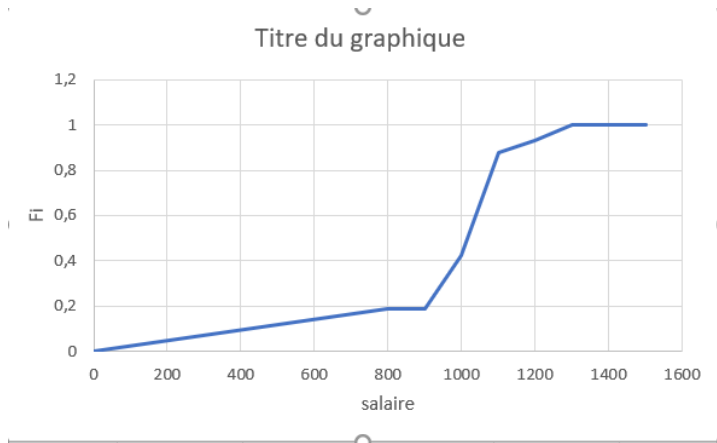
### Exemple

Salaire mensuel net des ouvriers d'un établissement industriel.

Classes de salaire	$n_i$	$f_i$	$a_i$	$h_i = f_i / a_i$ (échelle x100)	$F_i$
[800,1000[	26	0.186	200	0.096	0.186
[1000,1100[	33	0.235	100	0.235	0.421
[1100,1200[	64	0.458	100	0.458	0.879
[1200,1300[	7	0.05	100	0.05	0.929
[1300,1500[	10	0.071	200	0.0355	1
total	140	1	700		



## Courbe cumulative



## Description numérique d'une variable

### Paramètres de position

#### Médiane

##### Définition

La médiane  $\mu$  (grec) ( $u$ ) est la valeur de la variable  $x$  pour laquelle la moitié au moins des observations sont supérieures ou égales et la moitié au moins des observations inférieures ou égales.

##### Remarque

Dans le cas continu on détermine d'abord la classe médiane avant de calculer le point médian par la méthode d'interpolation linéaire.

##### Exemple

1. (nombre d'enfants par famille)  $u=2$  car ( $F_2=0.32$  et  $F_3=0.56$ )
2. (salaires des ouvriers)

Classe médiane =  $[1100, 1200[$

Calcul de la médiane par la méthode d'interpolation linéaire

$$\begin{aligned} \Leftrightarrow \frac{u - 1100}{0.5 - 0.421} &= \frac{1200 - 1100}{0.879 - 0.421} \\ \Leftrightarrow \frac{u - 1100}{0.079} &= \frac{100}{0.458} \\ \Leftrightarrow u - 1100 &= 0.079 * \frac{100}{0.458} = 17.249 \\ \Leftrightarrow u &= 1117.249 \end{aligned}$$

#### Le mode

##### Définition

Le mode est la valeur de la variable  $x$  ayant la plus grande fréquence

##### Remarque

Certaines séries statistiques peuvent avoir plusieurs modes. Dans le cas continu on parle de classe modale, on veillera cependant à tenir compte de l'amplitude des classes. La classe modale correspond à la classe ayant la plus grande hauteur  $h_i$ .

##### Exemple

1. (salaire des ouvriers)

Classe modale =  $[1100, 1200[$

Equation (M M2)

$$\begin{aligned} a &= \frac{0.458 - 0.235}{1200 - 1100} = 0.223 * 10^{-2} \\ b &= 0.548 - 0.223 * 10^{-2} * 1200 = -2.218 \\ y &= 0.223 * 10^{-2}x - 2.218 \end{aligned}$$

Equation(M3 M4)

$$a = \frac{0.05 - 0.458}{1200 - 1100} = -0.408 * 10^{-2}$$

$$b = 0.458 + 0.408 * 10^{-2} * 1100 = 4.946$$

$$y = -0.408 * 10^{-2}x + 4.946$$

$$0.223 * 10^{-2}x - 2.218 = -0.408 * 10^{-2}x + 4.946$$

$$\Leftrightarrow x = \frac{7.164}{0.631} * 10^{-2} = 135.34$$

*la moyenne*

Soit  $x$  une variable prenant les valeurs  $x_1, \dots, x_k$  avec les effectifs  $n_1, \dots, n_k$  respectivement (avec les fréquences  $f_1, \dots, f_k$ ). Alors la moyenne de la variable  $x$  est donnée par

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

Où  $n = \sum_{i=1}^k n_i$

Dans le cas continu, la moyenne d'une variable  $x$  est définie par

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

Où les  $C_i$  sont les milieux des classes  $[e_i, e_{i+1}[$

*Propriétés*

## 1) Linéarité

Si on considère la transformation  $Y = aX + b$

Alors la moyenne de  $Y$  est  $\bar{Y} = a\bar{X} + b$

$$2) \sum_{i=1}^k n_i (x_i - \bar{x}) = 0$$

3) Si on définit la fonction

$$L(c) = \sum_{i=1}^k n_i (x_i - c)^2$$

Alors  $L(c)$  prend son minimum pour  $c = \bar{X}$

4) Si  $P = P_1 \cup P_2$ 

Où la moyenne de  $X$  sur  $P_1$  est  $\bar{X}_1$ , et l'effectif est  $n_1$ ; la moyenne de  $X$  sur  $P_2$  est  $\bar{X}_2$  et l'effectif de  $P_2$  est  $n_2$ , alors la moyenne de  $X$  sur  $P$  est :

$$\bar{X} = \frac{n_1 * \bar{X}_1 + n_2 * \bar{X}_2}{n_1 + n_2}$$

*Paramètre de dispersion**Etendu*

L'étendu est la différence entre la plus grande et la plus petite valeur de la série statistique.

*L'écart moyen absolu*

On considère une série de  $n$  observations dont les modalités sont  $x_1, \dots, x_k$ .

On note  $\bar{X}$  sa moyenne empirique.

Alors l'écart-type moyen absolu est défini par

$$e_n = \frac{1}{n} \sum_{i=1}^k n_i |x_i - \bar{x}|$$

*Variance (écart quadratique moyen)*

Si  $x_1, \dots, x_n$  sont les modalités d'une variance  $X$  observée  $n$  fois, alors la variance de  $X$  est définie par

$$\text{Var}(X) = \frac{1}{2} \sum_{i=1}^k n_i (x_i - \bar{x})^2$$

*Propriété*

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

$$2) \text{ Si } Y = aX + b, \text{ alors } \text{Var}(Y) = a^2 \text{Var}(X)$$



## Preuve

$$\begin{aligned}
Var(Y) &= \frac{1}{n} \sum_{i=1}^k n_i (y_i - \bar{Y})^2 \\
&= \frac{1}{n} \sum_{i=1}^k n_i (ax_i + b - (a\bar{X} + b))^2 \\
&= \frac{1}{n} \sum_{i=1}^k n_i (ax_i - a\bar{X})^2 \\
&= (a^2) \frac{1}{n} \sum_{i=1}^k n_i (x_i - \bar{X})^2 \\
&= a^2 Var(X)
\end{aligned}$$

3) Si  $P_1$  et  $P_2$  sont 2 sous populations d'une population  $P$  tel que les moyennes, les variances et les effectifs de  $X$  sont :

Pour  $P_1$  :  $\bar{X}_1, Var_1(X), n_1$

Pour  $P_2$  :  $\bar{X}_2, Var_2(X), n_2$

Alors la variance de  $X$  sur l'ensemble de  $P$  est

$$Var(X) = \frac{n_1 Var_1(X) + n_2 Var_2(X)}{n_1 + n_2} + \frac{n_1 (\bar{X}_1 - \bar{X})^2 + n_2 (\bar{X}_2 - \bar{X})^2}{n_1 + n_2}$$

Où  $\bar{X} = \frac{n_1 \bar{X}_1 + n_2 \bar{X}_2}{n_1 + n_2}$  est la moyenne de  $X$  sur  $P$

$Var(X) = \text{moyenne des variances} + \text{variances des moyennes} = \text{variances intra} + \text{variance inter}$

## Ecart type

L'écart type est défini par :

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

## Intérêt :

$X$  est en  $Km \Rightarrow$  en  $Km^2$

L'écart type exprime la dispersion dans la même unité de mesure que la variable  $X$

## Coefficient de variation

Le coefficient de variation est défini par

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

## Intérêt :

Ce coefficient est indépendant de l'unité de mesure. Il permet alors de comparer les dispersions de sens statistiques exprimés dans des unités de mesure différents.

## Les moments

Définition

On appelle moment d'ordre  $t$  ( $t$  appartient à  $N$ ) par rapport à une constante  $a$  d'une variable statistique  $x$

$$m_t(a) = \frac{1}{n} \sum_{i=1}^k n_i (x_i - a)^t$$

Les moments non centrés correspondent à

$$a = 0, ie : m_t = \frac{1}{n} \sum_{i=1}^k n_i x_i^t$$

Les moments centrés correspondent à  $a = \bar{X}$

ie :  $\mu_t = m_t(\bar{X}) = \frac{1}{n} \sum_{i=1}^k n_i (x_i - \bar{X})^t$

Remarque

$$m_1(0) = m_1 = \frac{1}{n} \sum_{i=1}^k n_i x_i = \bar{X}$$

$$m_2(\bar{X}) = \frac{1}{n} \sum_{i=1}^k n_i (x_i - \bar{X})^2 = \text{Var}(X)$$

Caractéristiques de forme

Coefficient d'asymétrie

Le coefficient d'asymétrie de Fisher est défini par

$$\gamma^1 = \frac{\mu_3}{\sigma_x^3}$$

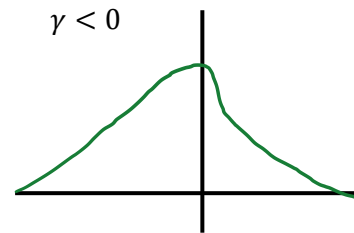
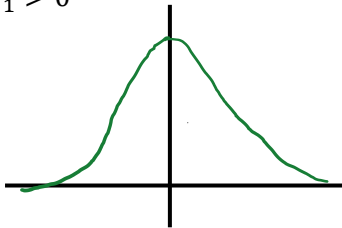
Si  $A \Rightarrow B$  alors  $\bar{B} \Rightarrow \bar{A}$

Remarque :

Si la distribution est symétrique, alors

$$\gamma_1 = 0$$

Si  $\gamma_1 \neq 0$  alors la distribution n'est pas symétrique  
 $\gamma_1 > 0$

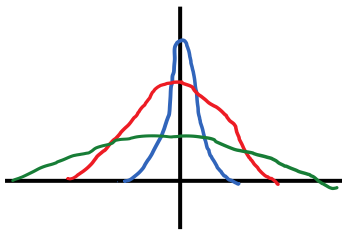


Distribution non symétrique, oblique à gauche et étroite à droite

Coefficient d'aplatissement (Kurtosis)

$$\gamma_a = \frac{\mu_4}{\sigma_x^4} - 3$$

Il compare l'aplatissement de courbe statistique à la courbe de la loi  $N(0,1)$ .



$$\gamma_2 < 0$$

$$\gamma_2 < 0$$

$$\gamma_2 < 0$$

Si la courbe est identique à celle de la loi  $N(0,1)$  alors  $\gamma_2 = 0$ .

Si  $\gamma_2 \neq 0$ , la courbe ne présente pas le même aplatissement que la courbe de la loi  $N(0,1)$ .

Box-plot (boîte à moustache)

Définition :

Premier quartile  $q_1$  est la valeur de la variable  $x$  pour laquelle au moins  $1/4$  des observations lui sont inférieures ou égales et au moins  $3/4$  des observations supérieures ou égales.

Le troisième quartile  $q_3$  est la valeur de la variable  $x$  pour laquelle au moins  $3/4$  des observations lui sont inférieures ou égales et au moins  $1/4$  des observations supérieures ou égales.

Remarque :

Dans le cas d'une variable continue, on détermine d'abord les classes contenant le premier quartile  $q_1$  et le 3ème quartile  $q_3$  avant de procéder par la méthode de l'interpolation linéaire au calcul de  $q_1$  et  $q_3$

Le box-plot :



## Distributions à deux dimensions

On considère un couple de variables (X,Y) observé sur une population P de taille n. on suppose que la variable X possède les modalités  $x_1, \dots, x_k$  et que les modalités de la variable y sont  $y_1, \dots, y_l$ . Ces modalités peuvent être des valeurs discrètes ou des classes de modalités (cas où les variables sont continues). L'observation du vecteur (X,Y) donne lieu à la table de contingence qui se présente sous la forme :

X\Y	$y_1$	$y_2$	...	$Y_j$	...	$Y_l$	
$x_1$	$n_{11}$	$n_{12}$		$n_{1j}$		$n_{1l}$	$n_{1\cdot}$
$x_2$	$n_{21}$	$n_{22}$		$n_{2j}$		$n_{2l}$	$n_{2\cdot}$
...							
$x_i$	$n_{i1}$	$n_{i2}$		$n_{ij}$		$n_{il}$	$n_{i\cdot}$
...							
$x_k$	$n_{k1}$	$n_{k2}$		$n_{kj}$		$n_{kl}$	$n_{k\cdot}$
	$n_{\cdot 1}$	$n_{\cdot 2}$		$n_{\cdot j}$		$n_{\cdot l}$	$n_{\cdot\cdot} = n$

Où  $n_{ij}$  est l'effectif connu portant un nombre de fois où la modalité  $x_i$  de x et la modalité de  $y_j$  de Y sont observé, simultanément sur les individus de la population P.

### Notations

$$n_{i\cdot} = \sum_{j=1}^l n_{ij}, n_{\cdot j} = \sum_{i=1}^k n_{ij}$$

$$n_{\cdot\cdot} = \sum_{i=1}^k n_{i\cdot} = \sum_{j=1}^l n_{\cdot j} = \sum_{i=1}^k \sum_{j=1}^l n_{ij} = n$$

La table de contingence donne la distribution des effectifs du couple (X,Y).

### Distribution des fréquences du couple (X,Y)

La fréquence de la modalité  $(x_i, y_j)$  de (X,Y) est définie par

$$f_{ij} = \frac{n_{ij}}{n}, \quad 1 \leq i \leq k \text{ et } 1 \leq j \leq l$$

### Notation

Pour  $1 \leq i \leq k$  et  $1 \leq j \leq l$ , on note  $f_{i\cdot} = \sum_{j=1}^l f_{ij} = \frac{n_{i\cdot}}{n}$

$$f_{\cdot j} = \sum_{i=1}^k f_{ij} = \frac{n_{\cdot j}}{n}$$

### Remarque

$$\sum_{i=1}^k f_{i\cdot} = \sum_{j=1}^l f_{\cdot j} = \sum_{i=1}^k \sum_{j=1}^l f_{ij} = 1$$

## Distribution marginale

La distribution marginale des effectifs de X est de définir par  $(n_1, \dots, n_k)$ . la distribution marginale des fréquences de X est  $(f_1, \dots, f_k)$ .

Les distributions marginales des effectifs et des fréquences de la variable Y sont respectivement  $(n_1, \dots, n_l)$  et  $(f_1, \dots, f_l)$ .

## Moyennes et variances marginales

La moyenne et la variance marginales de X sont définies par

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

$$Var(X) = \frac{1}{n} \sum_{i=1}^k n_i \cdot (x_i - \bar{X})^2 = \sum_{i=1}^k f_i \cdot (x_i - \bar{X})^2 = \frac{1}{n} \sum_{i=1}^k n_i \cdot x_i^2 - \bar{X}^2 = \sum_{i=1}^k f_i \cdot x_i^2 - \bar{X}^2$$

La moyenne de la variance marginale de Y sont :

$$\bar{Y} = \frac{1}{n} \sum_{j=1}^l n_j \cdot y_j = \sum_{j=1}^l f_j \cdot y_j$$

$$Var(Y) = \frac{1}{n} \sum_{j=1}^l n_j \cdot (y_j - \bar{Y})^2 = \sum_{j=1}^l f_j \cdot (y_j - \bar{Y})^2 = \frac{1}{n} \sum_{j=1}^l n_j \cdot y_j^2 - \bar{Y}^2 = \sum_{j=1}^l f_j \cdot y_j^2 - \bar{Y}^2$$

## Distributions conditionnelles

### Définition

La fréquence conditionnelle de  $X=x_i$  sachant  $Y=y_j$  est définie, pour  $1 \leq i \leq k$  et  $1 \leq j \leq l$ , par

$$f_{\frac{i}{j}} = \frac{n_{ij}}{n_j}$$

La fréquence conditionnelle de  $Y=y_j$  sachant  $X=x_i$  est définie pour tous  $1 \leq i \leq k$  et  $1 \leq j \leq l$  par

$$f_{\frac{j}{i}} = \frac{n_{ij}}{n_i}$$

La distribution conditionnelle de X sachant  $Y=y_j$  est

$$\left( \frac{n_{1j}}{n_j}, \frac{n_{2j}}{n_j}, \dots, \frac{n_{kj}}{n_j} \right)$$

La distribution conditionnelle de Y sachant  $X=x_i$  est

$$\left( \frac{n_{i1}}{n_i}, \frac{n_{i2}}{n_i}, \dots, \frac{n_{il}}{n_i} \right)$$

### Remarque

$$f_{\frac{i}{j}} = \frac{n_{ij}}{n_j} = \frac{\frac{n_{ij}}{n}}{\frac{n_j}{n}} = \frac{f_{ij}}{f_j}$$

Par conséquent

$$f_{ij} = f_{.j} * f_{i.}$$

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

Conséquence

$$f_{ij} = f_{i.} * f_{.j}$$

### Indépendance des variables X et Y

Les variables X et y sont indépendantes lorsqu'on a

$$f_{i.} = f_{i.} \text{ ou } f_{.j} = n_{.j} \quad 1 \leq i \leq k \text{ et } 1 \leq j \leq l$$

### Conséquence

Si X et Y sont indépendants alors

$$f_{ij} = f_{i.} * f_{.j}$$

Cela peut s'exprimer aussi sous la forme

$$\frac{n_{ij}}{n} = \frac{n_{i.}}{n} * \frac{n_{.j}}{n} \leftrightarrow n_{ij} = \frac{n_{i.} * n_{.j}}{n}$$

### Moyennes et variances conditionnelles

La moyenne conditionnelle de X sachant Y=y<sub>j</sub> est définie par

$$\bar{X}_j = \sum_{i=1}^k f_{i.} x_i = \frac{1}{n_{.j}} \sum_{i=1}^k n_{ij} x_i$$

$$1 \leq j \leq l$$

La moyenne conditionnelle de Y sachant X=x<sub>i</sub> est

$$\bar{Y}_i = \sum_{j=1}^l f_{.j} y_j = \frac{1}{n_{i.}} \sum_{j=1}^l n_{ij} y_j$$

$$1 \leq i \leq k$$

### Résultat

Les moyens marginales et les moyennes conditionnelles sont liées par les relations suivantes :

$$\bar{X} = \sum_{j=1}^l f_{.j} \bar{X}_j$$

$$\bar{Y} = \sum_{i=1}^k f_{i.} \bar{Y}_i$$

Les variances conditionnelles de X sachant  $Y=y_j$  et sachant  $X=x_i$  sont définies respectivement par

$$Var_j(X) = \sum_{i=1}^k f_{ij} (x_i - \bar{X}_j)^2$$

$$1 \leq j \leq l$$

Et

$$Var_i(Y) = \sum_{j=1}^l f_{ij} (y_j - \bar{Y}_i)^2$$

$$1 \leq i \leq k$$

Les variances marginales de X et Y peuvent être décomposées dans les termes suivants :

$$Var(X) = \sum_{j=1}^l f_{.j} Var_j(X) + \sum_{j=1}^l f_{.j} (\bar{X}_j - \bar{X})^2$$

= moyennes des variances conditionnelles et variances des moyennes conditionnelles

$$Var(Y) = \sum_{i=1}^k f_{i.} Var_i(Y) + \sum_{i=1}^k f_{i.} (\bar{Y}_i - \bar{Y})^2$$

## Notion de corrélation

### Coefficient de corrélation linéaire

#### Définition

Le coefficient de corrélation linéaire entre deux variables X et Y est défini par

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

Où

$$cov = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^l n_{ij} x_i y_j - \bar{X} \bar{Y} = \sum_{i=1}^k \sum_{j=1}^l f_{ij} x_i y_j - \bar{X} \bar{Y} = \sum_{i=1}^k \sum_{j=1}^l f_{ij} (x_i - \bar{X})(y_i - \bar{Y})$$

#### Propriétés

- $-1 \leq \rho_{X,Y} \leq 1$
- si X et Y sont indépendants alors  $\rho_{X,Y} = 0$

#### Remarque

Si X et Y ne sont pas corrélées, cela n'implique pas forcément que X et Y sont indépendants.

Il y a identité entre indépendance et non corrélation uniquement dans le cas de variable gaussiennes.

#### Propriétés (suite)

- Si X et Y présentent une liaison fonctionnelle, i.e., il existe une fonction f tant que  $f(X) = Y$ , alors on a  $\rho_{X,Y}^2 = 1$

## Rapports de corrélation

On considère deux variables X et Y prennent les valeurs  $x_1, \dots, x_k$  et  $y_1, \dots, y_l$  respectivement.

On note  $\bar{X}$  et  $\bar{Y}$  les moyennes empiriques de X et Y respectivement.

On note  $var(X)$  et  $var(Y)$  les variances marginales de X et Y.

Si  $\bar{X}_j$  est la moyenne de X sachant  $Y = y_j$ ,  $\bar{Y}_i$  la moyenne de Y sachant  $X = x_i$ ,

$$On \text{ note } var_{inter}(X) = \frac{1}{n} \sum_{j=1}^l n_{.j} (\bar{X}_j - \bar{X})^2 = \sum_{j=1}^l f_{.j} (\bar{X}_j - \bar{X})^2$$

$$Et \text{ } var_{inter}(Y) = \frac{1}{n} \sum_{i=1}^k n_{i.} (\bar{Y}_i - \bar{Y})^2 = \sum_{i=1}^k f_{i.} (\bar{Y}_i - \bar{Y})^2$$

#### Définition

Le rapport de corrélation de X en Y est défini par

$$\eta_{X/Y}^2 = \frac{var_{inter}(X)}{var(X)} = \frac{\sum_{j=1}^l f_{.j} (\bar{X}_j - \bar{X})^2}{\sum_{i=1}^k f_{i.} (x_i - \bar{X})^2}$$

Le rapport de corrélation de Y en X est défini par

$$\eta_{\frac{Y}{X}}^2 = \frac{Var_{inter}(Y)}{Var(Y)} = \frac{\sum_{i=1}^k f_i (\bar{Y}_i - \bar{Y})^2}{\sum_{j=1}^l f_j (Y_j - \bar{Y})^2}$$

### Propriété

- 1)  $0 \leq \eta_{\frac{X}{Y}}^2 \leq 1$  ET  $0 \leq \eta_{\frac{Y}{X}}^2 \leq 1$
- 2) Si  $\eta_{\frac{X}{Y}}^2 = 0$ , cela veut dire que  $\bar{X}_j = \bar{X}$  pourtant  $1 \leq j \leq l$   
Le conditionnement de X par Y ne donne aucune liaison significative
- 3) Si  $\eta_{\frac{X}{Y}}^2 = 1$ , cela veut dire que  $\sum_{j=1}^l f_j (\bar{X}_j - \bar{X})^2 = Var(X)$   
X est lié de façon complète à Y
- 4) Plus  $\eta_{\frac{X}{Y}}^2$  est proche de 1 plus la liaison de X à Y est importante

### Courbes de régression

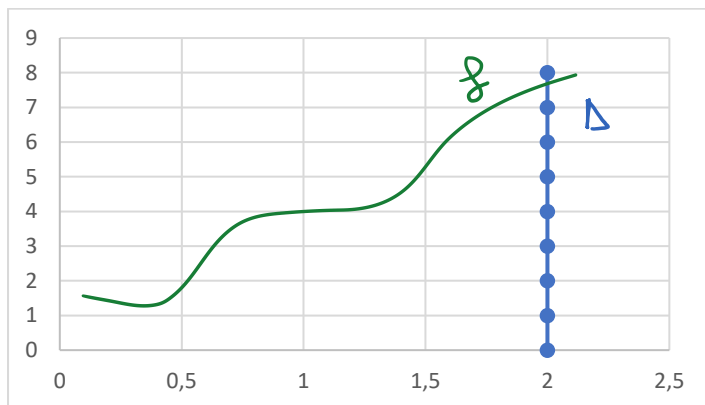
Les courbes de régression ont pour objet de donner une représentation graphique sur le plan de la distribution conjointe du vecteur (X,Y). Il y a deux courbes de régression. La courbe de régression de X en Y et la courbe de régression de Y en X.

#### Définition

Si la variable X prend ses valeurs dans l'ensemble  $\{x_1, \dots, x_k\}$  et Y dans l'ensemble  $\{y_1, \dots, y_l\}$ . On note  $\bar{X}_j$  la moyenne conditionnelle de X sachant  $Y=y_j$  et  $\bar{Y}_i$  la moyenne conditionnelle de Y sachant  $X=x_i$ . Alors la courbe de régression de X en Y notée  $C_{courbe \frac{X}{Y}}$  est la courbe passant par les points de coordonnées  $(\bar{X}_j, y_j)$   $1 \leq j \leq l$ . La courbe passait par les points de coordonnées  $(x_i, \bar{Y}_i)$   $1 \leq i \leq k$ .

#### Commentaire

Lorsque  $X=x_i$ , la valeur qui synthétise le plus la variable Y est la moyenne conditionnelle de Y sachant  $X = x_i$   $\bar{Y}_i$



$y = ax + b$

Pour ajuster une droite à une courbe on utilise le critère des moindres carrés donné dans l'exemple par

$$C \cap C(C_{courbe}, D) = \sum_{i=1}^n (f(x_i) - y_i)^2 = \sum_{i=1}^n (f(x_i) - ax_i - b)^2$$

Il faudra alors minimiser CMC(C,D) par rapport à a et b pour trouver la meilleure droite d'ajustement.

### Liaison entre deux variables

#### Liaison nulle ente variables X et Y

Cela signifie qu'il n'y a pas d'influence d'une variable sur l'autre

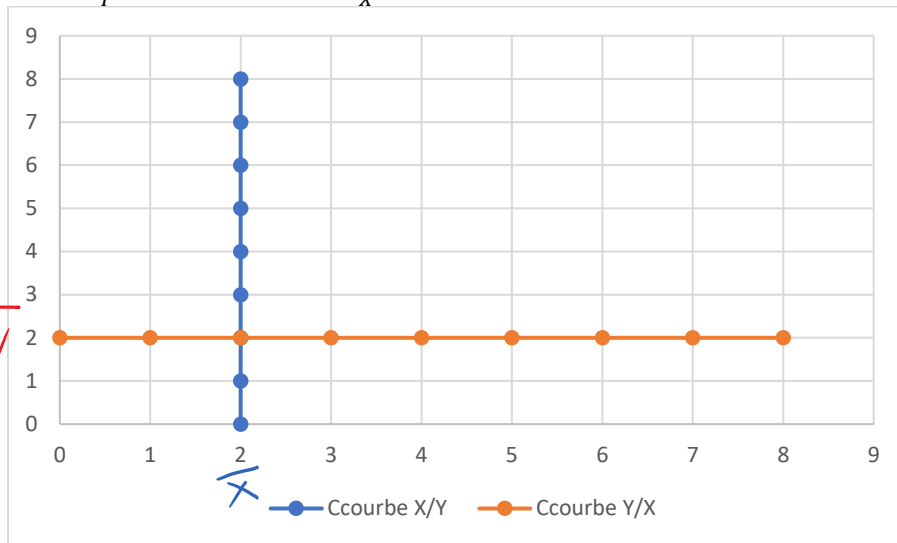
Les variables X et Y sont indépendantes, signifie que la variation de l'une des variables n'entraîne pas d'effets sur l'autre variable.

Autrement dit,  $f_i = f_i$  avec  $1 \leq i \leq k$  et  $f_j = f_j$  avec  $1 \leq j \leq l$

Ainsi on a  $\bar{X}_j = \bar{X}$ ,  $1 \leq j \leq l$  et  $\bar{Y}_i = \bar{Y}$ ,  $1 \leq i \leq k$

Les courbes de régression ont pour équations alors :

$$C_{\text{courbe } \frac{X}{Y}} = x = \bar{X} \text{ et } C_{\text{courbe } \frac{Y}{X}} = y = \bar{Y}$$



### Liaison fonctionnelle

La liaison fonctionnelle est donnée par :  $f(X) = Y$

### Remarque

Lorsque X et Y présentent une liaison fonctionnelle, les courbes de régression  $C_{\text{courbe } \frac{X}{Y}}$  et  $C_{\text{courbe } \frac{Y}{X}}$  sont confondues.

### Liaison relative

Dans ce cas le nuage de points du vecteur (X,Y) est résumé par les courbes  $C_{\text{courbe } \frac{X}{Y}}$  et  $C_{\text{courbe } \frac{Y}{X}}$  qui se coupent au point centre de gravité  $G(\bar{X}, \bar{Y})$ .

Ces courbes donnent des informations sur la nature de la liaison entre les variables X et Y.

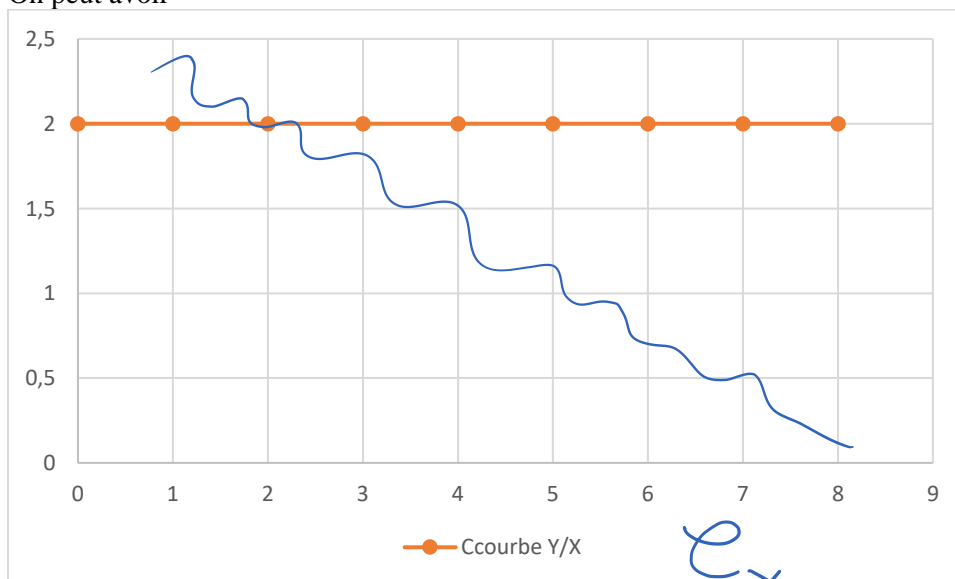
- On dira que la corrélation est positive si les deux variables varient dans le même sens.
- On dira que la corrélation est négative si les varient dans les sens opposés.
- On dira que la corrélation entre X et Y est linéaire si les courbes de régression sont des droites parallèles aux axes.

La corrélation est d'autant plus grande que l'angle formé par les courbes  $C_{\text{courbe } \frac{X}{Y}}$  et  $C_{\text{courbe } \frac{Y}{X}}$  est petit

### Remarque

A la différence avec l'indépendance la corrélation n'est pas propriété réciproque.

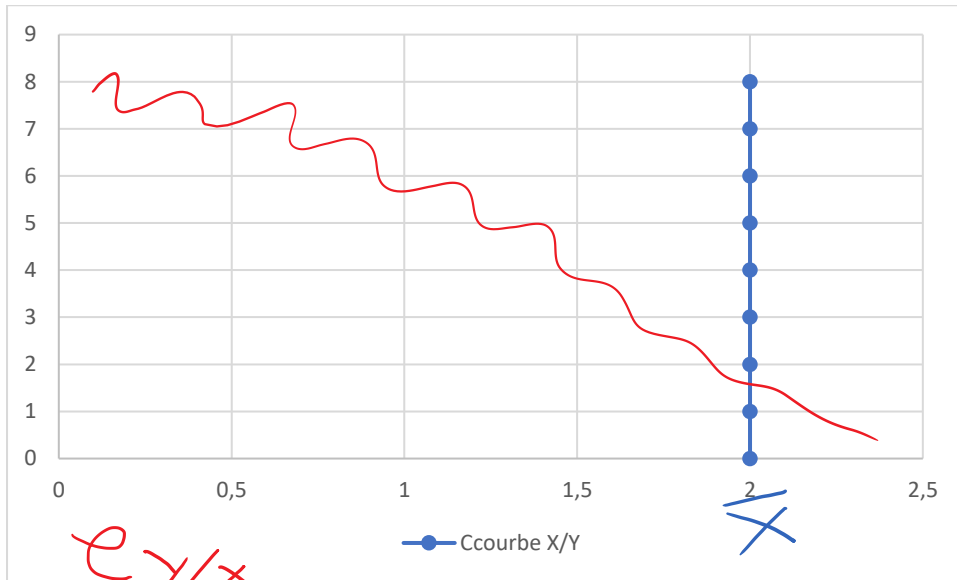
On peut avoir



X est corrélé avec Y et Y est non corrélé avec X



# MINFO 0401



Y dépend de X et X ne dépend pas de Y.