

# STATISTIQUE DESCRIPTIVE

## I Historique

- Questions de dénombrement (problèmes d'inventaire, de recensement) remontant à la Préhistoire.
- Plus récemment :
  - \* Probabilités : 17<sup>e</sup> siècle, Pascal et Fermat  
puis 19<sup>e</sup> siècle, famille Bernoulli, Gauss, Laplace.  
début 20<sup>e</sup> siècle : axiomatisation; Kolmogorov ...
  - \* Statistiques : 19<sup>e</sup> siècle, Quételet (recensement de Bruxelles qui fut un exemple à suivre) (Adolphe Quételet 1796-1874)  
puis 20<sup>e</sup> siècle, grand essor avec Pearson, Gosset (Student), Neyman, Fisher (de 1850 à 1965).

## II Traitement des données

Présentation : tableaux IC (individus x caractères)

$i$  : individus       $n$  individus  
 $j$  : variables       $p$  variables

$i \setminus j$	1	...	$j$	...	$p$
1	$x_{11}$	...	$x_{1j}$	...	$x_{1p}$
...					
$i$	$x_{i1}$	...	$x_{ij}$	...	$x_{ip}$
...					
$n$	$x_{n1}$	...	$x_{nj}$	...	$x_{np}$

$x_{ij}$  : caractères  
 $n$  : effectif total

cas univarié :  $\begin{array}{c|cccc} i & 1 & \dots & i & \dots & n \\ \hline x_i & x_1 & \dots & x_i & \dots & x_n \end{array}$ , cas bivarié :  $\begin{array}{c|cccc} i & 1 & \dots & i & \dots & n \\ \hline x_i & x_{i1} & \dots & x_{ij} & \dots & x_{in} \\ y_i & y_1 & \dots & y_i & \dots & y_n \end{array}$ , etc...

Organisation : série observée - série ordonnée

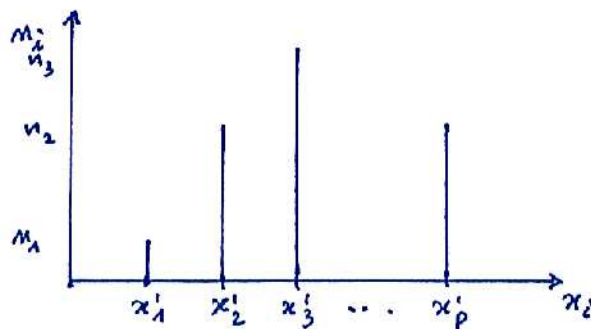
observations :  $\{x_1, \dots, x_n\}$

on ordonne :  $\{x_{(1)}, \dots, x_{(n)}\}$  avec  $i \leq j \Rightarrow x_{(i)} \leq x_{(j)}$ .

ordre strict :  $\{x'_1, \dots, x'_p\}$  avec  $i < j \Rightarrow x'_i < x'_j$   
et nombre d'observations identiques  $x'_i$  :  $m_i$  (effectif de  $x'_i$ )

$\{(x'_i, m_i), 1 \leq i \leq p\}$  est appelé série statistique.  $\sum_{i=1}^p m_i = n$ .

diagramme en bâtons :



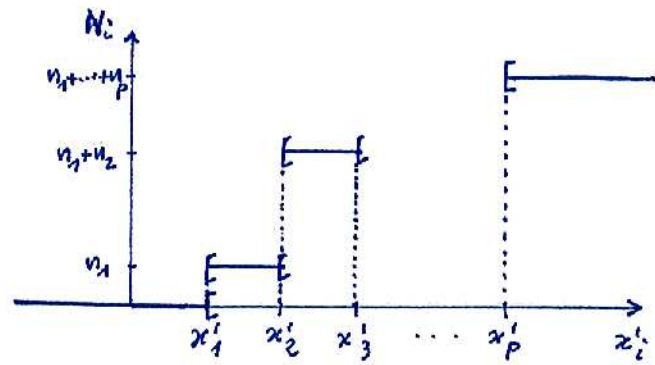
fréquences, effectifs cumulés, fréquences cumulées

$$f_i = \frac{n_i}{n}, \quad \sum_{i=1}^p f_i = 1$$

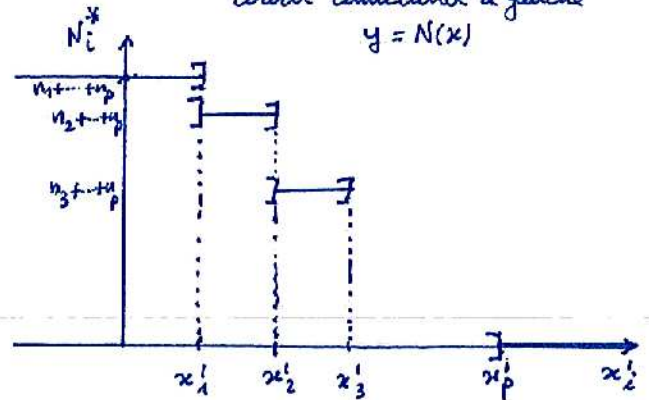
$$\left\{ \begin{array}{l} \text{effectifs cumulés à gauche } N_i = \sum_{k=1}^i n_k \\ \text{fréquences cumulées à gauche } F_i = \frac{N_i}{n} \end{array} \right.$$

$$\left\{ \begin{array}{l} \text{effectifs cumulés à droite } N_i^* = \sum_{k=i}^p n_k \\ \text{fréquences cumulées à droite } F_i^* = \frac{N_i^*}{n} \end{array} \right.$$

$$\text{On a } N(x) + N^*(x) = \begin{cases} n & \text{si } x \notin \{x'_1, \dots, x'_p\} \\ n + n_i & \text{si } x = x'_i \end{cases} \geq n$$



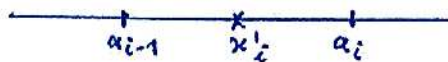
courbe cumulative à gauche  
 $y = N(x)$



courbe cumulative à droite  
 $y = N^*(x)$

regroupement des données en classes

Si on dispose de trop de données, le diagramme en bâtons est surchargé.  
On regroupe alors les données en classes  $[a_0, a_1[$ ,  $[a_1, a_2[$ , ...,  $[a_{p-1}, a_p]$ .



centre de la classe :  $x'_i$   
longueur de la classe :  $l_i = a_i - a_{i-1}$   
effectif de la classe :  $n_i$

histogramme - polygone des effectifs (ou fréquences)

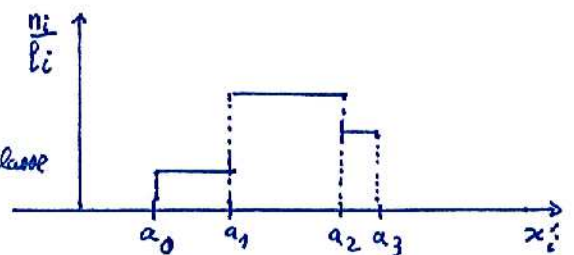
• histogramme des effectifs (ou fréquences)

$\frac{n_i}{l_i}$  = densité d'effectif = effectif par unité de longueur de classe

(idem avec les fréquences,  $\frac{f_i}{l_i}$ )

aire de chaque rectangle :  $\frac{n_i}{l_i} \times l_i = n_i =$  effectif de la classe  $[a_{i-1}, a_i[$ .

aire sous l'histogramme :  $n =$  effectif total.



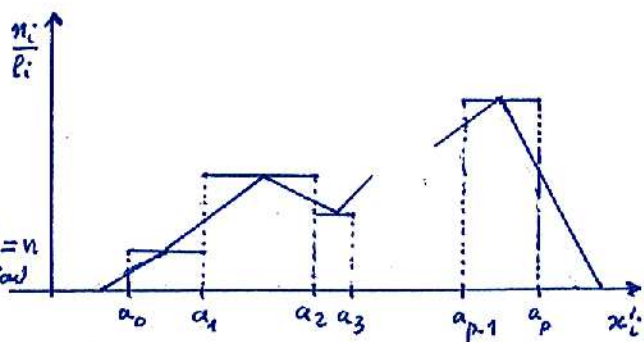
- "image plus continue" de l'histogramme

polygone des effectifs = interpolation linéaire entre les centres des classes.

On suppose l'uniformité des répartitions dans chaque classe.

L'aire sous la courbe polygonale est inchangée = n (les 2 bouts extrêmes sont placés pour vérifier cette condition)

Lorsque les  $h_i$  deviennent petits, on obtient un "graphe de densité locale".

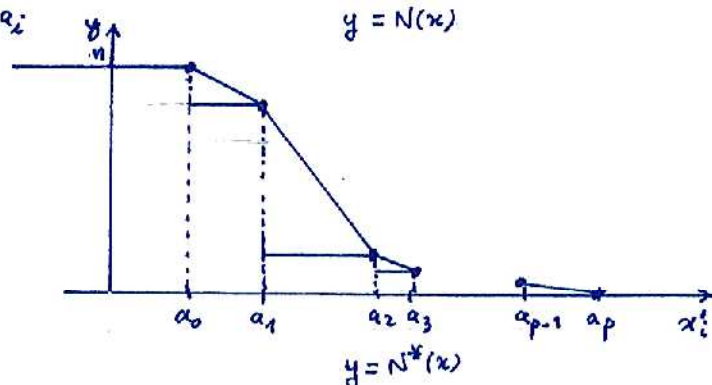
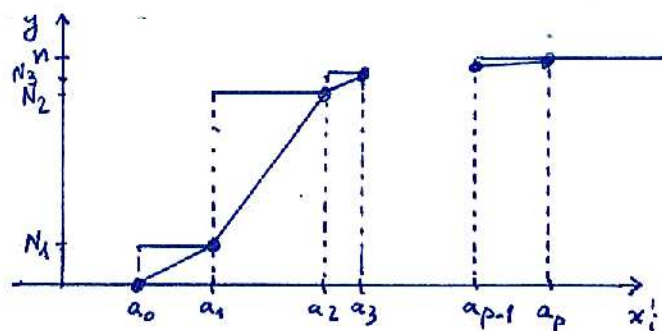


- courbes cumulées

courbe cumulée à gauche = ligne brisée joignant les points  $(a_i, N_i)$

$$N(x) = \begin{cases} 0 & \text{si } x < a_0 \\ \frac{n_1}{l_1}(x-a_0) & \text{si } a_0 \leq x < a_1 \\ \vdots \\ N_{i-1} + \frac{n_i}{l_i}(x-a_{i-1}) & \text{si } a_{i-1} \leq x < a_i \\ \vdots \\ n & \text{si } x \geq a_p \end{cases}$$

On a  $\forall x, N(x) + N^*(x) = n.$



### III Paramètres de position

moyenne (pour les effectifs quantitatifs):  $\bar{x} = \frac{1}{n} \sum_{i=1}^p x_i = \frac{1}{n} \sum_{i=1}^p m_i x_i^j = \sum_{i=1}^p f_i x_i^j$

Pour les observations groupées:  $\bar{x} = \frac{1}{n} \sum_{i=1}^p m_i x_i^j$ ,  $x_i^j$ : centre de la classe  $[a_{i-1}, a_i]$ .

C'est le paramètre le plus répandu, mais pas nécessairement le plus adéquat.

Un inconvénient: e.g. s'il y a une donnée aberrante, elle affecte sensiblement la moyenne. On peut donc être amené à tronquer la moyenne, ou à la pondérer.

$$\frac{1}{p-2} \sum_{i=2}^{p-1} m_i x_i^j, \quad \frac{\sum_{i=1}^p w_i x_i^j}{\sum_{i=1}^p w_i} \text{ etc. ...}$$

Si on a deux séries statistiques  $\{x_i, 1 \leq i \leq n_x\}$  et  $\{y_j, 1 \leq j \leq n_y\}$ , la série complète  $\{z_k, 1 \leq k \leq n_x + n_y\} = \{x_i, 1 \leq i \leq n_x\} \cup \{y_j, 1 \leq j \leq n_y\}$  a pour moyenne

$$\bar{z} = \frac{n_x \bar{x} + n_y \bar{y}}{n_x + n_y}$$

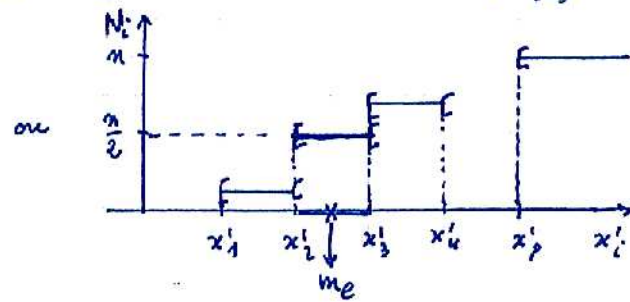
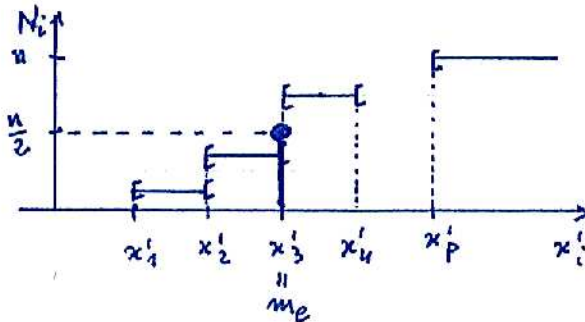
médiane :  $m_e$  est la valeur "centrale" de la série ordonnée  $x_{(1)}, x_{(2)}, \dots, x_{(n)}$  (avec répétition)

• si  $n$  est impair  $m_e = x_{(\frac{n+1}{2})}$

• si  $n$  est pair, il y a un intervalle médian  $[x_{(\frac{n}{2})}, x_{(\frac{n}{2}+1)}]$ . On pose  $m_e = \frac{x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}}{2}$

ex : série ordonnée 1 2 2 3 4 4 4 5 6  $\rightarrow n=9$   $m_e = x_{(5)} = 4$   
 1 2 2 3 4 4 5  $\rightarrow n=8$   $m_e = \frac{x_{(4)} + x_{(5)}}{2} = 3,5$

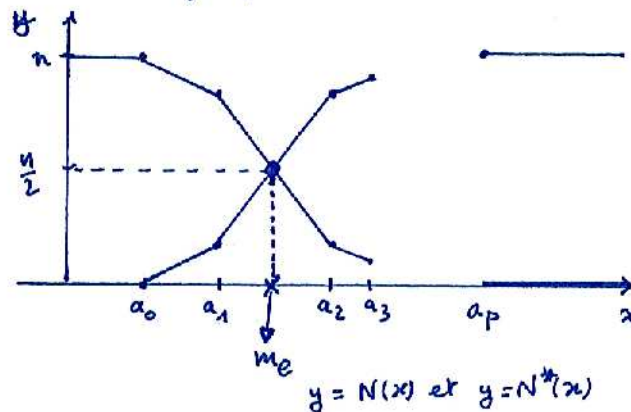
Détermination de  $m_e$  à l'aide des courbes cumulatives : on a  $N(m_e) \geq \frac{n}{2}$  et  $N^*(m_e) \geq \frac{n}{2}$  (i.e. le nombre de valeurs de la série  $< m_e$  est égal au nombre de valeurs  $> m_e$  :  $\lfloor \frac{n}{2} \rfloor$ )



Si il existe un  $i$  tel que  $N_{i-1} < \frac{n}{2} < N_i$  alors  $m_e = x'_i$

Si il existe un  $i$  tel que  $N_i = \frac{n}{2}$  alors  $[x'_i, x'_{i+1}]$  est l'intervalle médian et  $m_e = \frac{x'_i + x'_{i+1}}{2}$ .

Cas des observations groupées : on utilise les courbes cumulatives continues :



$$N + N^* = n$$

$$N(m_e) = N^*(m_e) = \frac{n}{2}$$

$$m_e = a_{i-1} + \frac{h_i}{m_i} \left( \frac{n}{2} - N_{i-1} \right)$$

$i$  étant l'indice de la classe contenant  $m_e$ .

quantiles (ou fractiles) : Soit  $p \in ]0, 1[$ .  $Q_p$  :  $p$ -quantile (ou quantile d'ordre  $p$ )

$Q_p$  est la valeur du caractère telle que  $N(Q_p) \geq np$  et  $N^*(Q_p) \geq n(1-p)$  (i.e. le nombre de valeurs de la série  $< Q_p$  est égal à  $\lfloor np \rfloor$ , et  $> Q_p$  :  $\lfloor n(1-p) \rfloor$ )

Si  $np \notin \mathbb{N}$   $Q_p = x_{(\lfloor np \rfloor + 1)}$  ; si  $np \in \mathbb{N}$  donc  $Q_p \in [x_{(np)}, x_{(np+1)}]$

ex : médiane :  $p = \frac{1}{2}$ ,  $m_e = Q_{1/2}$

cas d'observations groupées :

quantiles :  $p = \frac{1}{4}, \frac{2}{4}, \frac{3}{4}$   $Q_{1/4}, Q_{2/4} = m_e, Q_{3/4}$

$$Q_p = a_{i-1} + \frac{h_i}{m_i} (np - N_{i-1})$$

déciles :  $p = \frac{1}{10}, \frac{2}{10}, \dots, \frac{9}{10}$

percentiles :  $p = \frac{1}{100}, \frac{2}{100}, \dots, \frac{99}{100}$

mode : valeur observée qui a le plus grand effectif  
 distribution unimodale (plurimodale) : qui contient un seul (plusieurs) mode(s).

Cas des observations groupées : si les classes ont même longueur, la classe modale est celle d'effectif maximal. Sinon, considérer l'histogramme des densités  $\frac{n_i}{l_i}$ .

autres paramètres de position :  $c_1 = \frac{x_{(1)} + x_{(n)}}{2}$ ,  $c_2 = \frac{Q_{1/4} + Q_{3/4}}{2}$ ,  $c_3 = \frac{Q_{1/4} + 2Q_{2/4} + Q_{3/4}}{4}$  etc...

remarque : tous les paramètres considérés ici ont même dimension que les caractères.

#### IV Paramètres de dispersion

étendue (empan) :  $x_{(n)} - x_{(1)}$ . Inconvénient : l'étendue ne tient pas compte de toutes les observations.

ex : 

ces trois séries ont même étendue.

écart interquartile :  $[Q_p, Q_{1-p}]$  pour  $0 < p < \frac{1}{2}$ . Cet intervalle contient un effectif égal ou juste supérieur à  $n(1-2p)$  :  $\begin{cases} n(1-2p) & \text{si } np \in \mathbb{N} \\ [n(1-2p)] + 1 & \text{si } np \notin \mathbb{N} \end{cases}$ .

ex : intervalle interquartile  $[Q_{1/4}, Q_{3/4}]$  : contient  $\approx 50\%$  d'observations  
 intervalle interdécile  $[Q_{1/10}, Q_{9/10}]$  : contient  $\approx 80\%$  d'observations.

ces intervalles sont utilisés pour éliminer les valeurs extrêmes (éventuellement) aberrantes.

écart moyen absolu :  $e_m = \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}| = \frac{1}{n} \sum_{i=1}^p n_i |x_i' - \bar{x}| = \sum_{i=1}^p f_i |x_i' - \bar{x}|$

Pour les observations groupées :  $e_m = \frac{1}{n} \sum_{i=1}^p n_i |x_i' - \bar{x}|$ ,  $x_i'$  : centre de classe

écart médian absolu :  $e_m^* = \frac{1}{n} \sum_{i=1}^n |x_i - m_e|$ .  $m_e$  minimise  $m \mapsto \sum_{i=1}^n |x_i - m|$ .

Inconvénient : le calcul de  $e_m$  et  $e_m^*$  est peu maniable à cause de la valeur absolue, et ils ont peu de propriétés mathématiques.

écart-type et variance :  $\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^p n_i (x_i' - \bar{x})^2 = \sum_{i=1}^p f_i (x_i' - \bar{x})^2$

Théorème de König-Huygens :  $\frac{1}{n} \sum_{i=1}^n (x_i - m)^2 = \sigma^2 + (\bar{x} - m)^2$

Donc  $\sigma^2$  minimise  $m \mapsto \sum_{i=1}^n (x_i - m)^2$  et  $\sigma^2 = \overline{x^2} - \bar{x}^2$ .

autre écriture :  $\sigma^2 = \frac{1}{2n^2} \sum_{1 \leq i, j \leq n} (x_i - x_j)^2$  ( $\rightarrow$  évite de calculer tous les écarts  $x_i - \bar{x}$ )

Valeurs centrées réduites :  $y_i = \frac{x_i - \bar{x}}{\sigma}$ ,  $\bar{y} = 0$ ,  $\sigma_y = 1$  ( $y$  est sans dimension)

si on réunit deux séries statistiques  $\{x_i, 1 \leq i \leq n_x\}$  et  $\{y_j, 1 \leq j \leq n_y\}$  en une seule  $\{z_k, 1 \leq k \leq n_x + n_y\}$ , les écarts-typé sont reliés par

$$\sigma_z^2 = \frac{n_x \sigma_x^2 + n_y \sigma_y^2}{n_x + n_y} + \frac{n_x (\bar{x} - \bar{z})^2 + n_y (\bar{y} - \bar{z})^2}{n_x + n_y}$$

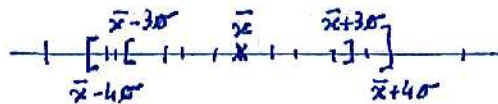
variance DANS les groupes =  $\frac{n_x n_y}{n_x + n_y} (\bar{x} - \bar{y})^2$ , variance ENTRE les groupes  
 démonstration: utiliser König:  $\sigma_z^2 = \frac{1}{n_x + n_y} \left[ \sum_{i=1}^{n_x+n_y} (z_i - \bar{z})^2 \right] = \frac{n_x}{n_x+n_y} [\sigma_x^2 + (\bar{x} - \bar{z})^2] + \frac{n_y}{n_x+n_y} [\sigma_y^2 + (\bar{y} - \bar{z})^2]$

Inégalité de Bienaymé - Tchebycheff:  $\left\{ \begin{array}{l} \text{la fréquence des } x_i \text{ tels que } |x_i - \bar{x}| \geq k \text{ est } \leq \frac{\sigma^2}{k^2}, \\ \text{la fréquence des } x_i \text{ tels que } |x_i - \bar{x}| < k \text{ est } \geq 1 - \frac{\sigma^2}{k^2}. \end{array} \right.$

démonstration:  $\sigma^2 = \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|^2 = \frac{1}{n} \sum_{i: |x_i - \bar{x}| \geq k} |x_i - \bar{x}|^2 + \frac{1}{n} \sum_{i: |x_i - \bar{x}| < k} |x_i - \bar{x}|^2$   
 $\geq \frac{k^2}{n} \text{card} \{ i : |x_i - \bar{x}| \geq k \}$

Puis  $\frac{1}{n} \text{card} \{ i : |x_i - \bar{x}| < k \} = 1 - \frac{1}{n} \text{card} \{ i : |x_i - \bar{x}| \geq k \} \geq 1 - \frac{\sigma^2}{k^2} . \square$

application:  $k = 3\sigma$ : la fréquence des  $x_i$  tels que  $|x_i - \bar{x}| \geq 3\sigma$  est  $\leq \frac{1}{9} \approx 11,2\%$   
 $k = 4\sigma$ : la fréquence des  $x_i$  tels que  $|x_i - \bar{x}| \geq 4\sigma$  est  $\leq \frac{1}{16} \approx 6,3\%$ .



Autres paramètres de dispersion:

• coefficient de variation:  $CV = \frac{\sigma}{\bar{x}}$ : nombre sans dimension, se mesure en %  
 (mesure relative de dispersion)

• paramètres de forme:  $m_3 = \overline{(x - \bar{x})^3} = \frac{1}{n} \sum_{i=1}^n m_i (x_i - \bar{x})^3$  moment centré d'ordre 3  
 c'est un coefficient de dissymétrie.

coefficient de Fisher  $g_1 = \frac{m_3}{\sigma^3}$  (sans dimension)

coefficient empirique de Pearson:  $S_k = \frac{\bar{x} - m_0}{\sigma}$

coefficient empirique de Yule-Kendall:  $Y_k = \frac{Q_{3/4} + Q_{1/4} - 2Q_2}{Q_{3/4} - Q_{1/4}}$

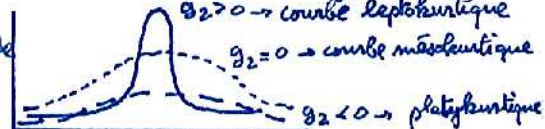
Attention: l'étude des signes de  $S_k$  et  $Y_k$  peut donner des contradictions.

\*  $m_4 = \overline{(x - \bar{x})^4} = \frac{1}{n} \sum_{i=1}^n m_i (x_i - \bar{x})^4$

coefficient d'aplatissement de Pearson:  $b_2 = \frac{m_4}{\sigma^4}$

coefficient d'aplatissement de Fisher:  $g_2 = \frac{m_4}{\sigma^4} - 3$

$g_2 > 0 \rightarrow$  courbe leptokurtique  
 $g_2 = 0$  pour la cloche  $g_2 = 0 \rightarrow$  courbe mésokurtique  
 $g_2 < 0 \rightarrow$  platykurtique



## II Analyse bivariate

On part de deux caractères  $x_i, y_i$   $1 \leq i \leq n$ . On ordonne strictement  $x'_1 < \dots < x'_p$  et  $y'_1 < \dots < y'_q$ .  
Série statistique double  $\{(x'_i, y'_j, n_{ij}), 1 \leq i \leq p, 1 \leq j \leq q\}$

Tableau de contingence

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

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

$$n = \sum_{i=1}^p n_{i \cdot} = \sum_{j=1}^q n_{\cdot j}$$

$i \setminus j$	$y'_1$	$\dots$	$y'_j$	$\dots$	$y'_q$	effectifs marginaux
$x'_1$	$n_{11}$	$\dots$	$n_{1j}$	$\dots$	$n_{1q}$	$n_{1 \cdot}$
$\vdots$	$\vdots$		$\vdots$		$\vdots$	
$x'_i$	$n_{i1}$	$\dots$	$n_{ij}$	$\dots$	$n_{iq}$	$n_{i \cdot}$
$\vdots$	$\vdots$		$\vdots$		$\vdots$	
$x'_p$	$n_{p1}$	$\dots$	$n_{pj}$	$\dots$	$n_{pq}$	$n_{p \cdot}$
effectifs marginaux	$n_{\cdot 1}$		$n_{\cdot j}$		$n_{\cdot q}$	$n$

covariance :  $\sigma_{xy} = \overline{(x-\bar{x})(y-\bar{y})} = \frac{1}{n} \sum_{1 \leq i, j \leq n} n_{ij} (x'_i - \bar{x})(y'_j - \bar{y}) = \overline{xy} - \bar{x}\bar{y}$

coefficient de corrélation de Bravais-Pearson :  $\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} \in [-1, 1]$

matrice de covariance  $\begin{pmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{pmatrix}$

## Régression linéaire

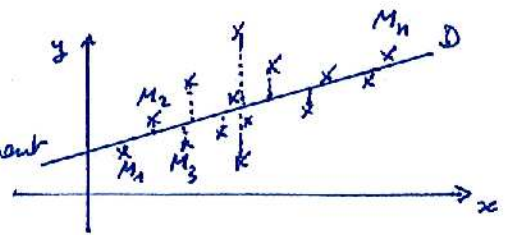
nuage de points  $\{(x_i, y_i), 1 \leq i \leq n\}$

méthode des moindres carrés : droite d'ajustement

$$y = ax + b \text{ minimisant } \sum_{i=1}^n d(M_i, D)^2$$

régression de  $y$  par rapport à  $x$  :

$d_V(M_i, D)$  : distance "verticale" de  $M_i$  à  $D$ .



$$\begin{cases} a = \frac{\sigma_{xy}}{\sigma_x^2} & \text{coefficient de régression (pente } a) \\ b = \bar{y} - a\bar{x} \end{cases}$$

$$y = \bar{y} + \frac{\sigma_{xy}}{\sigma_x^2} (x - \bar{x})$$

variance résiduelle :  $\sigma_{y \cdot x}^2 = \frac{1}{n} \sum (y_i - ax_i - b)^2 = \sigma_y^2 (1 - \rho_{xy}^2)$

régression de  $x$  par rapport à  $y$  :  $x = a'y + b'$  minimisant  $\sum_{i=1}^n d_H(M_i, D)^2$ ,  $d_H$  distance "horizontale".

$$\begin{cases} a' = \frac{\sigma_{xy}}{\sigma_y^2} & \text{(pente } \frac{1}{a'}) \\ b' = \bar{x} - a'\bar{y} \end{cases}$$

$$x = \bar{x} + \frac{\sigma_{xy}}{\sigma_y^2} (y - \bar{y})$$

En général ces deux droites sont distinctes.  $aa' = \frac{\sigma_{xy}^2}{\sigma_x^2 \sigma_y^2} = \rho_{xy}^2$ .

Les pentes se rapprochent de plus en plus lorsque  $|\rho_{xy}|$  est proche de 1.

On considère l'ajustement valable pour  $0,7 \leq |\rho_{xy}| \leq 1$ , excellent si  $|\rho_{xy}| \geq 0,95$ .

autres ajustements :  $y = ax^a \rightarrow$  régression linéaire en  $(\ln x, \ln y)$  :  $\ln y = a \ln x + \ln a$

$y = ae^{bx} \rightarrow$  régression linéaire en  $(x, \ln y)$  :  $\ln y = bx + \ln a$ .

Bibliographie : J.-J. Dooebeke : éléments de statistiques 1988, chapitres 2, 3, 10

H. Egon : Statistique et probabilités 1992, chapitre A

G. Demongel : probabilités, statistique inférentielle, fiabilité, 1997, chapitres 1, 2

M. Gaultier : Statistique 1997, chapitre 2, 7.