

Université de Batna

2005/2006

Faculté de Médecine

Département de Pharmacie

Cours de Statistique

1^{ère} Année Pharmacie

Chapitre II : Les séries statistiques doubles

D'après le cahier de :

I. Hadeef

Chapitre II: Les séries statistiques doubles.

Definition:

Une série statistique double ou série à 2 variables est un ensemble de résultats issus de l'observation de 2 caractères numériques sur une même population

exp:

- 1/ La taille et le poids d'un groupe d'enfants.
- 2/ le salaire est la qualification d'un ensemble de salariés.
- 3/ la température est la pression d'un milieu à différentes heures.
- * la série statistique des couples (X_i, Y_i) sera notée (X, Y)

Notation et représentation des séries statistiques doubles.

Une série statistique double peut être donnée comme l'énumération d'un certain nombre de résultats, on distingue deux cas:

- 1/ les effectifs des valeurs égaux à 1.
- 2/ Les effectifs des couples (X_i, Y_j) égaux n_{ij}

1/ L'effectif des valeurs (X_i, Y_j) égaux à 1:

Dans ce cas une série statistique se présente comme suit :

X	x_1	x_2	...	x_i
Y	y_1	y_2	...	y_i

la moyenne, la variance et l'écart type.

Les variables statistiques X et Y ont chacune une moyenne, variance et l'écart type

$$\bar{X} = \sum_{i=1}^n \frac{x_i}{n}, \quad V(X) = \sum_{i=1}^n \frac{x_i^2}{n} - \bar{X}^2$$

$$s = \sqrt{V}$$

$$\bar{Y} = \sum_{i=1}^n \frac{y_i}{n}, \quad V(Y) = \sum_{i=1}^n \frac{y_i^2}{n} - \bar{Y}^2$$

$$s(Y) = \sqrt{V(Y)}$$

2/ Les effectifs des valeurs (des couples) (x_i, y_j) égaux n_{ij} .

Soit X et Y deux séries statistiques définies sur la même population, prenant les valeurs x_1, x_2, \dots, x_n et y_1, y_2, \dots, y_n .

Définition : On appelle effectifs partiels n_{ij} des couples (x_i, y_j) le nombre d'individus de

la population pour les quelles la série X prend la valeur X_i et la série Y prend la valeur Y_j .

Tableau d'effectif:

$X \backslash Y$	Y_1	Y_2	...	Y_j	...	Y_s	total
X_1	n_{11}	n_{12}	...	n_{1j}	...	n_{1s}	$n_{1\cdot}$
\vdots							
X_i	n_{i1}	n_{i2}	...	n_{ij}	...	n_{is}	$n_{i\cdot}$
\vdots							
X_r	n_{r1}	n_{r2}	...	n_{rj}	...	n_{rs}	$n_{r\cdot}$
totale	$n_{\cdot 1}$	$n_{\cdot 2}$...	$n_{\cdot j}$...	$n_{\cdot s}$	n

Dans n_{ij} le 1^{er} indice i correspond à la ligne.
 - Le 2^{ème} indice j correspond à la colonne.

La définition: Si n_{ij} est l'effectif partiel d'un couple (X_i, Y_j) et n l'effectif total. le rapport $f_{ij} = \frac{n_{ij}}{n}$ s'appelle la fréquence partielle du

couple (X_i, Y_j) .

Remarque: $\sum_{j=1}^s \left(\sum_{i=1}^r n_{ij} \right) = n$.

La somme des effectifs partiels correspondant à la valeur X_i est égale à l'effectif des individus pour lesquels X prend la valeur X_i et on note $n_{i.}$

$$n_{i.} = \sum_{j=1}^p m_{ij} = m_{i1} + m_{i2} + \dots + m_{ip}$$

De la même manière on définit $n_{.j}$

$$n_{.j} = \sum_{i=1}^n m_{ij} = m_{1j} + m_{2j} + m_{3j} + \dots + m_{nj}$$

Caractéristiques d'une série à deux variables quantitatives :

Soit (X, Y) une série statistique double quantitative définies sur une même population

la moyenne de X : $\bar{X} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p m_{ij} X_i$

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n m_{i.} X_i$$

la moyenne de Y

$$\bar{Y} = \frac{1}{n} \sum_{j=1}^p \sum_{i=1}^n m_{ij} Y_j$$

$$\bar{Y} = \frac{1}{n} \sum_{j=1}^p n_{.j} Y_j$$

La variance de X :

$$V(X) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p m_{ij} (x_i - \bar{X})^2$$
$$= \frac{1}{n} \sum_{i=1}^n m_{i\cdot} (x_i - \bar{X})^2$$

Après le développement de la forme :

$$V(X) = \frac{1}{n} \sum_{i=1}^n m_{i\cdot} x_i^2 - \bar{X}^2$$

L'écart type de X : $\sigma_X = \sqrt{V(X)}$

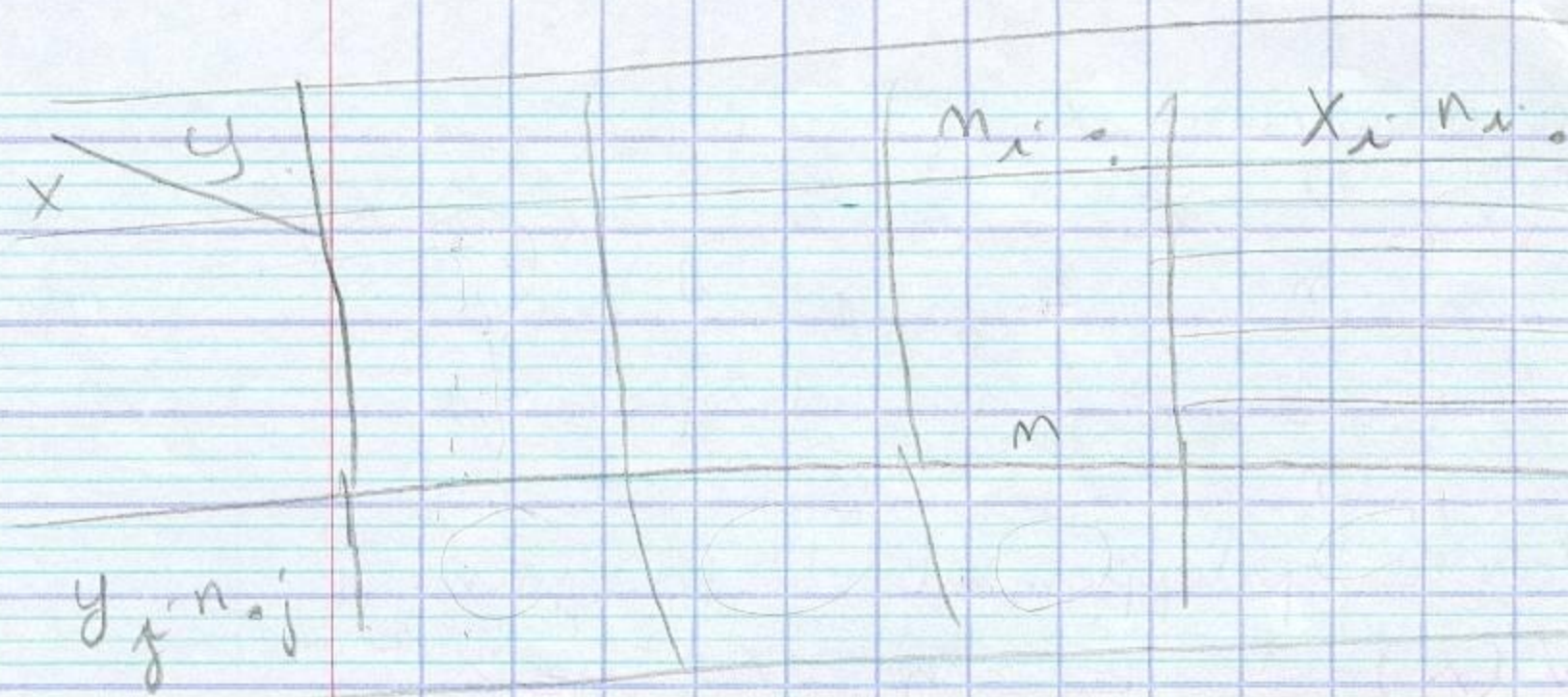
La variance Y : $V(Y) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p m_{ij} (y_j - \bar{Y})^2$

$$V(Y) = \frac{1}{n} \sum_{j=1}^p m_{\cdot j} (y_j - \bar{Y})^2$$

$$V(Y) = \left(\frac{1}{n} \sum_{j=1}^p m_{\cdot j} y_j^2 \right) - (\bar{Y})^2$$

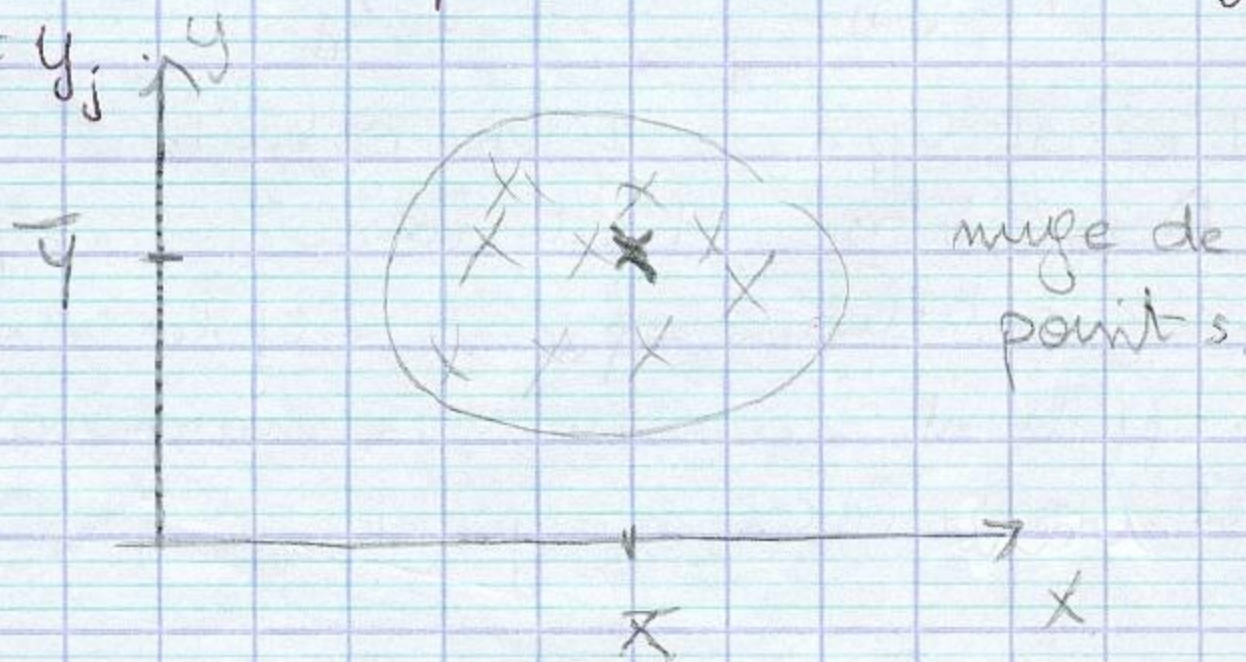
L'écart type de Y se calcule : $\sigma_Y = \sqrt{V(Y)}$

Definition : étant donné un couple (X, Y) de série statistique échantillonnées les moyennes \bar{X} , \bar{Y} s'appellent les moyennes marginales, $V(X)$ et $V(Y)$ les variances marginales.



Nuage de point.

Dans un repère cartésien on représente la série statistique à 2 variables par les points M_{ij} de coordonnées x_i et y_j . L'ensemble de ces points s'appelle le nuage de points de la série double (x, y) . C'est une représentation importants des liens qu'il y a entre les x_i et y_j .



Le point moyen : le point moyen de la série (X, Y) est le point $\bar{M}(\bar{X}, \bar{Y})$.

dont les coordonnées sont les moyennes \bar{X} et \bar{Y} des X_i et Y_j .

Définition de la covariance :

Si (X, Y) désigne un couple des séries statistiques quantitatives prenant respectivement :

$$\begin{matrix} X_1, X_2, \dots, X_n \\ Y_1, Y_2, \dots, Y_n \end{matrix}$$

n_{ij} l'effectif partiel du (X_i, Y_j)

On appelle covariance du couple (X, Y) et on note

$$\text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p n_{ij} (X_i - \bar{X})(Y_j - \bar{Y}).$$

$$= \frac{1}{n} \left(\sum_{i=1}^n \sum_{j=1}^p n_{ij} X_i Y_j \right) - (\bar{X} \cdot \bar{Y})$$

$$= 6 \times 4$$

Propriétés :

$$1/ \text{Cov}(X, X) = V(X), \quad \text{Cov}(Y, Y) = V(Y)$$

$$2/ \text{Cov}(X, Y) = \text{Cov}(Y, X)$$

Proposition : Si X et Y sont indépendantes \Rightarrow

$$\text{Cov}(X, Y) = 0.$$

L'ajustement = (x, y) une série double.

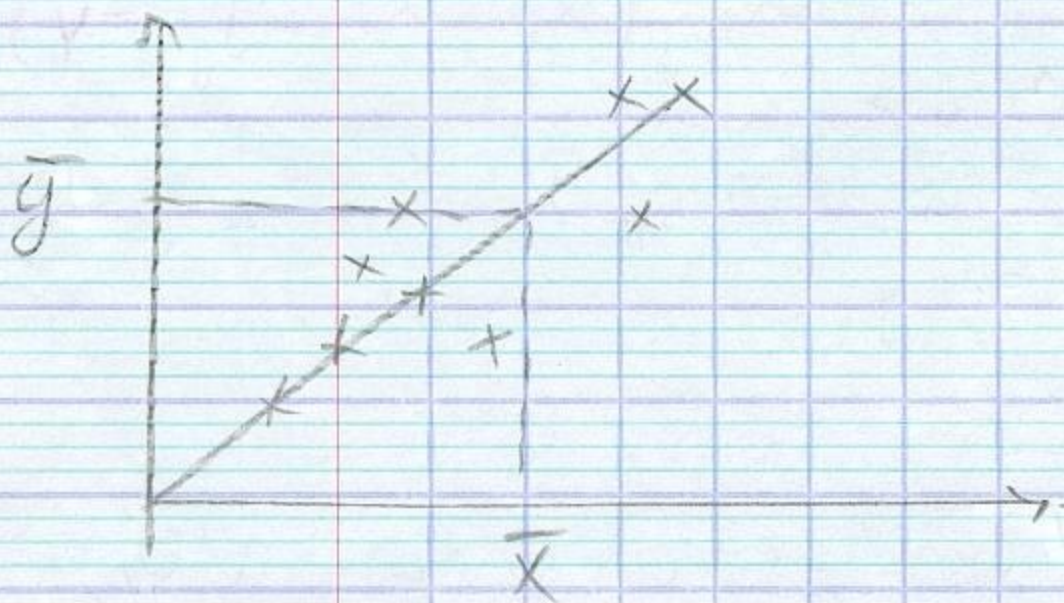
$M_{ij} = (x_i, y_j)$ L'ensemble des points de la série

Définition: Ajuster un ensemble de points consiste à déterminer une courbe simple (aussi proche que possible de l'ensemble M_{ij}).

L'ajustement linéaire: lorsque le nuage de points d'une série double semble avoir une forme rectiligne. On cherche à approcher par une fonction affine.

(ajustement linéaire)

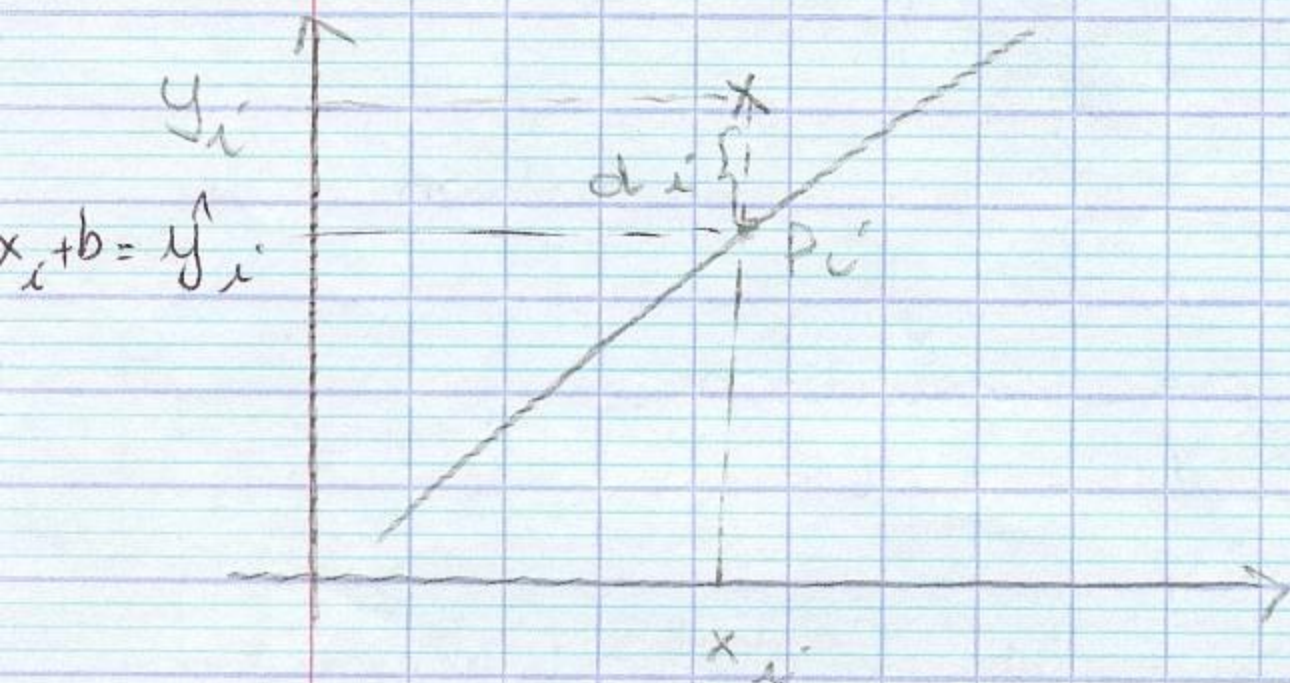
On dit que l'on effectue un ajustement linéaire de ce nuage de points et on note : $y = ax + b$.
L'équation de la droite recherchée.



Méthode de Moindres Carrés:

On cherche une équation $y = ax + b$ de la droite (D) telle que la somme de ses distances

au différentes points représentant soit minimal
 La distance choisie est le carré de la différence des ordonnées entre chaque point et la droite, de la droite ayant le même abscisse.



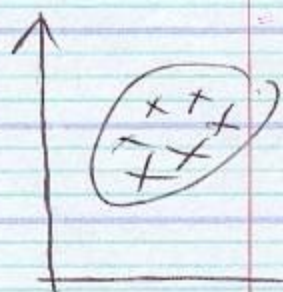
$$M_i(x_i, y_i), P_i(x_i, \hat{y}_i)$$

$$\sum_i d_i^2 = \sum_i (y_i - \hat{y}_i)^2 \text{ soit minimum.}$$

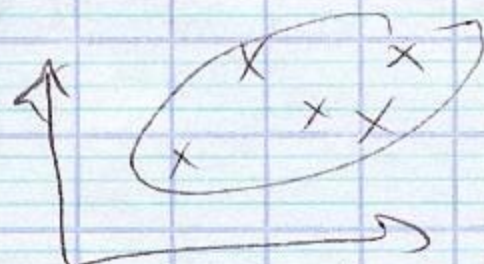
coefficient de corrélation linéaire

La méthode de Moindres carrés peut être utilisée pour n'importe quelle série double. il existe une droite d'estimation par la méthode moindres carrés pour s'assurer d'une façon objective et non risquée que l'ajustement est valide ou non. On calcule le coefficient r de corrélation linéaire.

Corrélation : Soit (x, y) une série double. On peut voir sur le nuage de points si les 2 variables sont corrélées ou non.



corrélation forte



corrélation nulle.

Le Coefficient de corrélation :

$$r = \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y} \quad -1 \leq r \leq 1$$

s'il est voisin en valeur de 1 l'ajustement est valide. $(0,70 \leq |r| \leq 1)$

Application de la méthode de moindres-carrés.

Regression Y en X : (A) : $\hat{y} = ax + b$

le problème est donc de calculer a et b.

pour que la somme des termes $(y_i - \hat{y}_i)^2$ est minimum. Cela revient à chercher le minimum de la fonction

Φ de deux variables a, b.

$$\begin{aligned} \Phi(a, b) &= \sum (y_i - ax_i - b)^2 \\ &= \sum (y_i - \hat{y}_i)^2 \end{aligned}$$

Φ admet un minimum si $\frac{\partial \Phi}{\partial a} = 0$ et $\frac{\partial \Phi}{\partial b} = 0$

$$\left\{ \begin{aligned} \frac{\partial \Phi}{\partial a} &= 2 \sum_{i=1}^n -x_i (y_i - ax_i - b) = 0 \\ \frac{\partial \Phi}{\partial b} &= 2 \sum_{i=1}^n -(y_i - ax_i - b) = 0 \end{aligned} \right.$$

$$\left\{ \begin{aligned} \sum_i -x_i y_i + a \sum_i x_i^2 + b \sum_i 1 &= 0 \quad \dots (1) \\ -\sum_i y_i + a \sum_i x_i + \sum_i b &= 0 \quad \dots (2) \end{aligned} \right.$$

$$(2) \Leftrightarrow -n\bar{y} + a n\bar{x} + nb = 0$$

$$\Leftrightarrow b = \bar{y} - a\bar{x} \Leftrightarrow \bar{y} = b + a\bar{x}$$

$$(1) \Leftrightarrow -\sum_i x_i y_i + a \sum_i x_i^2 + (\bar{y} - a\bar{x}) \sum_i x_i = 0$$

$$\Rightarrow -\sum_i x_i y_i + a \sum_i x_i^2 + (\bar{y} - a\bar{x}) n\bar{x} = 0$$

$$a = \frac{\sum x_i y_i - n\bar{x}\bar{y}}{\sum x_i^2 - n\bar{x}^2}$$

$$a = \frac{\text{cov}(X, Y)}{V(X)} \quad (\text{D}) \text{ la droite de regression } Y \text{ en } X$$

$$X, Y = aX + b \quad \text{d'où } a = \frac{\text{cov}(X, Y)}{V(X)}, \quad b = \bar{y} - a\bar{x}$$

En permutant les rôles de X et Y on obtient :

la droite (X') de regression X en Y : $X = a'Y + b'$

$$\begin{cases} b' = \bar{x} - a'\bar{y} \\ a' = \frac{\text{cov}(Y, X)}{V(Y)} \end{cases}$$

$$a \cdot a' = r^2$$

$$\Rightarrow |r| = \sqrt{aa'}$$

Récapitulons le coefficient de corrélation linéaire :

Les 2 droites de régression trouvées sont :

$$y = ax + b \quad a = \frac{\text{cov}(x, y)}{V(x)}$$

$$x = \bar{a}y + b' \quad \bar{a} = \frac{\text{cov}(x, y)}{V(y)}$$

$$a \bar{a} = \frac{\text{cov}^2(x, y)}{V(x) \cdot V(y)} = \frac{\text{cov}^2(x, y)}{\sigma^2(x) \cdot \sigma^2(y)} = r^2$$

Lorsque les droites sont identiques les coefficients a et \bar{a} sont égaux : $\frac{1}{a} = \bar{a} \Rightarrow a \bar{a} = 1 \Rightarrow r^2 = 1 \Rightarrow |r| = 1$

Si les droites sont proches $|r|$ est voisin de 1, ce qui correspond à un ajustement valide. Si $|r|$ n'est pas très différent de zéro, c'est que a et \bar{a} sont loin d'être inverse l'une de l'autre, et par conséquent l'ajustement est non valide.

L'ajustement non linéaire :

Il peut arriver que le nuage de points représentant une série double ne soient pas alignés, mais soient voisins d'une courbe connue.

1/ L'ajustement à l'aide d'une fonction exponentielle :

$$y = B \cdot A^x \quad (x \text{ : abscisse}, y \text{ : ordonnée})$$

Si a des valeurs x_i en progression arithmétique

correspondant à des valeurs y_i sensiblement en progression géométrique.

On peut ajuster la série par une fonction exponentielle

$$Y = B \cdot A^X \iff \ln Y = \ln(B \cdot A^X)$$

$$\ln Y = \ln B + X \ln A \text{ de la forme : } y = b + ax$$

$$\text{avec : } \begin{cases} y = \ln Y \\ a = \ln A \end{cases}$$

$$\begin{cases} a = \frac{\text{cov}(X, Y)}{V(X)} \\ b = \bar{y} - a \bar{x} \end{cases}$$

$$\begin{cases} a = \ln A \\ b = \ln B \end{cases} \iff \begin{cases} A = e^a \\ B = e^b \end{cases}$$

$$Y = B \cdot A^X = e^b (e^a)^X$$

2/ l'ajustement à l'aide d'une fonction puissance

$$Y = B X^a$$

Si une série statistique varie de tel sorte que les valeurs x_i et y_i sont sensiblement en progression géométrique. On peut ajuster par une fonction puissance d'équation : $Y = B X^a$.

$$\ln Y = \ln B + a \ln X \text{ de la forme : } y = b + ax$$

$$\text{avec } \begin{cases} y = \ln Y \\ x = \ln X \\ b = \ln B \end{cases}$$

$$\text{d'où } a = \frac{\text{cov}(X, y)}{V(X)} = \frac{\text{cov}(\ln X, \ln Y)}{V(\ln X)}$$

$$b = \bar{y} - a \bar{x} = \ln \bar{y} - a \ln \bar{X}$$

$$\ln B = b \Leftrightarrow B = e^b$$

$$y = B x^a = e^b x^a$$

3/ L'ajustement à l'aide d'une fonction parabole.

Si la présentation graphique de la série semble permettre l'ajustement par une parabole d'équation

$$y = a + bX + cX^2$$

le problème consisterait à rechercher a, b, c de telle sorte que $\sum (y_i - a - bx_i - cx_i^2)^2$ soit minimum /

on définit Φ une fonction de 3 variables.

$$\Phi(a, b, c) = \sum (y_i - a - bx_i - cx_i^2)^2$$

Φ est minimum lorsque :

$$\frac{\partial \Phi}{\partial a} = 0, \quad \frac{\partial \Phi}{\partial b} = 0, \quad \frac{\partial \Phi}{\partial c} = 0 \text{ alors :}$$

$$2 \sum (y_i - a - bx_i - cx_i^2) = 0$$

$$2 \sum x_i (y_i - a - bx_i - cx_i^2) = 0$$

$$2 \sum x_i^2 (y_i - a - bx_i - cx_i^2) = 0$$

$$na + b \sum x_i + c \sum x_i^2 = \sum y_i$$

$$a \sum x_i + b \sum x_i^2 + c \sum x_i^3 = \sum x_i y_i$$

$$a \sum x_i^2 + b \sum x_i^3 + c \sum x_i^4 = \sum x_i^2 y_i$$

c'est un système de 3 équations à 3 variables