

Statistique: Théorie :

Guillaume DELACOLLETTE

Table des matières

1	Lien entre l'ogive et l'aire sous l'histogramme :	4
1.1	Propriété :	4
1.2	Preuve :	4
2	Effet d'un changement d'origine et d'échelle sur la moyenne :	5
2.1	Propriété :	5
2.2	Preuve :	5
3	Effet d'un changement d'origine et d'échelle sur la médiane :	6
3.1	Propriété :	6
3.2	Preuve :	6
4	Effet d'un changement d'origine et d'échelle sur la variance et l'écart-type :	8
4.1	Propriété :	8
4.2	Preuve :	8
5	Définition et calcul de la moyenne de la série des valeurs centrées :	8
5.1	Propriété :	8
5.2	Preuve :	8
6	Propriété d'optimalité de la moyenne :	9
6.1	Propriété :	9
6.2	Preuve :	9
7	Propriété d'optimalité de la médiane :	10
7.1	Propriété :	10
8	Théorème de König-Huygens :	10
8.1	Propriété :	10
8.2	Preuve :	10
9	Formule de décomposition de la variance pour une série décomposée en k sous-séries :	11
9.1	Propriété :	11
9.2	Preuve :	11
10	Démonstration de la formule équivalente $s^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$:	12
10.1	Preuve :	12
11	Propriété de Tchebychev :	12
11.1	Propriété :	12
11.2	Preuve :	12
12	Effet d'un changement d'origine et d'échelle sur la covariance et la corrélation :	14
12.1	Hypothèse :	14
12.2	Thèse :	14
12.3	Démonstration :	14
13	Formule équivalente $s_{xy} = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}$:	14
13.1	Hypothèse :	14
13.2	Thèse :	14

13.3	Démonstration :	14
14	Démonstration de la forme des moyennes et variance de la série des résidus :	15
14.1	Hypothèse :	15
14.2	Thèse :	15
14.3	Démonstration :	15
14.3.1	La moyenne :	15
14.3.2	La variance :	16
15	Démonstration de la forme des moyennes et variances de la série des valeurs ajustées :	16
15.1	Hypothèse :	16
15.2	Thèse :	16
15.3	Démonstration :	17
15.3.1	Moyenne :	17
15.3.2	Variance :	17
16	Décomposition de la variance :	17
17	Borne sur la covariance :	18
18	Equation de la droite de régression par la technique des moindres carrés :	18

1 Lien entre l'ogive et l'aire sous l'histogramme :

1.1 Propriété :

Soit X une variable continue dont les valeurs sont groupées en J classes et dont la distribution des fréquences est décrite par un histogramme de surface unitaire ainsi que par l'ogive des fréquences cumulées $y = F(x)$. Pour toute valeur x^* , la surface délimitée par l'histogramme et l'axe des abscisses et située à gauche de x^* est égale à l'ordonnée $F(x^*)$.

1.2 Preuve :

Soient les classes $[e_0, e_1], [e_1, e_2], \dots, [e_{J-1}, e_J]$. Notons $A(x^*)$ la surface à gauche de x^* ,

- si $x^* < e_0$, $A(x^*) = 0$ et $F(x^*) = 0$.
- si $x^* \geq e_J$, $A(x^*) = 1$ et $F(x^*) = 1$.
- si $x^* \in [e_{j-1}, e_j]$,

$$A(x^*) = \sum_{k=1}^{j-1} a_k \frac{f_k}{a_k} + (x^* - e_{j-1}) \frac{f_j}{a_j} = F_{j-1} + \frac{f_j}{a_j} (x^* - e_{j-1}).$$

La fonction $F(x)$ est une droite entre les points (e_{j-1}, F_{j-1}) et (e_j, F_j) dont l'équation s'écrit

$$y - F_{j-1} = \frac{F_j - F_{j-1}}{e_j - e_{j-1}} (x - e_{j-1}).$$

Le point $(x^*, F(x^*))$ appartenant à cette droite, on obtient

$$F(x^*) = F_{j-1} + \frac{f_j}{a_j} (x^* - e_{j-1}).$$

La conclusion $F(x^*) = A(x^*)$ est immédiate.

2 Effet d'un changement d'origine et d'échelle sur la moyenne :

2.1 Propriété :

Si un changement d'échelle et d'origine est effectué sur les observations x_1, \dots, x_n pour obtenir la nouvelle série $S' = \{x'_1, \dots, x'_n\}$ avec $x'_i = ax_i + b$ où a et b sont des constantes réelles, alors la moyenne arithmétique \bar{x}' de la série S' est donnée par

$$\bar{x}' = a\bar{x} + b,$$

où \bar{x} est la moyenne arithmétique de $S = \{x_1, \dots, x_n\}$.

2.2 Preuve :

Hypothèse :

- x_1, \dots, x_n n observations de moyenne \bar{x}
- $x'_i = ax_i + b$ où a et b sont des constantes réelles, $a \neq 0$. x'_1, \dots, x'_n a pour moyenne \bar{x}' .

Thèse :

$$\bar{x}' = a\bar{x} + b$$

Démonstration : Par définition,

$$\begin{aligned}\bar{x}' &= \frac{1}{n} \sum_{i=1}^n x'_i \\ &= \frac{1}{n} \sum_{i=1}^n (ax_i + b) \\ &= \frac{1}{n} \left(\sum_{i=1}^n ax_i + \sum_{i=1}^n b \right) \\ &= \frac{1}{n} \left(a \sum_{i=1}^n x_i + nb \right) \\ &= a \frac{1}{n} \sum_{i=1}^n x_i + b \\ &= a\bar{x} + b\end{aligned}$$

3 Effet d'un changement d'origine et d'échelle sur la médiane :

3.1 Propriété :

Si un changement d'échelle et d'origine est effectué sur les observations x_1, \dots, x_n pour obtenir la nouvelle série $S' = \{x'_1, \dots, x'_n\}$ avec $x'_i = ax_i + b$ où a et b sont des constantes réelles, alors la médiane \tilde{x}' de la série S' est donnée par

$$\tilde{x}' = a\tilde{x} + b,$$

où \tilde{x} est la médiane de $S = \{x_1, \dots, x_n\}$.

3.2 Preuve :

La transformation va soit conserver soit inverser l'ordre des observations de départ. La médiane étant basée sur les observations centrales sera définie à partir des versions transformées de ces observations centrales.

Hypothèse :

- x_1, \dots, x_n n observations de médiane \tilde{x}
- $x'_i = ax_i + b$ où a et b sont des constantes réelles, $a \neq 0$. x'_1, \dots, x'_n a pour médiane \tilde{x}' .

Thèse :

$$\tilde{x}' = a\tilde{x} + b$$

Démonstration :

A partir de la série $\{x_1, \dots, x_n\}$, on a

$$x_{(1)} \leq \dots \leq x_{(j)} \leq x_{(j+1)} \leq \dots \leq x_{(n)}$$

D'où, pour tout $a > 0$ et tout b

$$x_{(1)} \leq \dots \leq x_{(j)} \leq x_{(j+1)} \leq \dots \leq x_{(n)}$$

$$ax_{(1)} \leq \dots \leq ax_{(j)} \leq ax_{(j+1)} \leq \dots \leq ax_{(n)}$$

$$ax_{(1)} + b \leq \dots \leq ax_{(j)} + b \leq ax_{(j+1)} + b \leq \dots \leq ax_{(n)} + b$$

On a donc $x'_{(j)} = ax_{(j)} + b$ pour tout j .

En particulier :

Si $n = 2k + 1$, l'observation centrale, $x'_{(k+1)}$, vérifie

$$\tilde{x} = x'_{(k+1)} = ax_{(k+1)} + b \Rightarrow \tilde{x}' = a\tilde{x} + b$$

Et si $n = 2k$, les deux observations du milieu, $x'_{(k)}$ et $x'_{(k+1)}$, sont telles que

$$\tilde{x}' = \frac{x'_{(k)} + x'_{(k+1)}}{2} = \frac{ax_{(k)} + b + ax_{(k+1)} + b}{2} = a \frac{x_{(k)} + x_{(k+1)}}{2} + b = a\tilde{x} + b$$

Et si $a < 0$
A partir de

$$x_{(1)} \leq \dots \leq x_{(j)} \leq x_{(j+1)} \leq \dots \leq x_{(n)}$$

Il vient

$$ax_{(1)} \geq \dots \geq ax_{(j)} \geq ax_{(j+1)} \geq \dots \geq ax_{(n)}$$

Ou encore

$$ax_{(1)} + b \geq \dots \geq ax_{(j)} + b \geq ax_{(j+1)} + b \geq \dots \geq ax_{(n)} + b$$

L'ordre est inversé : $x'_{(j)} = ax_{(n-j+1)} + b$ mais le milieu reste au milieu :
Si $n = 2k + 1$

$$ax_{(1)} + b \geq \dots \geq \underbrace{ax_{(k+1)} + b}_{= x'_{(k+1)}} \geq \dots \geq ax_{(n)} + b$$

Si $n = 2k$

$$ax_{(1)} + b \geq \dots \geq \underbrace{ax_{(k)} + b}_{= x'_{(k)}} \geq \underbrace{ax_{(k+1)} + b}_{= x'_{(k+1)}} \geq \dots \geq ax_{(n)} + b$$

4 Effet d'un changement d'origine et d'échelle sur la variance et l'écart-type :

4.1 Propriété :

Si on effectue un changement d'échelle sur les observations initiales x_1, \dots, x_n pour obtenir la série $S' = \{x'_1, \dots, x'_n\}$, avec $x'_i = ax_i + b$, $a, b \in \mathbb{R}$, alors la variance s'^2 de S' est donnée par $s'^2 = a^2 s^2$ où s^2 est la variance de S . Pour les écarts-types, la relation devient $s' = |a| s$. Les changements d'origine n'ont donc aucun effet sur la variance et l'écart-type.

4.2 Preuve :

Une propriété similaire a été établie que la moyenne de la nouvelle série suit la même transformation que les données : $\bar{x}' = a\bar{x} + b$. Dès lors, par définition de la variance, on a

$$s'^2 = \frac{1}{n} \sum_{i=1}^n (x'_i - \bar{x}')^2 = \frac{1}{n} \sum_{i=1}^n (ax_i + b - a\bar{x} - b)^2 = a^2 s^2$$

5 Définition et calcul de la moyenne de la série des valeurs centrées :

5.1 Propriété :

Définissons la série des valeurs centrées S_c en prenant les différences entre les valeurs observées de la série $S = \{x_1, \dots, x_n\}$ et la moyenne arithmétique de la série S . Cette nouvelle série $S_c = \{x_1 - \bar{x}, \dots, x_n - \bar{x}\}$ a une moyenne nulle.

5.2 Preuve :

Il s'agit d'un cas particulier de la propriété 1.

La propriété 1 nous dit ceci, Si un changement d'échelle et d'origine est effectué sur les observations x_1, \dots, x_n pour obtenir la nouvelle série $S' = \{x'_1, \dots, x'_n\}$ avec $x'_i = ax_i + b$ où a et b sont des constantes réelles, alors la moyenne arithmétique \bar{x}' de la série S' est donnée par

$$\bar{x}' = a\bar{x} + b,$$

où \bar{x} est la moyenne arithmétique de $S = \{x_1, \dots, x_n\}$.

En prenant $a = 1$ et $b = -\bar{x}$ dans la propriété 1, on obtient directement $\bar{x}_c = \bar{x} - \bar{x}$.

6 Propriété d'optimalité de la moyenne :

6.1 Propriété :

La somme des carrés des écarts des éléments d'une série par rapport à la moyenne arithmétique de la série est inférieure ou égale à la somme des carrés des écarts par rapport à toute autre valeur $a \in \mathbb{R}$. Cette propriété s'écrit aussi

$$\sum_{i=1}^n (x_i - \bar{x})^2 \leq \sum_{i=1}^n (x_i - a)^2 \quad \forall a \in \mathbb{R}.$$

6.2 Preuve :

Il suffit de rechercher le minimum sur \mathbb{R} de la fonction $f(a) = \sum_{i=1}^n (x_i - a)^2$.

Hypothèse :

- x_1, \dots, x_n de moyenne \bar{x}

Thèse :

$$\sum_{i=1}^n (x_i - \bar{x})^2 \leq \sum_{i=1}^n (x_i - a)^2 \quad \forall a \in \mathbb{R}$$

où \bar{x} est la valeur qui minimise la fonction $f(a) = \sum_{i=1}^n (x_i - a)^2$.

Démonstration :

Recherchons le minimum de $f(a)$

Point Stationnaire :

$$\begin{aligned} f'(a) &= \sum_{i=1}^n 2(x_i - a)(-1) = 0 \\ \Leftrightarrow \cancel{2} \sum_{i=1}^n (x_i - a) &= 0 \\ \Leftrightarrow \sum_{i=1}^n x_i - na &= 0 \\ \Leftrightarrow a &= \bar{x} \end{aligned}$$

et $f''(a) = 0 \Leftrightarrow \bar{x}$ est le minimum.

7 Propriété d'optimalité de la médiane :

7.1 Propriété :

La somme des écarts absolus des éléments d'une série par rapport à la médiane de la série est inférieure ou égale à la somme des écarts absolus par rapport à toute autre valeur $a \in \mathbb{R}$. Cette propriété s'écrit aussi

$$\sum_{i=1}^n |x_i - \tilde{x}| \leq \sum_{i=1}^n |x_i - a| \quad \forall a \in \mathbb{R}$$

8 Théorème de König-Huygens :

8.1 Propriété :

La moyenne des carrés des écarts entre les observations d'une série et un paramètre a se décompose de la façon suivante :

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

où s^2 est la variance de la série.

8.2 Preuve :

Hypothèse :

x_1, \dots, x_n de moyenne \bar{x} et de variance s^2 .

Thèse :

$$\frac{1}{n} \sum_{i=1}^n (x_i - a)^2 = s^2 + (\bar{x} - a)^2, \quad \forall a \in \mathbb{R}$$

Démonstration :

$$\begin{aligned} s^2(a) &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - a)^2 \\ &= \frac{1}{n} \sum_{i=1}^n \left((x_i - \bar{x}) + (\bar{x} - a) \right)^2 \\ &= \frac{1}{n} \sum_{i=1}^n \left((x_i - \bar{x})^2 + (\bar{x} - a)^2 + 2(x_i - \bar{x})(\bar{x} - a) \right) \\ &= \underbrace{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}_{\text{variance}} + \underbrace{\frac{1}{n} \sum_{i=1}^n (\bar{x} - a)^2}_{\text{indépendant}} + \underbrace{\frac{1}{n} \sum_{i=1}^n 2(x_i - \bar{x})(\bar{x} - a)}_{2(\bar{x} - a) \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})} \\ &= s^2 \quad (\bar{x} - a)^2 \quad \text{moyenne de la série centrée} \\ &= s^2 \quad (\bar{x} - a)^2 \quad \text{qui est égal à 0} \\ &= s^2 + (\bar{x} - a)^2 \end{aligned}$$

9 Formule de décomposition de la variance pour une série décomposée en k sous-séries :

9.1 Propriété :

Soit une population P de n individus décomposée en k sous-populations P_1, \dots, P_k d'effectifs n_1, \dots, n_k (avec $\sum_{i=1}^k n_i = n$), de moyennes $\bar{x}_1, \dots, \bar{x}_k$ et de variances s_1^2, \dots, s_k^2 . La variance globale s^2 de la population est donnée par

$$s^2 = \frac{\sum_{i=1}^k n_i s_i^2}{n} + \frac{\sum_{i=1}^k n_i (\bar{x}_i - \bar{x})^2}{n},$$

où \bar{x} est la moyenne globale de P.

9.2 Preuve :

Hypothèse :

- P de n individu
- moyenne \bar{x} et variance s^2
- k sous population et effectifs n_1, \dots, n_k avec $n_1 + \dots + n_k = n$
- moyenne "local" $\bar{x}_1, \dots, \bar{x}_k$
- variance "local" s_1^2, \dots, s_k^2

Thèse :

$$s^2 = \frac{1}{n} \sum_{i=1}^k n_i s_i^2 + \frac{1}{n} \sum_{i=1}^k n_i (\bar{x}_i - \bar{x})^2$$

Démonstration :

Par définition, $s^2 = \frac{1}{n} \sum_{i=1}^k (x_i - \bar{x})^2$

Soit : $I = \{1, \dots, n\}$

Soit : $I_j = \{i \in I, x_i \in P_j\}$, on a $I_j \cap I_k = \emptyset$ et $\cup_j I_j = I$

$$\rightarrow s^2 = \frac{1}{n} \sum_{j=1}^k \sum_{i \in I_j} (x_i - \bar{x})^2$$

Plaçons nous dans P_j (\rightarrow jème population) et n_j individus, moyenne \bar{x}_j , variance s_j^2 .
Application de TKH :

$$s_j^2(a) = \frac{1}{n_j} \sum_{i \in I_j} (x_i - a)^2 = s_j^2 + (\bar{x}_j - a)^2$$

vrai $\forall a \in \mathbb{R}$ (en particulier pour $a = \bar{x}$).

$$\begin{aligned}
\rightarrow \frac{1}{n_j} \sum_{i \in I_j} (x_i - \bar{x})^2 &= s_j^2 + (\bar{x}_j - \bar{x})^2 \\
\sum_{i \in I_j} (x_i - \bar{x})^2 &= n_j s_j^2 + n_j (\bar{x}_j - \bar{x})^2 \\
\Rightarrow s^2 &= \frac{1}{n} \sum_{j=1}^k (n_j s_j^2 + n_j (\bar{x}_j - \bar{x})^2) \\
\Rightarrow s^2 &= \frac{1}{n} \sum_{j=1}^k n_j s_j^2 + \frac{1}{n} \sum_{j=1}^k n_j (\bar{x}_j - \bar{x})^2
\end{aligned}$$

10 Démonstration de la formule équivalente $s^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$:

10.1 Preuve :

Le théorème de König-Huygens nous donne la relation suivante : La moyenne des carrés des écarts entre les observations d'une série et un paramètre a se décompose de la façon suivante :

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

où s^2 est la variance de la série.

En prenant $a = 0$ dans la relation précédente, il vient

$$\frac{1}{n} \sum_{i=1}^n x_i^2 = s^2 + \bar{x}^2 \Rightarrow s^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$$

11 Propriété de Tchebychev :

11.1 Propriété :

Soit $S = \{x_1, \dots, x_n\}$ une série de moyenne \bar{x} et d'écart-type s . La proportion d'observations s'écartant d'au moins t écarts-types de la moyenne est inférieure ou égale à $\frac{1}{t^2}$.

11.2 Preuve :

Hypothèse :

x_1, \dots, x_n de moyenne \bar{x} et d'écart-type s et $t > 0$.

Thèse :

$$\frac{k}{n} \leq \frac{1}{t^2}$$

Démonstration :

Toute observation x_i est soit tel que $x_i \in]\bar{x} - ts, \bar{x} + ts[$, soit tel que $x_i \notin]\bar{x} - ts, \bar{x} + ts[$
 x_i est dans l'intervalle si et seulement si $|x_i - \bar{x}| < ts$
 x_i est en dehors de l'intervalle si $|x_i - \bar{x}| \geq ts$

Notons I comme l'ensemble de tous les indices :

$$I = \{1, 2, \dots, n\}$$

$$I_1 = \{i \in I : |x_i - \bar{x}| < ts\}$$

$$I_2 = \{i \in I : |x_i - \bar{x}| \geq ts\}$$

$$I_1 \cap I_2 = \emptyset$$

$$I_1 \cup I_2 = I$$

Par définition,

$$\begin{aligned} s^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= \frac{1}{n} \sum_{i \in I_1} (x_i - \bar{x})^2 + \frac{1}{n} \sum_{i \in I_2} (x_i - \bar{x})^2 \end{aligned}$$

C'est I_2 qui nous intéresse

$$\begin{aligned} &\frac{1}{n} \sum_{i \in I_2} \underbrace{(x_i - \bar{x})^2}_{\geq t^2 s^2} \\ &\geq t^2 s^2 \end{aligned}$$

$$\begin{aligned} \frac{1}{n} \sum_{i \in I_2} (x_i - \bar{x})^2 &\geq \frac{1}{n} \underbrace{\sum_{i \in I_2} t^2 s^2}_{= \frac{kt^2 s^2}{n}} \\ &= \frac{kt^2 s^2}{n} \end{aligned}$$

k est le nombre d'indices appartenant à I_2 ou encore le nombre d'observations s'écartant d'au moins t écarts-types de la moyenne.

En réarrangeant les termes, on obtient,

$$\begin{aligned} \frac{k}{n} &\geq \frac{t^2}{s^2} \\ \Rightarrow \frac{k}{n} &< \frac{1}{t^2} \end{aligned}$$

CQFD

12 Effet d'un changement d'origine et d'échelle sur la covariance et la corrélation :

12.1 Hypothèse :

$$\begin{aligned} (x_i, y_i) \quad 1 \leq i \leq n \text{ de covariance } s_{xy} \\ x_i' = ax_i + b \quad \bar{x}' = a\bar{x} + b \\ y_i' = cy_i + d \quad \bar{y}' = c\bar{y} + d \end{aligned}$$

12.2 Thèse :

$$s'_{xy} = acs_{xy}$$

12.3 Démonstration :

Par définition,

$$\begin{aligned} s'_{xy} &= \frac{1}{n} \sum_{i=1}^n (x'_i - \bar{x}')(y'_i - \bar{y}') \\ &= \frac{1}{n} \sum_{i=1}^n (ax_i + b - a\bar{x} - b)(cy_i + d - c\bar{y} - d) \\ &= \frac{1}{n} \sum_{i=1}^n a(x_i - \bar{x})c(y_i - \bar{y}) \\ &= acs_{xy} \end{aligned}$$

13 Formule équivalente $s_{xy} = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}$:

13.1 Hypothèse :

$$\begin{aligned} (x_i, y_i) \quad 1 \leq i \leq n \\ \text{moyenne } \bar{x} \quad \text{moyenne } \bar{y} \\ \text{covariance } s_{xy} \end{aligned}$$

13.2 Thèse :

$$s_{xy} = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}$$

13.3 Démonstration :

Par définition,

$$\begin{aligned} s_{xy} &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ &= \frac{1}{n} \sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n x_i \bar{y} - \frac{1}{n} \sum_{i=1}^n \bar{x} y_i + \frac{1}{n} \sum_{i=1}^n \bar{x} \bar{y} \\ &= \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y} - \bar{x} \bar{y} + \bar{x} \bar{y} \end{aligned}$$

14 Démonstration de la forme des moyennes et variance de la série des résidus :

14.1 Hypothèse :

$$(x_i, y_i) \quad 1 \leq i \leq n$$

droite de régression par les moindres carrés

$$y = \hat{a}x + \hat{b} \text{ avec } \hat{a} = \frac{s_{xy}}{s_x^2} \text{ et } \hat{b} = \bar{y} - \hat{a}\bar{x}$$

$$\delta_i = y_i - \hat{y}_i$$

14.2 Thèse :

La moyenne de la série des résidus est nulle : $\bar{\delta} = 0$

La variance (appelée variance résiduelle de y par rapport à x), vaut $s_{\delta}^2 = s_y^2(1 - r^2)$ où $r = r_{xy}$

14.3 Démonstration :

14.3.1 La moyenne :

$$\begin{aligned} \bar{\delta} &= \frac{1}{n} \sum_{i=1}^n \delta_i \\ &= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \\ &= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{a}x_i - \hat{b}) \\ &= \frac{1}{n} \sum_{i=1}^n y_i - \frac{1}{n} \sum_{i=1}^n \hat{a}x_i - \frac{1}{n} \sum_{i=1}^n \hat{b} \\ &= \bar{y} - \hat{a}\bar{x} - \hat{b} \\ &\rightarrow \text{par hypothèse } \hat{b} = \bar{y} - \hat{a}\bar{x} \\ &= \bar{y} - \hat{a}\bar{x} - \bar{y} + \hat{a}\bar{x} \\ &= 0 \end{aligned}$$

14.3.2 La variance :

$$\begin{aligned}
s_{\delta}^2 &= \frac{1}{n} \sum_{i=1}^n (\delta_i - \bar{\delta})^2 \\
&= \frac{1}{n} \sum_{i=1}^n \delta_i^2 \\
&= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\
&= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{a}x_i - \hat{b})^2 \\
&\rightarrow \text{par hypothèse } \hat{b} = \bar{y} - \hat{a}\bar{x} \\
&= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y} - \hat{a}(x_i - \bar{x}))^2 \\
&= \frac{1}{n} \sum_{i=1}^n \{(y_i - \bar{y})^2 - 2(y_i - \bar{y})\hat{a}(x_i - \bar{x}) + \hat{a}^2(x_i - \bar{x})^2\} \\
&= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 - \frac{1}{n} \sum_{i=1}^n 2(y_i - \bar{y})\hat{a}(x_i - \bar{x}) + \frac{1}{n} \sum_{i=1}^n \hat{a}^2(x_i - \bar{x})^2 \\
&= s_y^2 - 2\hat{a}s_{xy} + \hat{a}^2 s_x^2 \\
&\rightarrow \text{par hypothèse } \hat{a} = \frac{s_{xy}}{s_x^2} \\
&= s_y^2 - 2\frac{s_{xy}}{s_x^2} s_{xy} + \frac{s_{xy}^2}{s_x^4} s_x^2 \\
&= s_y^2 - \frac{s_{xy}^2}{s_x^2} \\
&= s_y^2 \left(1 - \frac{s_{xy}^2}{s_x^2 s_y^2}\right) \\
&= s_y^2 (1 - r_{xy}^2)
\end{aligned}$$

15 Démonstration de la forme des moyennes et variances de la série des valeurs ajustées :

15.1 Hypothèse :

$(x_i, y_i) \quad 1 \leq i \leq n$

droite de régression par les moindres carrés

$y = \hat{a}x + \hat{b}$ avec $\hat{a} = \frac{s_{xy}}{s_x^2}$ et $\hat{b} = \bar{y} - \hat{a}\bar{x}$

$\delta_i = y_i - \hat{y}_i$

15.2 Thèse :

$$s_{\hat{y}}^2 = r^2 s_y^2 \text{ et } \bar{y} = \bar{\hat{y}}$$

15.3 Démonstration :

15.3.1 Moyenne :

$$\begin{aligned}0 = \bar{\delta} &= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \\&= \frac{1}{n} \sum_{i=1}^n y_i - \frac{1}{n} \sum_{i=1}^n \hat{y}_i \\&\Rightarrow \overline{\hat{y}_i} = \bar{y}\end{aligned}$$

15.3.2 Variance :

$$\begin{aligned}s_{\hat{y}}^2 &= \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - \overline{\hat{y}})^2 \\&= \frac{1}{n} \sum_{i=1}^n (\hat{a}x_i + \hat{b} - \bar{y})^2 \\&\rightarrow \text{par hypothèse } \hat{b} = \bar{y} - \hat{a}\bar{x} \\&= \frac{1}{n} \sum_{i=1}^n (\hat{a}x_i - \hat{a}\bar{x})^2 \\&= \hat{a}^2 \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \\&= \hat{a}^2 s_x^2 \\&\rightarrow \text{par hypothèse } \hat{a} = \frac{S_{xy}}{S_x^2} \\&= \frac{S_{xy}^2}{(S_x^2)^2} S_x^2 \\&= \frac{S_{xy}^2}{S_x^2} \\&= r^2 S_y^2 \\&\rightarrow \text{On a, } r^2 = \frac{S_{xy}^2}{S_x^2 S_y^2} \\&\Leftrightarrow r^2 S_y^2 = \frac{S_{xy}^2}{S_x^2}\end{aligned}$$

16 Décomposition de la variance :

$$\begin{aligned}s_y^2 &= 1s_y^2 \\&= (1 - r^2 + r^2)s_y^2 \\&= (1 - r^2)s_y^2 + r^2s_y^2 \\&= s_{\delta}^2 + s_{\hat{y}}^2\end{aligned}$$

17 Borne sur la covariance :

La covariance de la série double S est toujours, en valeur absolue, inférieure ou égale au produit des écarts-types marginaux s_x, s_y : $|s_{xy}| \leq s_x s_y$. L'égalité n'est possible que si et seulement si tous les points observés sont situés sur une même droite.

18 Equation de la droite de régression par la technique des moindres carrés :

Soit $S = \{(x_i, y_i), 1 \leq i \leq n\}$ une série statistique (quantitative) bivariée. L'équation de la droite de régression obtenue par la méthode des moindres carrés est donnée par

$$y = \hat{a}x + \hat{b} \quad \text{où} \quad \hat{a} = \frac{s_{xy}}{s_x^2} \quad \text{et} \quad \hat{b} = \bar{y} - \hat{a}\bar{x}$$