

(Suite du cours statistique MI)

L'ajustement linéaire par la Méthode des moindres carrés

1) Introduction

Une situation courante est d'avoir à sa disposition deux ensembles de données de taille n , $\{y_1, y_2, \dots, y_n\}$ et $\{x_1, x_2, \dots, x_n\}$, obtenus expérimentalement ou mesurés sur une population.

Le problème de l'ajustement linéaire ou de la régression consiste à rechercher une relation pouvant éventuellement exister entre les x et les y , par exemple de la forme $Y = f(X)$.

Lorsque la relation recherchée est affine, c'est-à-dire de la forme $y = ax + b$, on parle de régression linéaire.

Même si une telle relation est effectivement présente, les données mesurées ne vérifient pas en général cette relation exactement. Pour tenir compte dans le modèle mathématique des erreurs observées, on considère les données $\{y_1, y_2, \dots, y_n\}$ comme autant de réalisations d'une variable aléatoire Y et parfois aussi les données $\{x_1, x_2, \dots, x_n\}$ comme autant de réalisations d'une variable aléatoire X . On dit que la variable Y est la variable dépendante ou variable expliquée et que la variable X est la variable explicative.

2) La droite d'ajustement linéaire

Les données $\{(x_i, y_i), i = 1, \dots, n\}$ peuvent être représentées par un nuage de n points dans le plan (x, y) , le diagramme de dispersion. Le centre de gravité de ce nuage peut se calculer facilement : il s'agit du point aux coordonnées (\bar{X}, \bar{Y}) .

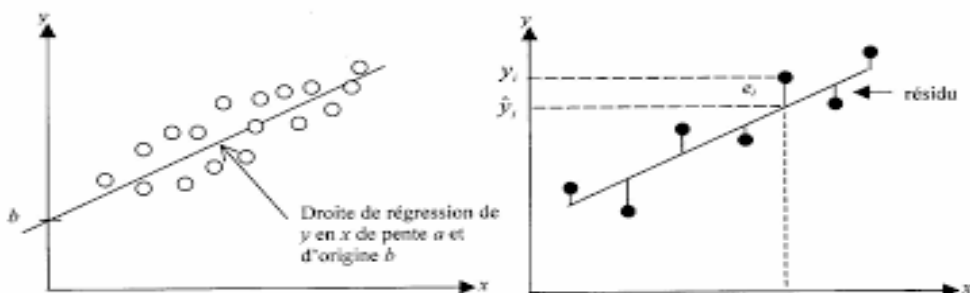
Rechercher une relation affine entre les variables X et Y revient à rechercher une droite qui s'ajuste le mieux possible à ce nuage de points. Parmi toutes les droites possibles, on retient celle qui rend **minimale la somme des carrés des écarts** des valeurs observées « y_i » à la droite $\hat{y}_i = ax_i + b$.

On pose $\varepsilon_i = y_i - \hat{y}_i$, qui n'est autre que la valeur qui représente cet écart (erreur), appelé aussi résidu.

Le principe des moindres carrés (MC) consiste à choisir les valeurs de « a » et de « b » qui minimisent la somme des erreurs commises sur chaque point.

Donc, on cherchera de minimiser cette fonction

$$E = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (ax_i + b))^2$$



Calculons a et b

Remarquons d'abord que notre fonction E a pour inconnues (a,b) les xi et yi étant connus expérimentalement. Les conditions initiales pour déterminer le minimum d'une fonction sont :

$$\frac{\partial E}{\partial a} = 0 \quad \text{et} \quad \frac{\partial E}{\partial b} = 0$$

$$\begin{aligned} \frac{\partial E}{\partial b} = 0 &\Leftrightarrow \sum_{i=1}^n 2 \cdot (-1)(y_i - ax_i - b) = 0 \\ &\Leftrightarrow \sum_{i=1}^n (y_i - ax_i - b) = 0 \\ &\Leftrightarrow \sum_{i=1}^n y_i - a \cdot \sum_{i=1}^n x_i - \sum_{i=1}^n b = 0 \\ &\Leftrightarrow n \cdot \bar{y} - a \cdot n \cdot \bar{x} - n \cdot b = 0 \\ &\Leftrightarrow \bar{y} - a \cdot \bar{x} - b = 0 \\ &\Leftrightarrow b = \bar{y} - a \cdot \bar{x} \end{aligned}$$

Premier résultat :

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

$$\begin{aligned} \frac{\partial E}{\partial a} = 0 &\Leftrightarrow \sum_{i=1}^n 2 \cdot (-a \cdot x_i)(y_i - ax_i - b) = 0 \\ &\Leftrightarrow \sum_{i=1}^n x_i \cdot (y_i - ax_i - b) = 0 \\ &\Leftrightarrow \sum_{i=1}^n x_i \cdot y_i - a \cdot \sum_{i=1}^n x_i^2 - \sum_{i=1}^n b \cdot x_i = 0 \\ &\Leftrightarrow n \cdot \bar{x}\bar{y} - a \cdot n \cdot \bar{x}^2 - n \cdot b \bar{x} = 0 \\ &\Leftrightarrow \bar{x}\bar{y} - a \cdot \bar{x}^2 - b \bar{x} = 0 \\ &\Leftrightarrow \bar{x}\bar{y} - a \cdot \bar{x}^2 - (\bar{y} - a \cdot \bar{x})\bar{x} = 0 \quad \text{car } b = \bar{y} - a \cdot \bar{x} \\ &\Leftrightarrow \bar{x}\bar{y} - a \cdot \bar{x}^2 - \bar{y} \cdot \bar{x} + a \cdot \bar{x}^2 = 0 \\ &\Leftrightarrow -a \cdot (\bar{x}^2 - \bar{x}^2) = -(\bar{x}\bar{y} - \bar{y} \cdot \bar{x}) \\ &\Leftrightarrow a \cdot V(X) = \text{cov}(X, Y) \\ &\Leftrightarrow a = \frac{\text{cov}(X, Y)}{V(X)} \end{aligned}$$

Deuxième résultat :

$$a = \frac{\text{cov}(X, Y)}{V(X)}$$

Conclusion :

La droite de régression Y/x, (y par rapport à x) est : $y = ax + b$

Avec

$$a = \frac{\text{cov}(X, Y)}{V(X)} \quad \text{et} \quad b = \bar{y} - a \cdot \bar{x}$$

$$\text{ie } y = \frac{\text{cov}(X, Y)}{V(X)} \cdot x + \bar{y} - \frac{\text{cov}(X, Y)}{V(X)} \cdot \bar{x}$$

Remarques :

- 1) La minimisation de E implique d'autres conditions sur les dérivées du second ordre, ceci n'étant pas le but de notre cours, on ne s'attarde pas dessus mais elles sont vérifiées.
- 2) Dans le calcul de la droite des moindres carrés, les variables X et Y ne jouent pas des rôles interchangeables. La variable dépendante Y prend, comme son nom l'indique, des valeurs qui dépendent de celles de X.
D'ailleurs si l'on échange les rôles de X et de Y, on aura une approximation linéaire de la forme $\hat{x}_i = a' y_i + b'$, le critère des MC est alors

$$E = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (x_i - \hat{x}_i)^2 = \sum_{i=1}^n (x_i - (a' y_i + b'))^2$$

et ce n'est plus le même et la droite que l'on obtient en général. Cette droite, tout comme la précédente, passe par le centre de gravité du nuage de point, mais c'est généralement leur seul point commun.

C'est le problème considéré qui indique s'il faut considérer Y ou plutôt X comme variable dépendante (et l'autre comme variable explicative).

Mais si l'on s'intéresse aux interactions entre deux variables X et Y dont ni l'une ni l'autre n'est clairement dépendante de l'autre, alors on pourra choisir de régresser Y en fonction de X ou bien le contraire. Mais on ne doit pas s'attendre à obtenir les mêmes résultats.

Ainsi la droite de régression de X/y (x par rapport à y) est : $x = a'y + b'$

Avec

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

$$\text{ie } x = \frac{\text{cov}(X,Y)}{V(Y)} \cdot y + \bar{x} - \frac{\text{cov}(X,Y)}{V(Y)} \cdot \bar{y}$$

(Exercice : montrez que les deux droites de régression sont égales si et seulement si $a \cdot a' = 1$.)

3) Le coefficient de corrélation linéaire

Pour évaluer la qualité de la régression linéaire et ainsi mesurer la qualité de l'approximation du nuage de points, on calcule son coefficient de corrélation linéaire défini par :

$$\rho_{xy} = r_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y}$$

Avec $-1 \leq |\rho_{xy}| \leq +1$,

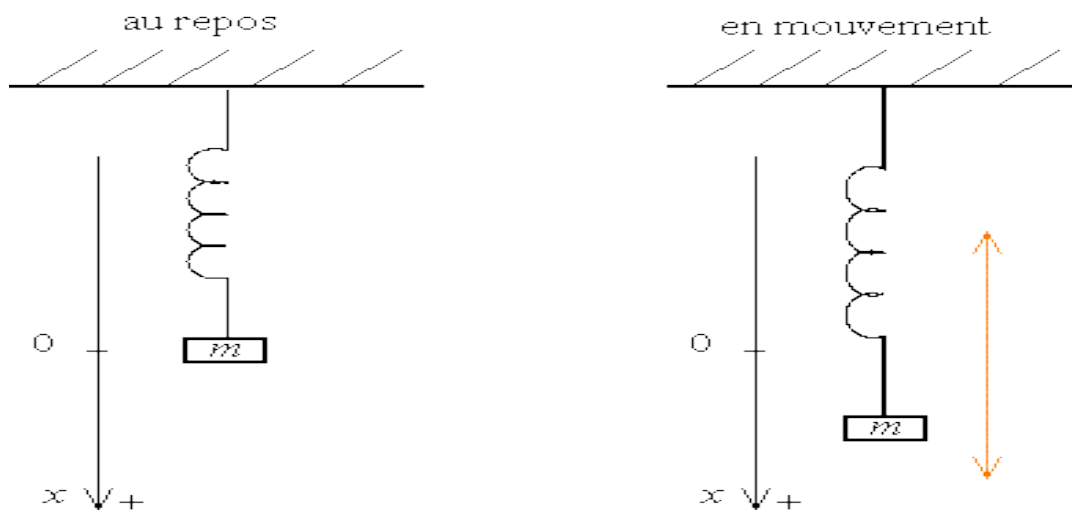
Il vaut +1 (resp. -1) si les points du nuage sont exactement alignés sur une droite de régression.

Ce coefficient est une mesure de la dispersion du nuage.

En pratique,

On estimera souvent la régression acceptable (et ainsi la prédiction acceptable) lorsque $|\rho_{xy}| \geq 0.85$ par analogie la régression sera inacceptable lorsque $|\rho_{xy}| < 0.85$.

Exemple



Dans un TP de physique à chaque masse m_i (un poids) on obtiendra une extension x_i , le tableau des données brutes est le suivant :

| | | | | | |
|---------|---|-----|-----|-----|-----|
| Mi (kg) | 0 | 10 | 20 | 30 | 40 |
| Xi (cm) | 0 | 1.1 | 1.5 | 1.9 | 2.5 |

- 1) Donner la droite de régression de Y en fonction de X : $y=ax+b$. (avec Y la variable poids)
- 2) Donner la droite de régression de X en fonction de Y : $x=a'y+b'$. (avec Y la variable poids)
- 3) Cette droite est elle acceptable ? peut-on faire de la prédiction.
- 4) Si oui, quelle masse aurions nous si on veut une dilatation du ressort égale à $x_i= 3\text{cm}$.

Solution

1.

| | | | | | | Somme |
|----------------------------------|---|------|------|------|------|-------|
| Y :mi | 0 | 10 | 20 | 30 | 40 | 100 |
| X :xi | 0 | 0.5 | 1.1 | 1.5 | 1.9 | 5 |
| mi.xi | 0 | 5 | 22 | 45 | 76 | 148 |
| X ² :xi ² | 0 | 0.25 | 1.21 | 2.25 | 3.61 | 7.32 |
| Y ² : mi ² | 0 | 100 | 400 | 900 | 1600 | 3000 |

On obtient :

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{5} (0 + 0.5 + \dots + 1.9) = \frac{5}{5} = 1$$

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n m_i = \frac{1}{5} (0 + 10 + \dots + 40) = \frac{100}{5} = 20$$

$$\overline{XY} = \frac{1}{n} \sum_{i=1}^n x_i \cdot m_i = \frac{1}{5} (0 + 5 + \dots + 76) = \frac{148}{5} = 29.6$$

$$\overline{X^2} = \frac{1}{n} \sum_{i=1}^n x_i^2 = \frac{1}{5} (0 + 0.25 + \dots + 3.61) = \frac{7.32}{5} = 1.464$$

$$\overline{Y^2} = \frac{1}{n} \sum_{i=1}^n m_i^2 = \frac{1}{5} (0 + 100 + \dots + 1600) = \frac{3000}{5} = 600$$

Ainsi

$$V(X) = \overline{X^2} - (\bar{X})^2 = 1.464 - 1^2 = 0.464 \quad \Rightarrow \quad \sigma_x = \sqrt{V(X)} = 0.681$$

$$V(Y) = \overline{Y^2} - (\bar{Y})^2 = 600 - 20^2 = 200 \quad \Rightarrow \quad \sigma_y = \sqrt{V(Y)} = 14.142$$

$$\text{cov}(X, Y) = \overline{XY} - \bar{X} \cdot \bar{Y} = 29.6 - 1 \times 20 = 9.6$$

En utilisant les formules obtenues avec la méthode des MC, la droite de régression qu'on obtiendra sera :

$$a = \frac{\text{cov}(X, Y)}{V(X)} = \frac{9.6}{0.464} = 20.689$$

$$b = \bar{y} - a \cdot \bar{x} = 20 - 20.689 \times 1 = -0.689$$

Conclusion notre droite de régression (Δ) de y en fonction de x, sera la suivante :

$$\boxed{(\Delta) : y = (20.689)x - 0.689}$$

2. En utilisant les formules obtenues avec la méthode des MC, la droite de régression qu'on obtiendra sera :

$$a' = \frac{\text{cov}(X, Y)}{V(Y)} = \frac{9.6}{200} = 0.048$$
$$b' = \bar{x} - a' \cdot \bar{y} = 1 - 0.048 \times 20 = 0.04$$

Conclusion notre droite de régression (Δ') de x en fonction de y sera la suivante :

$$\boxed{(\Delta') : x = (0.048)y + 0.04}$$

3. pour savoir si la droite de régression est acceptable, calculons le coefficient de corrélation :

$$\rho_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y} = \frac{9.6}{0.681 \times 14.142} = 0.996$$
$$|\rho_{xy}| = 0.996 \geq 0.85 \Rightarrow \text{la droite de régression } (\Delta) \text{ est acceptée}$$

ainsi que les prédictions qui en découleront

4.

Selon la réponse précédente, on peut faire la prédiction du poids si l'extension $x_i = 3\text{cm}$

$$y = (20.689)x - 0.689 \Rightarrow y_i = (20.689)x_i - 0.689 = 20.689 \times 3 - 0.689 = 61.378\text{kg}$$

(Corrigé du TD2 sera sur le site)