

Chapitre 1

Statistique Descriptive

I- Vocabulaires:

- a- Observations: ce sont les données relatives à un phénomène obtenues au cours d'une enquête; ou à la suite d'expériences
- b- Population: C'est l'ensemble des éléments sur lequel porte l'étude (Etudiants, bovidés, pièces mécaniques...)
- c- Individu: Un élément de la population est appelé individu.
- d- Caractère: C'est le trait particulier, ou l'aspect des individus qui nous intéresse; par exemple: poids, tailles, Couleurs.

Notation: On note par:

$\Omega \equiv$ la population

$w \equiv$ individu ($w \in \Omega$)

$C \equiv$ l'ensemble des valeurs du caractère auquel on s'intéresse

Exemple:

$\Omega \equiv$ l'ensemble des étudiants(es) de 2^{ème} année Agronomie (année 2003-2004). (population).

$w \equiv$ étudiant (individu)

On s'intéresse à la note obtenue par les étudiants en 1^{ère} année

$C = \{ 10, 10, 25, \dots, 17 \}$ (Caractère).

Variable Statistique:

Soient Ω (population) et C (le caractère) les ensembles sur lesquels porte notre étude.

L'application: X de Ω vers C qui à chaque individu ω associe la valeur du caractère $X(\omega)$ est appelée variable statistique notée V.S.

$$X: \Omega \longrightarrow C$$
$$\omega \longmapsto X(\omega)$$

Exemple:

$$X: \left\{ \begin{array}{l} \text{Promotion de 2}^{\text{ème}} \text{ année Agronomie} \\ \omega \end{array} \right\} \longrightarrow \{10, \dots, 17\}$$
$$\longmapsto X(\omega) = \text{note de } \omega$$

Si Ω contient un nombre fini d'individus

$$\Omega = \{ \omega_1, \omega_2, \dots, \omega_N \} \quad (\text{card } \Omega = N)$$

On note $X(\omega_i) = x_i$.

$$X: \Omega \longrightarrow C$$
$$\{ \omega_1, \dots, \omega_N \} \longmapsto \{ x_1, x_2, \dots, x_N \}$$
$$x_1 < x_2 < \dots < x_N$$

Si: $X: \Omega \longrightarrow C$ et $\forall \omega \in \Omega$; $X(\omega)$ est un réel alors on dira que X est une V.S. quantitative.

par contre si $X(\omega)$ n'est pas un réel

$$X: \Omega = \left\{ \begin{array}{l} \text{ensemble des voitures} \\ \text{des enseignants} \end{array} \right\} \longrightarrow \{ \text{couleurs} \}$$

X est dite V.S. qualitative.

Etude d'une variable statistique:

V. S discrète

Si $X: \Omega \rightarrow \mathbb{C}$ est une v.a.s et si \mathbb{C} contient un nombre fini d'éléments on dira que X est une v.s. discrète.

Exemple: Dans un village habite 500 familles; on s'intéresse au nombre d'enfants de chaque famille
Après recensements nous avons les données

112	familles	ont	0 enfant
170	"	"	1 enfant
128	"	"	2 enfants
60	"	"	3 enfants
30	"	"	4 enfants.

$\Omega = \{ \text{familles du village} \}$ (population)

$w = \text{famille}$ (individu)

$\mathbb{C} = \{ \begin{matrix} 0, 1, 2, 3, 4 \\ x_1, x_2, x_3, x_4, x_5 \end{matrix} \}$ le caractère \equiv nombre d'enfants.

Définition: Le nombre d'éléments de Ω est appelé effectif total noté N .

Le nombre d'éléments de $\{ w \in \Omega; X(w) = x_i \}$ s'appelle effectif partiel de la valeur x_i

il est noté n_i

($n_i = \text{Card} \{ \omega \in \Omega ; X(\omega) = x_i \}$).

-4-

Exemple: effectif total $N = 500$.

effectif partiel de la valeur $x_4 = 3$ est $n_4 = 60$.

Proposition: La somme des effectifs partiels est égal à l'effectif

Total:
$$\sum_i n_i = N$$

Définition: Soit $X: \Omega \longrightarrow \{x_1, x_2, \dots, x_n\}$.

$N = \text{Card } \Omega$ effectif total

$n_i =$ effectif partiel de x_i

Le nombre $f_i = \frac{n_i}{N}$ s'appelle fréquence partielle de la valeur x_i

Exemple:

$$f_1 = \frac{112}{500} = 0,224 ; f_2 = \frac{170}{500} = 0,340 ; f_3 = \frac{128}{500} = 0,256 ; f_4 = \frac{60}{500} = 0,120 ; f_5 = \frac{30}{500} = 0,060$$

Proposition:
$$\sum_i f_i = 1.$$

$$\sum_i f_i = \sum_i \frac{n_i}{N} = \frac{\sum_i n_i}{N} = \frac{N}{N} = 1.$$

Définition: Soit $X: \Omega \longrightarrow \{x_1, x_2, \dots, x_n\}$ une v.s. discrète

soit f_i la fréquence partielle de x_i ; alors la loi de la v.s. X

est donnée par le tableau:

x_1	x_2	x_3	...	x_n
f_1	f_2	f_3	...	f_n

Exemple: la loi de notre exemple est donnée par:

x_i	0	1	2	3	4
f_i	0,224	0,340	0,256	0,120	0,060

on doit lire: il y'a 22,4 % de familles qui ont 0 enfant.
 il y'a 12 % de familles qui ont 3 enfants.

Effectifs Cumulés et fréquences Cumulées:

Soit $X : \Omega \rightarrow C = \{x_1, x_2, \dots, x_m\}$.

n_i = effectif partiel de x_i $i = 1, 2, \dots, m$.

f_i = fréquence partielle de x_i $i = 1, \dots, m$.

alors le nombre $N_k = n_1 + n_2 + \dots + n_k$ s'appelle l'effectif cumulé de la valeur x_k .

et le nombre $F_k = f_1 + f_2 + \dots + f_k$ s'appelle fréquence cumulée de la valeur x_k .
 ($F_k = \frac{N_k}{N}$)

Exemple:

x_i	0	1	2	3	4
n_i	112	170	128	60	30
N_i	112	282	410	470	500
f_i	0,224	0,340	0,256	0,120	0,06
F_i	0,224	0,564	0,82	0,94	1.

Interprétation: N_i = nombre d'individus dont la valeur du caractère est inférieure ou égale à x_i

$N_3 = 410$ = nombre de familles qui ont au plus 2 enfants

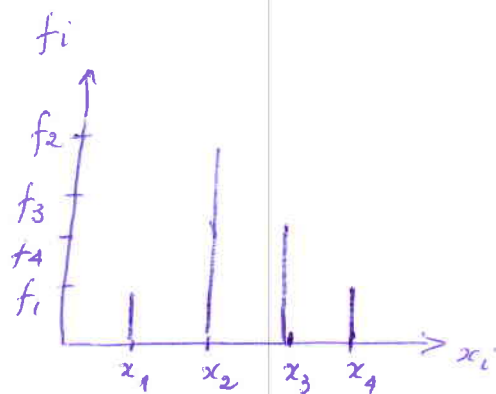
F_i = pourcentage d'individus dont la valeur du caractère est $\leq x_i$

$F_3 = 0,82 = 82\%$ de famille qui ont au plus 2 enfants.

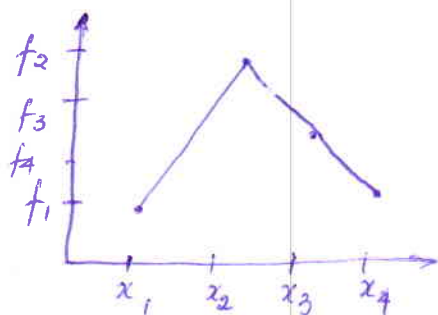
Représentation graphique:

Diagramme en batons:

On porte les fréquences partielles ou les effectifs partiels en ordonnée en fonction de x_i (ou les fréquences cumulées ou les effectifs cumulés).



Polygone des effectifs partiels.



Courbe cumulative des fréquences:

Soit $X: \Omega \rightarrow C$ une v.s. discrète

Soit l'application: $F: \mathbb{R} \rightarrow \mathbb{R}$
 $x \mapsto F(x)$

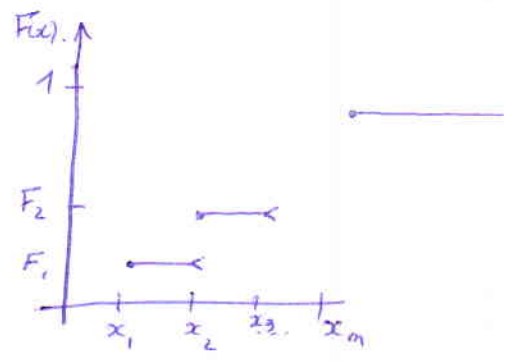
où $F(x)$ = pourcentage des individus dont la valeur du caractère est inférieure ou égale à x (F est appelée fonction de répartition).

La courbe représentative de cette fonction est appelée courbe cumulative des fréquences.

Propriétés:

1) $\forall i=1, 2, \dots, m \quad F(x_i) = F_i$

2)
$$F(x) = \begin{cases} 0 & \text{si } x < x_1 \\ F_i & \text{si } x \in [x_i, x_{i+1}[\\ 1 & \text{si } x \geq x_m \end{cases}$$



3) On note $\lim_{x \rightarrow x_i^-} F(x) = F(x_i - 0)$; $\lim_{x \rightarrow x_i^+} F(x) = F(x_i + 0)$.
 $F(x_i + 0) - F(x_i - 0) = f_i$

Variable statistique continue:

Soit: $X: \Omega \rightarrow C$ ($\text{Card}(\Omega) = N \equiv \text{effectif total}$)

Si C contient un nombre "important" de valeurs alors on dira que X est une v.s. continue

Exemple:

Ω = Ensemble des nouveaux nés du mois de Février 2004 à l'hôpital de Tlemcen. (10000)

ω = un nouveau né

$X(\omega)$ = poids de ω (Kg).

$C = \{2, 2,01, \dots, 4,5\}$.

dans ce cas on considère que $C = [2, 4,5]$ (un intervalle)

Lorsqu'il s'agira d'une V.S. continue C sera toujours considérée

comme un intervalle. $X: \Omega \longrightarrow [x_{min}, x_{max}] = [\alpha, \beta]$

$\alpha = x_{min}$ = la plus petite valeur de x

$\beta = x_{max}$ = la plus grande valeur de x

Définition: le nombre noté e donné par $x_{max} - x_{min} = \beta - \alpha = e$ est appelé étendue de la variable statistique

$$e = x_{max} - x_{min} = \beta - \alpha.$$

Pour étudier une V.S. continue $X: \Omega \longrightarrow C = [\alpha, \beta]$

nous allons diviser l'intervalle C en m classes (sous intervalles) de même longueur h .

On prendra par convention $m \approx \sqrt{N}$

et on prendra h tel que $h \times m > e$.

en pratique h est la valeur la plus simple telle que $h > \frac{e}{m}$.

$$C_1 = [\alpha, a_1[, C_2 = [a_1, a_2[, \dots , C_m = [a_{m-1}, a_m[\quad (\beta \in [a_{m-1}, a_m[)$$

Exemple:

Ω = Une récolte de pastèques (60 pastèques).

ω = une pastèque

$X(\omega)$ = diamètre de la pastèque (poids).

$$C = [10, 16, 1] \quad \text{Kgs}$$

$N = 60$; étendue $e = 16,1 - 10 = 6,1$.

m nombre de classes $m \approx \sqrt{N} = \sqrt{60}$

$m = 7$ classes.

On choisit h tel que $h > \frac{e}{m}$.

$$\frac{e}{m} = \frac{6,1}{7} = 0,87 \quad h > 0,87$$

On prendra $h = 1$ (par exemple).

On obtient ainsi les classes :

$$C_1 = [10, 11[; C_2 = [11, 12[, C_3 = [12, 13[, C_4 = [13, 14[, C_5 = [14, 15[\\ C_6 = [15, 16[, C_7 = [16, 17[.$$

$\underbrace{\hspace{1.5cm}}_{16,1}$

Effectif et fréquence d'une classe:

Soit $X: \Omega \longrightarrow C = [\alpha, \beta]$ une V.S. continue

Supposons qu'on divise C en m classes

$$C_1 = [a_0, a_1[, C_2 = [a_1, a_2[\dots C_m = [a_{m-1}, a_m[$$

$n_i = \text{Card} \{ \omega \in \Omega , X(\omega) \in C_i \}$ s'appelle effectif partiel de la classe C_i

$f_i = \frac{n_i}{N}$ s'appelle fréquence partielle de la classe C_i

$N_k = n_1 + n_2 + \dots + n_k$ s'appelle effectif cumulé de la classe C_k .

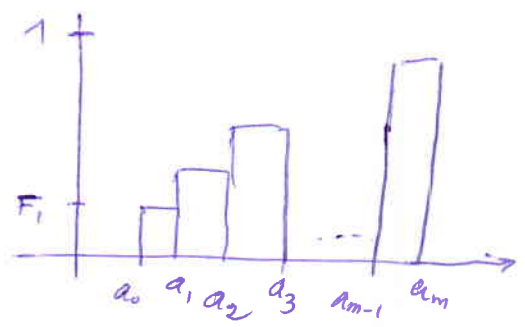
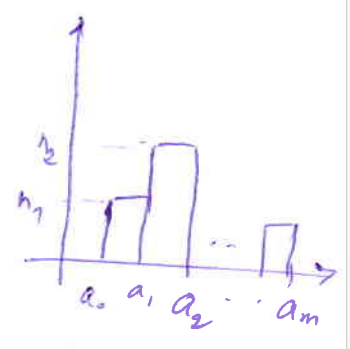
$F_k = f_1 + f_2 + \dots + f_k$ s'appelle fréquence cumulée de la classe C_k .

La loi de la variable statistique X est donnée par le tableau

C_i	C_1	C_2	...	C_m
f_i	f_1	f_2		f_m

Représentation graphique:

histogramme de effectifs (fréquences) partielles: (cumulés)



Fonction de répartition:

$X: \Omega \longrightarrow C = [a_0, a_1[\cup [a_1, a_2[\cup \dots \cup [a_{m-1}, a_m[$

Soit $F: \mathbb{R} \longrightarrow \mathbb{R}$

$F(x)$ = pourcentage des individus dont la valeur du

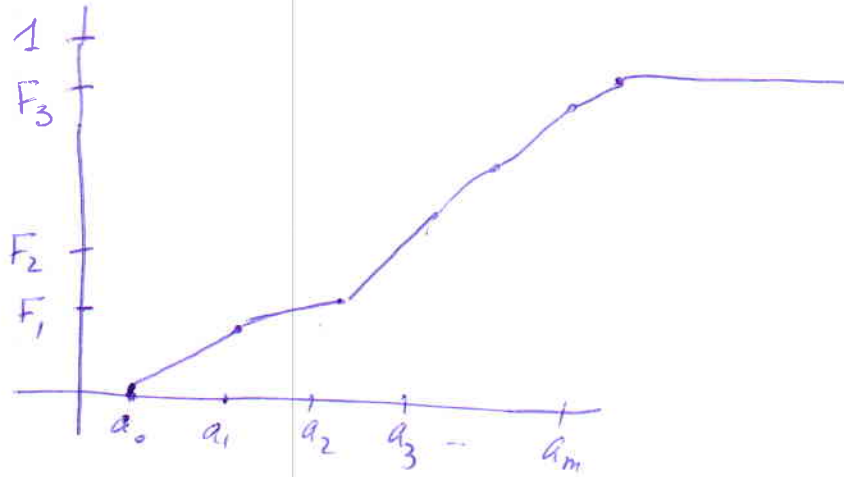
Caractère est inférieure ou égale à x . (fonction de répartition)

la courbe représentant F est appelée courbe cumulative des fréquences.

F est une fonction vérifiant:

- 1- nulle sur l'intervalle $] -\infty, a_0 [$
- 2- elle vaut 1 sur $[a_m, +\infty [$
- 3 elle passe par les points: $(a_0, 0), (a_1, F_1), (a_2, F_2), \dots, (a_m, F_m) = 1$

$$F(x) = \begin{cases} 0 & \text{if } x < a_0 \\ \frac{f_1}{h} (x - a_0) & \text{if } x \in [a_0, a_1[\\ F_i + \frac{f_{i+1}}{h} (x - a_i) & \text{if } x \in [a_i, a_{i+1}[\\ 1 & \text{if } x \geq a_m \end{cases}$$



• Indices de position et de dispersion:

• Indices de position (ou de tendances centrales).

• Mode: noté M_0 .

pour une V.S. discrète le mode noté M_0 est la valeur x_i du caractère qui a le plus grand effectif partiel (ou bien la plus grande fréquence partielle).

Ex(2) famille le mode est $x_2 = 1 = M_0$.

pour une V.S continue la classe modale est celle qui a le plus grand effectif partiel (ou la plus grande fréquence partielle).

le mode M_0 est le centre de la classe modale. c

Rq: le mode n'est pas nécessairement situé au centre

On peut avoir plusieurs modes.

• Médiane: noté M_e

La médiane est la valeur de la variable statistique qui partage la population en deux.

Ex: $M_e = x_2 = 1$.

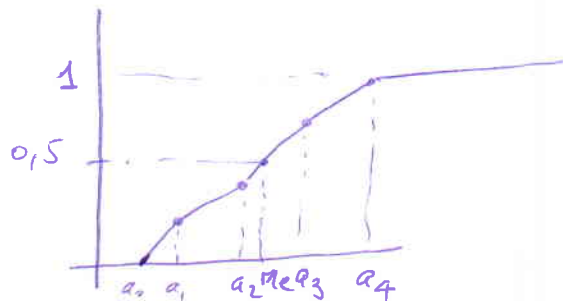
la médiane vérifie

$$F(M_e - 0) < \frac{1}{2} \leq F(M_e + 0).$$

si X est une V.S. continue

$$M_e \text{ vérifie } F(M_e) = \frac{1}{2} = 0,5$$

on retrouve M_e graphiquement



ou alors par interpolation:

Ait $C_k = [a_k, a_{k+1}[$ la première classe telle que la fréquence cumulée est $\geq 0,5$ (50%) alors

$$Me = a_k + \frac{0,5 - F_{k-1}}{F_k - F_{k-1}} (a_{k+1} - a_k)$$

F_k étant la fréquence cumulée.

$$Me = a_k + \frac{a_{k+1} - a_k}{f_k} (0,5 - F_{k-1})$$

b. Moyenne:

Ait $X: \Omega \rightarrow \mathbb{C}$ une v.s.

de loi : $\begin{array}{c|c|c|c} x_1 & x_2 & \dots & x_m \\ \hline f_1 & f_2 & & f_m \end{array}$ ou $\begin{array}{c|c|c|c} c_1 & c_2 & \dots & c_m \\ \hline f_1 & f_2 & & f_m \end{array}$

la moyenne de la variable statistique est donnée par:

$$\bar{X} = \sum_{i=1}^m x_i f_i = x_1 f_1 + x_2 f_2 + \dots + x_m f_m$$

$$= \sum_{i=1}^m x_i \frac{n_i}{N} = \frac{1}{N} \left[\sum_{i=1}^m x_i \cdot n_i \right] \text{ si } X \text{ est discrète.}$$

$$\text{et } \bar{X} = \sum_{i=1}^m f_i c_i = \frac{1}{N} \left[\sum_{i=1}^m n_i c_i \right] \text{ } c_i \text{ est le centre de } C_i$$

$$\underline{E}_X: \text{ familles: } \bar{X} = \sum_{i=1}^m x_i f_i = 0 \times 0,224 + 1 \times 0,340 + 2 \times 0,256 + 3 \times 0,120 + 4 \cdot 0,06 = 1,45$$

Propriétés:

$$1. \sum_{i=1}^m f_i (x_i - \bar{x}) = 0$$

$$2. \text{ si } y = ax + b \quad a, b \in \mathbb{R} \quad \text{alors } \bar{y} = a \bar{x} + b.$$

Preuve:

$$1. \sum_{i=1}^m f_i (x_i - \bar{x}) = \sum_{i=1}^m f_i x_i - \sum_{i=1}^m f_i \bar{x} = \bar{x} - \bar{x} \sum_{i=1}^m f_i = \bar{x} - \bar{x} = 0$$

$$2. y = ax + b, \bar{y} = \sum_{i=1}^m f_i y_i = \sum_{i=1}^m f_i (ax_i + b) = a \sum_{i=1}^m f_i x_i + b \sum_{i=1}^m f_i = a \bar{x} + b.$$

Indices de dispersion:

- 14 -

• Variance: la variable Statistique $X: \Omega \rightarrow \mathbb{C}$

ayant pour loi

$$\begin{array}{c} x_1 \quad x_2 \quad \dots \quad x_m \\ f_1 \quad f_2 \quad \dots \quad f_m \\ \hline c_1 \quad c_2 \quad \dots \quad c_m \\ f_1 \quad f_2 \quad \dots \quad f_m \end{array}$$

alors la variance de la variable X notée $\text{Var}(X)$ est donnée

par:

$$\begin{aligned} \text{Var}(X) &= \sum_{i=1}^m f_i (x_i - \bar{x})^2 \quad \text{ou} \quad \text{Var}(X) = \sum_{i=1}^m f_i (c_i - \bar{x})^2 \\ &= \sum_{i=1}^m f_i (x_i^2 - 2x_i\bar{x} + \bar{x}^2) \\ &= \sum_{i=1}^m f_i x_i^2 - 2\bar{x} \sum_{i=1}^m f_i x_i + \sum_{i=1}^m f_i \bar{x}^2 \\ &= \sum_{i=1}^m f_i x_i^2 - 2\bar{x}^2 + \bar{x}^2 \\ &= \sum_{i=1}^m f_i x_i^2 - \bar{x}^2 \end{aligned}$$

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

• Ecart type:

l'écart type d'une V.S est donnée par:

$$\sigma_x = \sqrt{\text{Var}(X)} = \sqrt{\sum_{i=1}^m f_i (x_i - \bar{x})^2} = \sqrt{\left(\sum_{i=1}^m f_i x_i^2 - \bar{x}^2 \right)}$$

Remarque:

σ_x donne une idée de la variation des valeurs de la V.S

autour de \bar{x} .

si σ_x est petit alors les valeurs de la V.S sont regroupées autour de \bar{x} .

si σ_x est grand alors les valeurs de la V.S sont dispersées autour de \bar{x} .

Proprietātes:

$$\text{Ja } y = ax + b, \quad a, b \in \mathbb{R}.$$

$$\text{Var}(y) = a^2 \text{Var}(x).$$

Preuve:

$$\begin{aligned} \text{Var}(y) &= \sum_{i=1}^m f_i y_i^2 - \bar{y}^2 = \sum_{i=1}^m f_i (ax_i + b)^2 - (a\bar{x} + b)^2 \\ &= \sum_{i=1}^m f_i (a^2 x_i^2 + b^2 + 2abx_i) - (a^2 \bar{x}^2 + b^2 + 2ab\bar{x}). \\ &= a^2 \sum_{i=1}^m f_i x_i^2 + b^2 \sum_{i=1}^m f_i + 2ab \sum_{i=1}^m f_i x_i - a^2 \bar{x}^2 - b^2 - 2ab\bar{x} \\ &= a^2 \sum_{i=1}^m f_i x_i^2 + b^2 + 2ab\bar{x} - a^2 \bar{x}^2 - b^2 - 2ab\bar{x} \\ &= a^2 \left(\sum_{i=1}^m f_i x_i^2 - \bar{x}^2 \right). \end{aligned}$$

à deux dimensions:

Définition: Soient X, Y deux Variables Statistiques définies sur la même population Ω ; le couple $Z = (X, Y)$ est appelé variable Statistique à deux dimensions. ($\text{Card}(\Omega) = N$ effectif total).

L'ensemble des valeurs prises par la V.S (X, Y)

est donné par $\{(X(\omega_i), Y(\omega_i)); i = 1, 2, \dots, N\}$.

ou bien par le tableau:

ω_i	ω_1	ω_2	ω_N
$X(\omega_i)$	$X(\omega_1)$	$X(\omega_2)$	$X(\omega_N)$
$Y(\omega_i)$	$Y(\omega_1)$	$Y(\omega_2)$	$Y(\omega_N)$

et on note $X(\omega_i) = x_i$
 $Y(\omega_i) = y_i$

Exemple: $\Omega =$ ensemble de 10 personnes.

$X(\omega) =$ taille de ω . (mètres)
 $Y(\omega) =$ poids (ω). (Kgs).

L'ensemble des valeurs de (X, Y) est donné par:

ω_i	ω_1	ω_2	ω_3	ω_4	ω_5	ω_6	ω_7	ω_8	ω_9	ω_{10}
$X(\omega_i)$	1,60	1,63	1,70	1,72	1,75	1,77	1,80	1,81	1,87	1,90
$Y(\omega_i)$	65	60	63	63	68	69	68	75	80	81

Représentation graphique en nuage de points:



Cette représentation s'appelle nuage de points (x_i, y_i) .

Rq: Cette représentation de limite aux cas où $\text{Card}(\Omega)$ fini.

2. Tableau de Contingence des effectifs de la V.S (X, Y):

Soient x_1, x_2, \dots, x_K les valeurs prises par la V.S X (ordonnées suivant les valeurs croissantes).

idem pour y_1, y_2, \dots, y_M les valeurs prises par la V.S Y.

Posons $n_{ij} = \text{Card} \{ \omega \in \Omega ; (X(\omega), Y(\omega)) = (x_i, y_j) \}$
= effectif de la valeur (x_i, y_j) .

le tableau de contingence des effectifs de la V.S (X, Y)

est donné par le tableau suivant:

$x_i \backslash y_j$	y_1	y_2	\dots	y_j	\dots	y_M
x_1	n_{11}	n_{12}	\dots	n_{1j}	\dots	n_{1M}
x_2	n_{21}	n_{22}	\dots	n_{2j}	\dots	n_{2M}
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_i	\vdots	\vdots	\vdots	n_{ij}	\vdots	\vdots
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_K	n_{K1}	\dots	\dots	\dots	\dots	n_{KM}

Exemple:

Ω = ensemble des familles d'une cité.

ω = une famille.

$X(\omega)$ = nbre d'enfants filles.

$Y(\omega)$ = nbre d'enfants garçons

le tableau de contingence des effectifs de la variable (X, Y)

est donnée par:

$x_i \backslash y_j$	y_1 0	y_2 1	y_3 2	y_4 3	y_5 4	
$x_1 = 1$	4	4	2	0	0	$\rightarrow n_{1.} = 10$
$x_2 = 2$	9	16	4	0	0	$\rightarrow n_{2.} = 29$
$x_3 = 3$	4	12	9	2	0	$\rightarrow n_{3.} = 27$
$x_4 = 4$	1	6	1	1	2	$\rightarrow n_{4.} = 11$
$x_5 = 5$	0	1	0	1	1	$\rightarrow n_{5.} = 3$
	$n_{.1} = 18$	$n_{.2} = 39$	$n_{.3} = 16$	$n_{.4} = 4$	$n_{.5} = 3$	

Lecture du tableau:

par exemple: il y'a 12 familles qui ont: $\underbrace{3 \text{ filles}}_{x_3}$ et $\underbrace{1 \text{ garçon}}_{y_2}$
 $n_{32} = 12$.

On note $n_{i.} = \sum_{j=1}^M n_{ij}$; $n_{i.}$ = effectif de la valeur x_i

$$n_{1.} = n_{11} + n_{12} + n_{13} + n_{14} + n_{15} = 4 + 4 + 2 + 0 + 0 = 10$$

$$n_{2.} = 9 + 16 + 4 + 0 + 0 = 29 ; \quad n_{3.} = 4 + 12 + 9 + 2 + 0 = 27$$

$$n_{4.} = 1 + 6 + 1 + 1 + 2 = 11 ; \quad n_{5.} = 0 + 1 + 0 + 1 + 1 = 3$$

$n_{3.} = 27$ = effectif de la valeur x_3
 = 27 familles ont 3 filles

de la même manière: $n_{.j} = \sum_{i=1}^K n_{ij}$ = $n_{.j}$ = effectif de la valeur y_j

$$n_{.1} = 4 + 9 + 4 + 1 + 0 = 18 , \quad n_{.2} = 4 + 16 + 12 + 6 + 1 = 39 , \quad n_{.3} = 2 + 4 + 9 + 1 + 0 = 16$$

$$n_{.4} = 0 + 0 + 2 + 1 + 1 = 4 , \quad n_{.5} = 0 + 0 + 0 + 2 + 1 = 3$$

$n_{.2} = 39$ = effectif de la valeur y_2
 = nombre de familles qui ont 1 garçon

Propriétés:

$$\sum_{i=1}^K n_{i.} = \sum_{j=1}^M n_{.j} = N \text{ effectif total.}$$

Exemple:

$$\sum n_{i.} = 80 \quad \text{80 familles dans la cité.}$$

$$\sum n_{.j} = 80$$

• Loi de la variable Statistique (X, Y) :

le nombre $f_{ij} = \frac{n_{ij}}{N}$ est appelé fréquence de la valeur (x_i, y_j) .

la loi de la variable Statistique (X, Y) notée $\mathcal{L}(X, Y)$ est donnée par le tableau:

$x_i \backslash y_j$	y_1	y_2	...
x_1	f_{11}	f_{12}	
...			

Remarque:

on note $f_{i.} = \sum_{j=1}^M f_{ij}$ et $f_{.j} = \sum_{i=1}^K f_{ij}$

alors $(f_{i.} = \frac{n_{i.}}{N} \text{ et } f_{.j} = \frac{n_{.j}}{N})$

$f_{i.}$ = la fréquence de la valeur x_i

$f_{.j}$ = la fréquence de la valeur y_j

Propriétés

$$\sum_{i=1}^K f_{i\cdot} = \sum_{j=1}^M f_{\cdot j} = 1$$

en effet:

$$\sum_{i=1}^K f_{i\cdot} = \sum_{i=1}^K \sum_{j=1}^M f_{ij} = \sum_{i=1}^K \sum_{j=1}^M \frac{n_{ij}}{N} = \sum_{i=1}^K \frac{n_{i\cdot}}{N} = \frac{N}{N} = 1$$

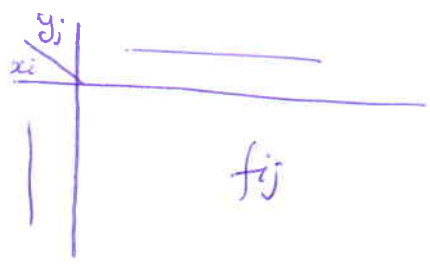
• Séries Marginales:

Soit (X, Y) une V.S à 2 dimensions.

La variable statistique X est appelée (série) V.S marginale.

_____ u _____ u Y _____ u _____ u _____

Par la suite on supposera que la loi du couple (X, Y) est donnée



alors:

• la loi marginale de la V.S X notée $\mathcal{L}(X)$ est donnée par:

x_i	x_1	x_2	\dots	x_K
$f_{i\cdot}$	$f_{1\cdot}$	$f_{2\cdot}$		$f_{K\cdot}$

• la loi marginale de la V.S Y notée $\mathcal{L}(Y)$ est donnée par:

y_j	y_1	y_2	\dots	y_M
$f_{\cdot j}$	$f_{\cdot 1}$	$f_{\cdot 2}$		$f_{\cdot M}$

Exemple: si l'on reprend l'exemple des familles.

La loi marginale de la v.s. X est donnée par:

x_i	1	2	3	4	5
$f_{i.} = \frac{n_{i.}}{N}$	0,125	0,3625	0,3375	0,1375	0,0375

y_j	0	1	2	3	4
$f_{.j}$	0,225	0,4875	0,2	0,05	0,0375

• Moyennes marginales:

1/ $\bar{x} = \sum_{i=1}^K f_{i.} x_i = \frac{1}{N} \sum_{i=1}^K x_i n_{i.}$ s'appelle moyenne marginale de la v.s. X . ($E_x = 2,6 = \bar{x}$)

2/ $\bar{y} = \sum_{j=1}^M f_{.j} y_j = \frac{1}{N} \sum_{j=1}^M n_{.j} y_j$ s'appelle la moyenne marginale de la v.s. Y . ($N = \sum n_{i.} = \sum n_{.j}$).

($E_x(\bar{y}) = 1,1875$).

• Covariance du couple (X, Y)

La covariance du couple (X, Y) notée $Cov(X, Y)$ est donnée par:

$$Cov(X, Y) = \sum_{i=1}^K \sum_{j=1}^M f_{ij} (x_i - \bar{x})(y_j - \bar{y})$$

où \bar{x} et \bar{y} sont respectivement les moyennes de x et y .

Remarque:

• $Cov(X, Y) = \frac{1}{N} \sum_{i=1}^K \sum_{j=1}^M n_{ij} (x_i - \bar{x})(y_j - \bar{y})$ ($N = \sum n_{i.} = \sum n_{.j} = \text{card}(\Omega)$).

• $Cov(X, Y) = Cov(Y, X)$.

• $Cov(X, Y) = \sum_{i=1}^K \sum_{j=1}^M f_{ij} (x_i y_j) - (\bar{x} \bar{y}) = \frac{1}{N} \sum_{i=1}^K \sum_{j=1}^M n_{ij} x_i y_j - (\bar{x} \bar{y})$.

Remarques:

$$1) \text{ si } X=Y \quad \text{Cov}(X, X) = \sum_{i=1}^k f_i x_i^2 - \bar{x}^2 = \text{Var}(X).$$

la notion de covariance généralise la notion de variance

2) La covariance de (X, Y) mesure le degré de dépendance des V.S X et Y .

Lois conditionnelles:

La distribution ou la loi conditionnelle de X sachant que $Y=Y_j$ notée $\mathcal{L}(X / Y=Y_j)$. (loi de X sachant que $Y=Y_j$); c'est la loi de la variable statistique qui prend les valeurs x_1, x_2, \dots, x_k avec les fréquences $\frac{f_{1j}}{f_{.j}}, \frac{f_{2j}}{f_{.j}}, \dots, \frac{f_{kj}}{f_{.j}}$

c'est donc la loi :

x_i	x_1	x_2	...	x_k
$\frac{f_{1j}}{f_{.j}}$	$\frac{f_{2j}}{f_{.j}}$	$\frac{f_{3j}}{f_{.j}}$		$\frac{f_{kj}}{f_{.j}}$

et la valeur: $\bar{x}_j = \sum_{i=1}^k \frac{f_{ij}}{f_{.j}} x_i =$ moyenne conditionnelle sachant que $Y=Y_j$.

de la même manière la distribution ou la loi de Y sachant que $X=x_i$ notée $\mathcal{L}(Y / X=x_i)$ est la loi de la V.S qui prend les valeurs y_1, y_2, \dots, y_m avec les fréquences $\frac{f_{i1}}{f_{i.}}, \frac{f_{i2}}{f_{i.}}, \dots, \frac{f_{im}}{f_{i.}}$

c'est donc la loi:

y_j	y_1	y_2	...	y_m
$\frac{f_{i1}}{f_{i.}}$	$\frac{f_{i2}}{f_{i.}}$	$\frac{f_{i3}}{f_{i.}}$		$\frac{f_{im}}{f_{i.}}$

et de la même manière la valeur :

$$\bar{y}_i = \sum_{j=1}^m \frac{f_{ij}}{f_{i.}} y_j \quad \text{est la moyenne conditionnelle}$$

sachant que $x = x_i$

Exemple récapitulatif :

le tableau de contingence d'une V.S à deux dimensions est donné par :

$x_i \setminus y_j$	$y_1 = -1$	$y_2 = 0$	$y_3 = 1$	
$x_1 = 1$	18	4	18	$\rightarrow n_{1.} = 40$
$x_2 = 2$	20	60	20	$\rightarrow n_{2.} = 100$
$x_3 = 3$	2	16	42	$\rightarrow n_{3.} = 60$
	$n_{.1} = 40$	$n_{.2} = 80$	$n_{.3} = 80$	$\Rightarrow N = 200$

1- Donner la loi de (x, y) , $\mathcal{L}(x, y)$, (les lois marginales de x et y les moyennes — " — " —)
 les lois conditionnelles $\mathcal{L}(x | y = y_j)$ les moyennes conditionnelles $\bar{x}_1, \bar{x}_2, \bar{x}_3$

$x_i \setminus y_j$	$y_1 = -1$	$y_2 = 0$	$y_3 = 1$	
$x_1 = 1$	0,09	0,02	0,09	$f_{1.} = 0,2$
$x_2 = 2$	0,1	0,3	0,1	$f_{2.} = 0,5$
$x_3 = 3$	0,01	0,08	0,21	$f_{3.} = 0,3$
	$f_{.1} = 0,2$	$f_{.2} = 0,4$	$f_{.3} = 0,4$	$\begin{matrix} 1 \\ 1 \\ = 1 \end{matrix}$

2- Donner la loi marginale de x .

x_i	$x_1 = 1$	$x_2 = 2$	$x_3 = 3$
$f_{i.}$	$0,2 = f_1$	$f_2 = 0,5$	$f_3 = 0,3$

• moyenne marginale de x .

$$\bar{x} = \sum_{i=1}^3 f_{i\cdot} x_i = 0,2 \times 1 + 0,5 \times 2 + 0,3 \times 3 = 2,1$$

$$\bar{x} = 2,1.$$

• loi marginale de y

y_j	$y_1 = -1$	$y_2 = 0$	$y_3 = +1$
$f_{\cdot j}$	$f_{\cdot 1} = 0,2$	$f_{\cdot 2} = 0,4$	$f_{\cdot 3} = 0,4$

• moyenne marginale de y

$$\bar{y} = \sum_{j=1}^3 f_{\cdot j} y_j = (-1) \cdot 0,2 + 0 \cdot 0,4 + 1 \cdot 0,4 = 0,2$$

$$\bar{y} = 0,2.$$

$$\bullet \text{Cov}(x, y) = \sum_{i=1}^k \sum_{j=1}^{kn} f_{ij} (x_i y_j) - (\bar{x} \bar{y})$$

$$= \sum_{i=1}^3 \sum_{j=1}^3 f_{ij} (x_i y_j) - (\bar{x} \bar{y}).$$

$$= f_{11} (x_1 y_1) + f_{12} (x_1 y_2) + f_{13} (x_1 y_3) + f_{21} (x_2 y_1) + f_{22} (x_2 y_2) + f_{23} (x_2 y_3) + f_{31} (x_3 y_1) + f_{32} (x_3 y_2) + f_{33} (x_3 y_3).$$

$$= 0,09(-1)1 + 0,02 \times 0 \times 1 + 0,09 \times 1 \times 1 + 0,1 \times 2 \times (-1) + 0,3 \cdot 0 \times 2 + 0,1 \times 1 \times 2 + 0,01 \times 3(-1) + 0,08 \times 3 \cdot 0 + 0,21 \times 3 \times 1 - 2,1 \times 0,2$$

$$= 0,18 - 0,42.$$

$$= -0,24$$

• Donner les lois conditionnelles: $\mathcal{L}(X|Y=1=y_1)$, $\mathcal{L}(X|Y=0=y_2)$ -25-

$\mathcal{L}(X|Y=1=y_2)$. et les moyennes conditionnelles correspondantes:

• $\mathcal{L}(X|Y=y_1=-1)$.

x_i	1	2	3
f_{i1}	$\frac{0,09}{0,2} = 0,45$	$\frac{0,1}{0,2} = 0,5$	$\frac{0,01}{0,2} = 0,05$

$\bar{x}_1 = 1 \times 0,45 + 2 \times 0,5 + 3 \times 0,05 = 1,6$ (moyenne conditionnelle sachant que $Y=y_1$)

• $\mathcal{L}(X|Y=y_2=0)$.

x_i	1	2	3
f_{i2}	$\frac{0,02}{0,4} = 0,05$	$\frac{0,3}{0,4} = 0,75$	$\frac{0,08}{0,4} = 0,2$

$\bar{x}_2 = 1 \times 0,05 + 2 \times 0,75 + 3 \times 0,2 = 2,15$. (moyenne conditionnelle sachant que $Y=y_2$)

• $\mathcal{L}(X|Y=y_3=+1)$.

x_i	1	2	3
f_{i3}	$\frac{0,09}{0,4} = 0,225$	$\frac{0,1}{0,4} = 0,25$	$\frac{0,21}{0,4} = 0,525$

$\bar{x}_3 = 1 \times 0,225 + 2 \times 0,25 + 3 \times 0,525 = 2,3$. (moyenne conditionnelle sachant que $Y=y_3$).

• Calculer: $\bar{x}_1 \cdot f_{.1} + \bar{x}_2 \cdot f_{.2} + \bar{x}_3 \cdot f_{.3} = 1,6 \cdot 0,2 + 2,15 \cdot 0,4 + 2,3 \cdot 0,4 = 2,1 = \bar{x}$.



Proposition:

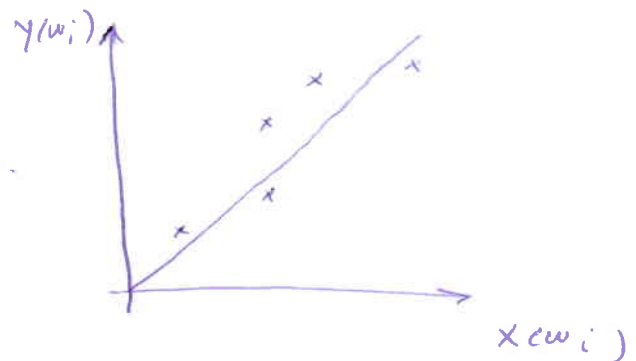
$\bar{x} = \sum_{j=1}^m \bar{x}_j \cdot f_{.j}$ et $\bar{y} = \sum_{i=1}^k \bar{y}_i \cdot f_{i.}$

\bar{y}_i, \bar{x}_j étant les moyennes conditionnelles et \bar{x}, \bar{y} les moyennes marginales.

Corrélation et droite de régression:

Il arrive que dans une étude statistique qu'on mesure sur chaque individu w 2 caractères (2 V.S); et par la suite on examine s'il existe une dépendance (ou corrélation) entre les 2 V.S

(Cela peut être remarqué après représentation en nuage de points).



Rappel:

$$\text{Cov}(x, y) = \sum_{i=1}^k \sum_{j=1}^m f_{ij} x_i y_j - \bar{x} \bar{y}$$

$$\bar{x} = \sum_{i=1}^k f_{i.} x_i, \quad \bar{y} = \sum_{j=1}^m f_{.j} y_j$$

$$\text{Var}(x) = \sum_{i=1}^k f_{i.} x_i^2 - \bar{x}^2, \quad \text{Var}(y) = \sum_{j=1}^m f_{.j} y_j^2 - (\bar{y})^2$$

$$\sigma_x = \sqrt{\text{Var}(x)}; \quad \sigma_y = \sqrt{\text{Var}(y)}$$

Définition: On appelle coefficient de corrélation entre les V.S x et y

le nombre ρ défini par : $\rho = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$

Propriétés:

- le coefficient ρ mesure le degré de liaison linéaire entre x et y .

- $|\rho| \leq 1$ (i.e. $-1 \leq \rho \leq 1$).

- si $\rho = 1$ ou $\rho = -1$ x et y sont liés par une relation linéaire
- si $\rho \approx +1$ ou $\rho \approx -1$ tous les points x_i, y_j se trouvent près d'une même droite
- si $\rho = 0$ il y'a indépendance entre x et y .

Droite de régression:

Nous cherchons à déterminer une droite d'équation:

$y = ax + b$ qui puisse approcher le nuage de points dans

le sens où $Q(\alpha, \beta) = \frac{1}{N} \sum_{i=1}^N (y(w_i) - \alpha x(w_i) - \beta)^2$ est minimale

$Q(\alpha, \beta)$ représente la moyenne de la somme des carrés de la distance des points $(x(w_i), y(w_i))$ et $(x(w_i), \alpha x(w_i) + \beta)$.

on cherche donc à minimiser la quantité $Q(\alpha, \beta)$

et on aura $Q(a, b) = \min Q(\alpha, \beta)$.

Définition: a et b sont tels que:

$a = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$ et $b = \bar{y} - a\bar{x} = \bar{y} - \frac{\text{Cov}(x, y)}{\text{Var}(x)} \bar{x}$.

la droite de régression de y en x est donnée par

$y = ax + b = \frac{\text{Cov}(x, y)}{\text{Var}(x)} \cdot x + \left(\bar{y} - \frac{\text{Cov}(x, y)}{\text{Var}(x)} \bar{x} \right)$