

Synthèse de cours (Terminale ES)

→ Séries statistiques à deux variables

Séries statistiques à deux variables

Définition

On appelle « série statistique à deux variables » la donnée de n couples (x_i, y_i) de valeurs réelles.

Nuage de points et point moyen

A chaque couple (x_i, y_i) on peut associer, dans un repère orthogonal, un point M_i de coordonnées (x_i, y_i) . L'ensemble des points ainsi obtenus est appelé « nuage de points associé à la série statistique ».

On considère alors les moyennes statistiques : $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ et $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$.

Le point, noté G, de coordonnées (\bar{x}, \bar{y}) est appelé « point moyen » associé à la série statistique.

Ajustement affine par la méthode des moindres carrés

Objectif

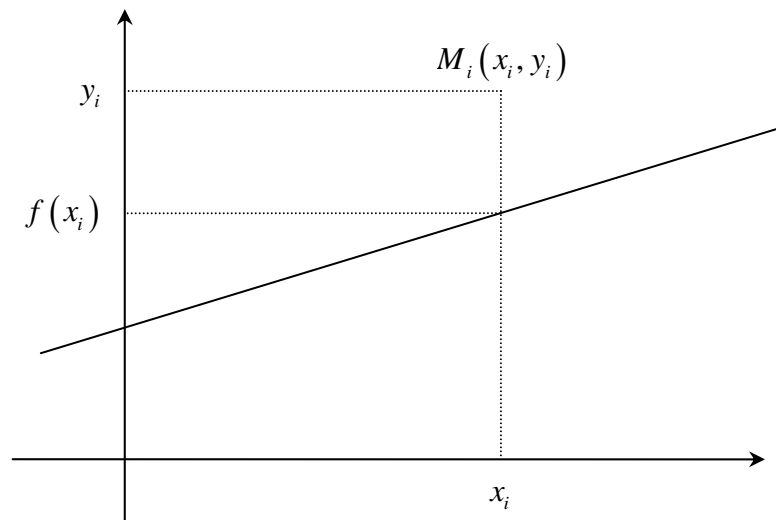
Effectuer un ajustement c'est chercher une fonction dont la représentation graphique décrit « au mieux » (le critère doit être précisé) le nuage de points associé à la série statistique considérée. Si la fonction cherchée est affine, on parle d'« ajustement affine ».

Si cette fonction affine est de la forme : $y = f(x) = ax + b$, on dit que l'on a effectué un « ajustement de y en x ». La variable « y » est alors appelée la « variable expliquée » et la variable « x », la « variable explicative ».

Disposant de f , on va pouvoir se donner une valeur de x quelconque et calculer $f(x)$:

- Si $x_{\min} < x < x_{\max}$, on dit que l'on effectue une « interpolation » ;
- Si $x < x_{\min}$ ou $x > x_{\max}$, on dit que l'on effectue une « extrapolation ».

Méthode des moindres carrés



→ Le principe général de cette méthode consiste à trouver une fonction f minimisant la quantité : $\sum_{i=1}^n (y_i - f(x_i))^2$ (d'où son nom ...)

Ajustement affine par la méthode des moindres carrés

On effectue un ajustement affine de y en x . On cherche f de la forme : $f(x) = ax + b$.

Dans ces conditions, on a :

$$a = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Remarque :

$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$ est la covariance des deux séries x_i et y_i . On la note : $\text{cov}(x, y)$;

$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ est la variance de la série y_i . On la note : $\text{var}(x)$

On peut donc écrire : $a = \frac{\text{cov}(x, y)}{\text{var}(x)}$.

Par ailleurs, la droite représentant la fonction f passe par le point moyen G . On a donc :

$$y - \bar{y} = a(x - \bar{x})$$

Et on en déduit :

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

La droite obtenue est appelée « droite de régression de y en x ».

Rappels sur la variance et la covariance

On a :

$$\text{var}(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$$

Pour retenir la seconde expression : « moyenne des carrés moins carré de la moyenne ».

On a :

$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n} \sum_{i=1}^n x_i \cdot y_i - \bar{x} \cdot \bar{y}$$

Pour retenir la seconde expression : « moyenne des produits moins produit de la moyenne ».