

Séries statistiques à 2 variables numériques. Nuage de pts associé  
ajustement affine par la méthode des moindres carrés. Droites  
de regression. Applications

Pré-requis : TES avec des outils de TS

Pré-requis - notion de var

-  $m \in \mathbb{N}^*$ ,  $\Omega = \{\omega_1, \dots, \omega_m\}$  population de  $m$  individus  
 $X$  var sur  $\Omega$ .  $\forall i \in \Delta_m, \alpha_i = X(\omega_i)$

$$\mu(X) = \bar{X} = \frac{1}{m} \sum_{i=1}^m \alpha_i$$

$$V(X) = \sigma^2(X) = \mu((X - \bar{X})^2) = \mu(X^2) - (\mu(X))^2$$

## 1. séries statistiques à 2 variables numériques

### 1.1. Exemple introductif

On considère 5 élèves  $1 \leq \alpha_i \leq 5$

Note $\alpha_i$ au bac	7	10	11	13	16
Note $y_i$ au concours	8	9	12	12	13

$$\bar{X} = 11,4 \quad V(X) = 9,04 \quad \sigma(X) = 3,01$$

$$\bar{Y} = 10,8 \quad V(Y) = 3,76 \quad \sigma(Y) = 1,94$$

### 1.2. Définition

$m \in \mathbb{N}^*$   $\Omega = \{\omega_1, \dots, \omega_m\}$  une population.

On appelle variable statistique quantitative (ou caractère quanti-  
tatif) toute application  $X: \Omega \rightarrow \mathbb{R}$ . On fait  $X$  est une var.

Une série statistique est la donnée d'un tableau où

$X$	$\alpha_1$	...	$\alpha_m$
$Y$	$y_1$	...	$y_m$



Rq Il existe des séries statistiques qui ont + de 2 caractères

### 1.3 Nuage de points associé

Definition Dans un repère orthogonal,  $\{M_i(x_i, y_i), 1 \leq i \leq m\}$  est le nuage de points associé à la série.

$G(\bar{x}, \bar{y})$  est le point moyen du nuage

### 1.4 Covariance

La covariance permet d'étudier la dépendance entre  $X$  et  $Y$

$$\text{Cov}(X, Y) = \mu[(X - \bar{X})(Y - \bar{Y})] = \mu(XY) - \mu(X)\mu(Y)$$

Dans l'ex.  $\text{Cov}(X, Y) = 5,28$

### 2. Ajustement affine par la méthode des moindres carrés

$X$  et  $Y$  non constantes ie  $\sigma(X) > 0$  et  $\sigma(Y) > 0$

#### 2.1 Principe

Soit  $f: \mathbb{R} \rightarrow \mathbb{R}$ ,  $f(x) = y - f(x)$  une fonction affine.  $f$  réalise un ajustement affine du nuage de points si  $f$  passe "le plus près possible" des points du nuage.

#### Méthode des moindres carrés

Cette méthode consiste à chercher  $a$  et  $b$  tels que  $\mu(Z)$  et  $\sigma^2(Z)$



$$\mu(Z^2) = \mu^2(Z) + \sigma^2(Z)$$

$$Z = Y - aX - b$$

$$\mu(Z^2) = \frac{1}{m} \sum_{i=1}^m (y_i - ax_i - b)^2$$

somme des résidus en  $Y$

en fait on veut minimiser  $\mu(Z^2)$

$$\mu(Z^2) = \sum_{i=1}^m M_i H_i^2$$

Rq Il existe d'autres méthodes pour minimiser notamment la méthode de la droite de Steyer



## 2.2 Théorème

Si  $Z = Y - aX - b$  alors :

$$\mu(Z) \text{ et } \sigma^2(Z) \text{ minimales} \iff a = \frac{\text{cov}(X,Y)}{\sigma^2(X)} \text{ et } b = \bar{Y} - a\bar{X}$$

### Démonstration

$$\mu(Z) = \mu(Y - aX - b) = \bar{Y} - a\bar{X} - b$$

$$Z - \bar{Z} = (Y - \bar{Y}) - a(X - \bar{X})$$

$$\text{donc } (Z - \bar{Z})^2 = (Y - \bar{Y})^2 - 2a(Y - \bar{Y})(X - \bar{X}) + a^2(X - \bar{X})^2$$

$$\begin{aligned} \sigma^2(Z) &= \mu[(Z - \bar{Z})^2] = \sigma^2(Y) - 2a \text{cov}(X,Y) + a^2 \sigma^2(X) \\ &= \underbrace{\left[ a\sigma(X) - \frac{\text{cov}(X,Y)}{\sigma(X)} \right]^2}_{\geq 0} + \underbrace{\sigma^2(Y) - \frac{\text{cov}^2(X,Y)}{\sigma^2(X)}}_{c \geq 0} \end{aligned}$$

$$\mu(Z^2) = (\bar{Y} - a\bar{X} - b)^2 + \left[ a\sigma(X) - \frac{\text{cov}(X,Y)}{\sigma(X)} \right]^2 + c$$

$$\mu(Z^2) \text{ minimale} \iff \begin{cases} a = \frac{\text{cov}(X,Y)}{\sigma^2(X)} \\ b = \bar{Y} - a\bar{X} \end{cases}$$

## 2.3 Droite de régression

D:  $y = ax + b$  est la droite de régression de Y en X. D passe par le pt moyen.

Cqqs: Elle permet d'estimer y en fct de x donné et inversement, on peut calculer la droite de régression de X en Y.

## 2.4 Coefficient de corrélation linéaire entre X et Y

Le coefficient permet de mesurer la validité de l'ajustement affine, il est égal à  $r(X,Y) = \frac{\text{cov}(X,Y)}{\sigma(X)}$

Rq On obtient une bonne corrélation si l'angle des droites de régression est  $\leq 80^\circ$



### Propriété

On a

$$-1 \leq r(x, y) \leq 1$$

$|r(x, y)| = 1 \iff$  les pts  $M_i$  sont alignés.

Dans l'ex:  $r(x, y) = 0,91$

forte corrélation car  $|r(x, y)| > 0,5$ .

### 3. Application

3.1 Déterminer les droites de  $Y$  en  $X$  et de  $X$  en  $Y$

- pour l'ex de la leçon

- si on considère un tron pour les pts  $A(1,1)$   $B(2,2)$   $D(3,1)$   $C(4,2)$

3.2 Le tableau donne l'évolution du taux d'équipement  $y_i$  (en  $Y$ ) en automobiles des familles françaises.

Année $x_i$	1975	1980	1985	1988	1991
taux $y_i$	64,1	69,3	73,4	74,6	76,8

Dans un repère représenter le nuage de pts associés à la série. La forme de ce nuage permet-elle d'envisager raisonnablement un ajustement affine?

Donner les eq. des droites de régression de  $Y$  en  $X$  et de  $X$  en  $Y$  les représenter.

Que concluez-vous expliquer?

Si l'évolution se poursuivait ainsi estimer le taux en 1993 et l'année à partir de laquelle ce taux sera  $\geq 85\%$ .