

People's Democratic Republic Of Algeria
Ministry Of Higher Education And Scientific Research



Centre Universitaire Abdelhafid Boussouf Mila



Regression linéaire et corrélation linéaire

Octobre , 2022

En statistique, le terme de corrélation est utilisé afin de désigner la liaison entre deux variables quantitatives (le plus souvent continues).

Coefficient de corrélation :

Le coefficient de corrélation permet de mesurer la dépendance linéaire entre deux variables quantitatives X et Y .

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \in [-1, 1].$$

Coefficient de détermination : Le coefficient de détermination est le carré du coefficient de corrélation :

$$\rho_{XY}^2 = \frac{\text{Cov}^2(X, Y)}{V(X)V(Y)}.$$

Remarques :

- X et Y sont indépendantes, alors $\rho = 0$. La réciproque est fautive, sauf cas particulier ; si X et Y sont distribuées normalement.
- Si $\rho > 0$, les valeurs prises par Y ont tendance à croître quand les valeurs de X augmentent.
- Si $\rho < 0$, les valeurs prises par Y ont tendance à décroître quand les valeurs de X augmentent.
- Si $|\rho| = 1$, alors il existe une relation linéaire parfaite entre X et Y .
- Un coefficient de corrélation nul ne signifie pas l'absence de toute relation entre les deux variables. Il peut exister une relation non linéaire entre elles.

Quelques exemples de corrélation (Ajustement linéaire :)

On considère n individus sur lesquels on mesure X et Y deux variables quantitatives. Pour chaque individu i ($1 \leq i \leq n$), on dispose d'un couple d'observations (x_i, y_i) qui représente les valeurs prises par X et Y pour l'individu i .

Les données (x_i, y_i) , $i = 1, \dots, n$ peuvent être représentées par un nuage de n points dans le plan (x, y)

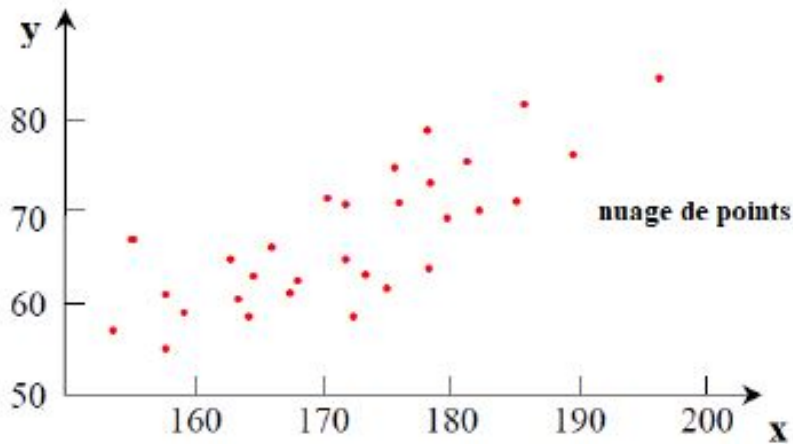


Figure – Nuage de points

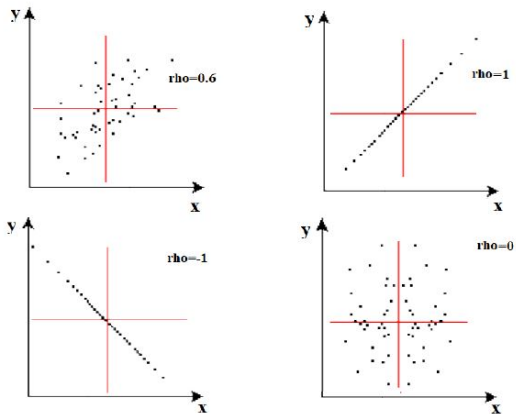


Figure – Nuage de points pour différentes valeurs de ρ

Selon l'allure du nuage, on a envie de remplacer ce nuage par le graphe d'une fonction f . Cette opération s'appelle un ajustement. La nature de l'ajustement dépend de la forme du nuage de points.

Nous étudions l'ajustement linéaire par la méthode des moindres carrés. On considère les données $\{y_1, y_2, \dots, y_n\}$ comme étant des réalisations d'une variable aléatoire Y et les données $\{x_1, x_2, \dots, x_n\}$ les réalisations d'une variable aléatoire X .

- La variable X est une variable, aléatoire ou contrôlée, dite explicative.
- Y est une variable aléatoire dite à expliquer.

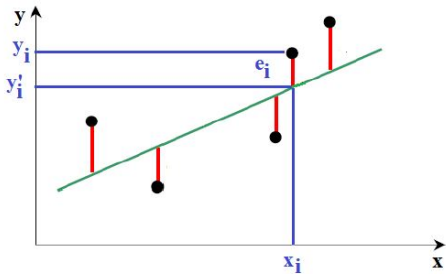


Figure – Nuage de points pour différentes valeurs de ρ

Ajustement linéaire au sens des moindres carrés :)

- Le problème de l'ajustement linéaire par la méthode des moindres carrés consiste à rechercher une relation affine entre les variables X et Y , ceci revient à trouver 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

$$y'_i = ax_i + b.$$

Si les coefficients a et b étaient connus, on pourrait calculer les résidus de la régression définis par :

$$e_i = y_i - ax_i - b.$$

Le résidu e_i est l'erreur que l'on commet en utilisant la droite de régression pour prédire y_i à partir de x_i . Les résidus peuvent être positifs ou négatifs.

Droite de régression au sens des moindres carrés :)

Afin de déterminer la valeur des coefficients a et b nous utilisons le principe des moindres carrés qui consiste à chercher la droite qui minimise la somme des carrés des résidus :

$$\varepsilon(a, b) = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - ax_i - b)^2 .$$

Calculons les coefficients a et b qui minimisent le critère des moindres carrés :

Pour déterminer le minimum la fonction ε , on doit avoir : $\frac{\partial \varepsilon}{\partial a} = 0$ et $\frac{\partial \varepsilon}{\partial b} = 0$.

$$\begin{aligned} \text{On } a \frac{\partial \varepsilon}{\partial b} = 0 &\Leftrightarrow \sum_{i=1}^n (-2)(y_i - ax_i - b) = 0 \Leftrightarrow \sum_{i=1}^n (y_i - ax_i - b) = 0 \\ &\Leftrightarrow \sum_{i=1}^n y_i - a \sum_{i=1}^n x_i - nb = 0 \Leftrightarrow n\bar{y} - an\bar{x} - nb = 0 \Leftrightarrow \bar{y} - a\bar{x} - b = 0 \\ &\Leftrightarrow b = \bar{y} - a\bar{x}. \end{aligned}$$

$$\begin{aligned}
\text{Et } \frac{\partial \varepsilon}{\partial a} = 0 &\Leftrightarrow \sum_{i=1}^n (-2x_i)(y_i - ax_i - b) = 0 \\
&\Leftrightarrow \sum_{i=1}^n x_i (y_i - ax_i - b) = 0 \Leftrightarrow \sum_{i=1}^n x_i y_i - a \sum_{i=1}^n x_i^2 - \sum_{i=1}^n b x_i = 0 \\
&\Leftrightarrow n\overline{xy} - a\overline{x^2} - bn\overline{x} = 0 \Leftrightarrow \overline{xy} - a\overline{x^2} - (\overline{y} - a\overline{x})\overline{x} = 0 \\
&\Leftrightarrow \overline{xy} - a\overline{x^2} - \overline{y}\overline{x} + a\overline{x}^2 = 0 \Leftrightarrow \overline{xy} - \overline{y}\overline{x} = a(\overline{x^2} - \overline{x}^2) \Leftrightarrow a = \frac{\overline{xy} - \overline{y}\overline{x}}{(\overline{x^2} - \overline{x}^2)} \\
&\Leftrightarrow a = \frac{\text{Cov}(X, Y)}{V(X)}. \text{ D'ou la droite de r\u00e9gression de } Y \text{ par rapport \u00e0 } X, \text{ est :}
\end{aligned}$$

$$y = ax + b$$

avec

$$a = \frac{\text{Cov}(X, Y)}{V(X)} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

et.

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

Remarque La droite de régression de X par rapport à Y est :

$$x = a'y + b',$$

avec

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

En pratique, on estime que la régression est acceptable lorsque $|\rho_{XY}| \geq 0.85$.

Exemple : On s'intéresse à la corrélation entre la taille X et le poids Y de

quatre individus :

Individus	Taille (cm)	Poids (kg)
1	180	79
2	170	75
9	175	85
4	160	59

1. Calculer le coefficient de corrélation.
2. Trouver la droite de régression.

1. Les moyennes des variables X et Y .

$$\bar{x} = \frac{1}{4}(180 + 170 + 175 + 160) = 171.25; \bar{y} = \frac{1}{4}(79 + 75 + 85 + 59) = 74.5$$

2. Les variances des variables X et Y .

$$V(X) = \frac{1}{n} \sum_{i=1}^4 x_i^2 - \bar{x}^2 = 54.68; V(Y) = \frac{1}{n} \sum_{i=1}^4 y_i^2 - \bar{y}^2 = 92.75$$

3. La covariance de X et Y :

$$\text{Cov}(X, Y) = \sum_{i=1}^4 x_i y_i - \bar{x} \bar{y} = \overline{xy} - \bar{x} \bar{y} = 63.12.$$

4. Le coefficient de corrélation :

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}} = 0.88.$$

5. La droite de régression de Y en X :

$$a = \frac{\text{Cov}(X, Y)}{V(X)} = 1.15$$

$$b = \bar{y} - a\bar{x} = -122.43.$$

D'où

$$y = ax + b = 1.15x - 122.43$$

Fin du partie 2