

Chapitre IV: Ajustement linéaire - Régression.

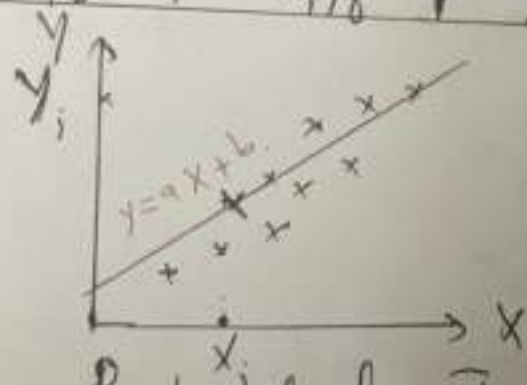
1) Ajustement à l'aide d'une droite: Supposons que nous ayons une population de N individus et que nous désirons l'étudier sur deux caractères simultanément. Notre but étant de déterminer et d'analyser la liaison entre ces deux caractères. Par exemple nous disposons d'un échantillon de 10 personnes sur chacune desquelles nous avons mesuré la taille X et le poids Y , nous pouvons classer les données dans un tableau comme suit:

Individus	Taille X (m)	Poids Y (kg)
1	1,55	58,2
2	1,60	58,1
3	1,62	61,3
4	1,64	61,3
5	1,65	69,5
6	1,70	69,7
7	1,72	70,3
8	1,73	75,4
9	1,76	74,2
10	1,78	82,0

Si nous présentons graphiquement ces données, on obtient un nuage de points qui peut être représenté par une droite d'équation:

$$Y = aX + b.$$

Il est demandé de déterminer a et b pour établir définitivement la relation entre X et Y .



$$\begin{cases} a = \frac{\text{cov}(X, Y)}{\text{Var}(X)} \\ \text{et} \\ b = \bar{Y} - a\bar{X} \end{cases}$$

Exemple de calcul: $\bar{X} = 1,675$ $\bar{Y} = 68,4$.

$$\text{Var}(X) = \frac{1}{N} \sum_{i=1}^N X_i^2 - \bar{X}^2 = 0,005.$$

$$\text{cov}(X, Y) = \frac{1}{N} \sum_{i=1}^N X_i Y_i - \bar{X} \bar{Y} = 0,4343.$$

$$\text{et donc: } \begin{cases} a = \frac{\text{cov}(X, Y)}{\text{Var}(X)} = 86,86 \\ b = \bar{Y} - a\bar{X} = -91,19 \end{cases}$$

La droite d'ajustement est alors $Y = 86,86X - 91,19$.

② Variance résiduelle: la méthode qui nous a permis de définir le paramètre de la droite de régression (a et b), s'appelle la méthode des moindres carrés. Cette méthode nous donne la meilleure droite au sens où elle est la plus proche des points, cela signifie qu'elle minimise la quantité:

$$\sum_{i=1}^N [y_i - (ax_i + b)]^2,$$

ce qui revient au même, celle qui minimise la quantité:

$$\text{Var}_r(Y) = \frac{1}{N} \sum_{i=1}^N [y_i - (ax_i + b)]^2.$$

Cette quantité s'appelle la variance résiduelle. Le calcul de la variance résiduelle dans le cas de la droite se fait par la formule:

$$\text{Var}_r(Y) = \text{Var}(Y) (1 - \rho^2).$$

En effet:
$$\begin{aligned} \text{Var}(Y) &= \frac{1}{N} \sum_{i=1}^N [y_i - a x_i - b]^2 \\ &= \frac{1}{N} \sum_{i=1}^N [y_i - a x_i - \bar{y} + a \bar{x}]^2 \\ &= \frac{1}{N} \sum_{i=1}^N [(y_i - \bar{y}) - a(x_i - \bar{x})]^2. \end{aligned}$$

En développant l'expression entre crochets on obtient:

$$\begin{aligned} \text{Var}_r(Y) &= \frac{1}{N} \sum_{i=1}^N [(y_i - \bar{y})^2 - 2a \sum_{i=1}^N (y_i - \bar{y})(x_i - \bar{x}) + a^2 \sum_{i=1}^N (x_i - \bar{x})^2] \\ &= \text{Var}(Y) - 2a \text{cov}(x, y) + a^2 \text{Var}(X). \end{aligned}$$

$$= \text{Var}(Y) - \frac{2 \text{cov}(x, y)}{\text{Var}(X)} + \frac{\text{cov}(x, y)^2}{\text{Var}(X)} \quad \left(\tan a = \frac{\text{cov}(x, y)}{\text{Var}(X)} \right)$$

$$\text{Var}_r(Y) = \text{Var}(Y) \left[1 - \frac{\text{cov}(x, y)^2}{\text{Var}(X) \text{Var}(Y)} \right] = \text{Var}(Y) [1 - \rho^2].$$