

## II - Séries Statistiques à deux variables.

### Chapitre 2 - Suite -

#### → la Régression linéaire Simple

#### 1. Introduction:

Dans cette partie, nous allons nous intéresser à des séries statistiques à deux variables  $x$  et  $y$ . On se demande s'il existe une relation de type linéaire entre ces deux variables, sachant que d'autres types de relations peuvent exister. L'objectif serait donc de trouver cette relation linéaire, mesurer le degré de liaison entre  $x$  et  $y$  et si cette liaison est assez forte, procéder à la prévision de  $y$  sachant  $x$  ou le contraire.

#### 2. Les séries statistiques à deux variables.

La série statistique est décrite par deux variables,  $x$  et  $y$  qui prennent les valeurs  $x_i$  et  $y_j$ . Les données sont regroupées dans un tableau à double entrée (tableau de contingence) ds lequel apparaissent les effectifs  $n_{ij}$  et les distributions marginales  $n_i$  et  $n_j$  de  $x$  et  $y$ .

#### 2.1. Moyennes

elles se calculent pour chacune des deux variables séparément de l'autre.

$$\bar{X} = \frac{1}{N} \sum_i n_i x_i$$

$$\bar{Y} = \frac{1}{N} \sum_j n_j y_j$$

#### 2.2. Variances et écarts-types

$$s_x^2 = \frac{1}{N} \sum_i n_i x_i^2 - \bar{X}^2$$

$$s_x = \sqrt{s_x^2}$$

$$s_y^2 = \frac{1}{N} \sum_j n_j y_j^2 - \bar{Y}^2$$

$$s_y = \sqrt{s_y^2}$$

### 2.3. Covariance entre variables

Elle mesure le sens de variation entre les variables

$$\text{Cov}(x, y) = \frac{1}{N} \sum_i \sum_j r_{ij} (x_i - \bar{x})(y_j - \bar{y})$$

ou

$$\text{Cov}(x, y) = \frac{1}{N} \sum_i \sum_j r_{ij} x_i y_j - \bar{x} \bar{y}$$

NB: La covariance peut être négative. Le signe moins (-) traduit une relation inverse entre les variables.

Par définition,  $\text{Cov}(x, x) = \sigma_x^2$  et  $\text{Cov}(y, y) = \sigma_y^2$ , et si les variables  $x$  et  $y$  sont indépendantes alors la covariance est nulle.

### 3. Détermination de la liaison entre $x$ et $y$ .

Sur un graphique, on représente les points  $(x_i, y_j)$ , on obtient un nuage statistique où tous les points ont le même poids. En fonction de la forme du nuage, on peut décider d'ajuster une droite ou de déterminer une relation linéaire entre les variables, tel que:

$$y = ax + b \quad \text{ou} \quad x = cy + d$$

Il reste que le sens de la relation dépend de l'existence d'un rapport logique entre les variables.

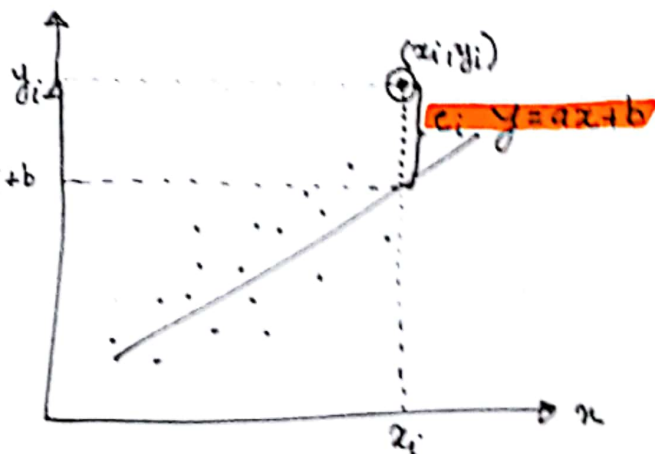
Plusieurs moyens sont disponibles pour déterminer la relation supposée linéaire entre  $x$  et  $y$ , ici nous focalisons l'attention sur la méthode des moindres carrés.

#### - La méthode des moindres carrés:

On considère une série statistique  $(x_i, y_j)$  de  $N$  observations, l'idée est d'ajuster une droite qui minimise les écarts entre les observations réelles et celles qui sont sur la droite.

Pour déterminer la droite d'ajustement  $y = ax + b$ , on minimise la somme des carrés <sup>des écarts</sup> entre les observations réelles et leurs projections

verticales sur la droite.



Soit la droite d'équation  $y = ax + b$

- les observations réelles  $x_i, y_i$
- les observations théoriques  $x_i, ax_i + b$

On minimise la somme des carrés des

$$e_i = y_i - ax_i - b$$

qui revient à annuler les dérivées partielles de cette expression par rapport à  $a$  et  $b$

On trouve les résultats suivants :

$$a = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

NB: de la même manière si  $x = cy + d$

$$c = \frac{\text{Cov}(x, y)}{\text{Var}(y)}, \quad d = \bar{x} - c\bar{y}$$

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

7. Ajustement se ramenant à un ajustement linéaire

4.1. Ajustement de la forme  $y = ba^x$

en prenant le log décimal

$$\log y = \log b + x \log a$$

$$Y = B + Ax$$

$$Y = Ax + B$$

$$Y = \log y$$

$$A = \log a$$

$$B = \log b$$

$$A = \frac{\text{Cov}(x, Y)}{\text{Var}(x)}$$

$$B = \bar{Y} - A\bar{x}$$

$$a = 10^A$$

$$b = 10^B$$



## 4.2. Ajustement se présentant sous la forme $y = b x^a$

En posant :

$$Y = \log y$$

$$B = \log b$$

$$X = \log x$$

$$y = b x^a \text{ devient } Y = a X + B$$

a et B sont déterminés en travaillant sur les séries ( $X = \log x, Y = \log y$ )

On obtient :

$$a = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

$$B = \bar{Y} - a \bar{X} \quad b = 10^B$$

NB : une telle relation se construit sur du papier log-log et met en évidence une relation de type puissance.

## 5. la Prédiction

Une fois la relation déterminée, entre  $x$  et  $y$  on peut faire des prévisions. Connaissant les valeurs de la variable explicative  $x$  on peut déterminer la valeur de la variable expliquée  $y$ .

## 6. la Corrélation

On appelle coefficient de corrélation linéaire, entre deux variables  $x$  et  $y$ , le réel  $r$  défini comme suit en utilisant la covariance entre  $x$  et  $y$  :

$$r = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$$

ou

$$r = \frac{\sum X_i Y_i}{\sqrt{\sum X_i^2 \sum Y_i^2}}$$

avec  $X = x - \bar{x}$  (écart par rapport à  $\bar{x}$ )

$Y = y - \bar{y}$  (écart " "  $\bar{y}$ )

(c'est la méthode des écarts à la moyenne)

$r$  peut être positif ou négatif selon le signe de la covariance  $\text{Cov}(x, y)$

-  $r = 0$  : il n'y a pas de relation entre  $x$  et  $y$ .

-  $|r| = 1$  : la relation linéaire entre  $x$  et  $y$  est parfaite.

-  $0.8 < |r| < 1$  : on admet qu'il y a une forte corrélation entre  $x$  et  $y$  et on accepte l'ajustement linéaire

-  $r^2$  est le coefficient de détermination et exprime le pourcentage de la variance expliquée par la droite de régression.

NB: avant tout calcul, il faut vérifier la qualité de la relation mesurée par  $r$ .  
une méthode rapide de le faire est celle des écarts à la moyenne :

1. Calculer  $\bar{x}$  et  $\bar{y}$

2. Calculer  $X_i = x_i - \bar{x}$  et  $Y_i = y_i - \bar{y}$ ,  $\sum_i X_i^2$  et  $\sum_i Y_i^2$

3. "  $\sum_i X_i Y_i$

4. Le coefficient de corrélation  $r$  est :

$$r = \frac{\sum_i X_i Y_i}{\sqrt{\sum_i X_i^2 \sum_i Y_i^2}}$$

## 7. La régression multiple.

Dans la régression simple, il s'agit de trouver la relation entre deux variables  $x$  et  $y$ . Pour un type linéaire, il s'agit de trouver les paramètres de la régression linéaire  $a$  et  $b$  tel que :

$$y = ax + b \quad (x: \text{variable explicative} \quad y: \text{var. expliquée})$$

pour estimer  $a$  et  $b$ , on résout un système de 2 équations et 2 inconnues en minimisant la somme des carrés des écarts entre les valeurs observées et celles théoriques.

même chose si on fait la régression de  $x$  en  $y$ , c'est-à-dire

$$x = cy + d, \quad \text{l'estimation de } c \text{ et } d \text{ obéit au même principe.}$$

( $x$ : variable expliquée,  $y$ : var. explicative.)