

Chapitre 5 : Statistiques descriptives

Table des matières

Chapitre 5 : Statistiques descriptives	1
Axel CARPENTIER	
1 Série statistique à une variable	3
1.1 Méthodes de représentation	3
1.2 Caractéristiques de positions	5
1.3 Caractéristiques de dispersion	8
2 Série statistique à deux variables	8
2.1 Ajustement affine	8
2.2 Méthode des moindres carrés	10
2.3 Coefficient de corrélation	12

1 Série statistique à une variable

Dans toute cette partie, on considèrera les 3 séries statistiques suivantes :

Série A :

Notes obtenues à un contrôle dans une classe de 30 élèves :

2 – 3 – 3 – 4 – 5 – 6 – 6 – 7 – 7 – 7 – 8 – 8 – 8 – 8 – 8 – 9 – 9 – 9 – 9 – 9 – 10 – 10 – 11 – 11 – 11 – 13 – 13 – 15 – 16

Série B :

Salaires en euros des employés d'une entreprise :

Salaires	[900; 1200]	[1200; 1400]	[1400; 1600]	[1600; 1800]	[1800; 2000]	[2000; 2400]	TOTAL
Effectif	30	30	60	80	40	40	280

Série C :

Proportion d'adhérents à un club sportif dans différentes sections :

- 17% jouent au handball,
- 25% jouent au rugby,
- 58% jouent au tennis.

1.1 Méthodes de représentation

1.1.1 Vocabulaire

La **population** est l'ensemble des individus sur lesquels portent l'étude statistique. (Par exemple la classe de BTS domotique, la population féminine, les fonctionnaires ...) dont chaque élément est appelé **individu**.

Un **échantillon** est une partie de la population considérée.

Le **caractère** (ou **variable**) d'une série statistique est une propriété étudiée sur chaque individu :

- Lorsque le caractère ne prend que des valeurs (ou **modalités**) numériques, il est **quantitatif** :
 - **discret** s'il ne peut prendre que des valeurs isolées (notes, âge ...)
 - **continu** dans le cas contraire (poids, taille ...). Dans ce cas on effectue souvent un regroupement des valeurs par **classes**.
- Sinon, on dit qu'il est **qualitatif** (couleur des yeux, sport pratiqué ...): les modalités ne sont pas des nombres.

A chaque valeur (ou classe) est associée un **effectif** n : c'est le nombre d'individus associés à cette valeur.

Faire des **statistiques**, c'est recueillir, organiser, synthétiser, représenter et exploiter des données, numériques ou non, dans un but de comparaison, de prévision, de constat ...

Les plus gros "consommateurs" de statistiques sont les assureurs (risques d'accidents, de maladie des assurés), les médecins (épidémiologie), les démographes (populations et leur dynamique), les économistes (emploi, conjoncture économique), les météorologues ...

1.1.2 Tableaux

Définition:

On considère une série statistique X à caractère quantitatif, dont les p valeurs sont données par x_1, x_2, \dots, x_p d'effectifs associés n_1, n_2, \dots, n_p avec $n_1 + n_2 + \dots + n_p = N$.

- A chaque valeur (ou classe) est associée une fréquence f_i : c'est la proportion d'individus associés à cette valeur.
- $f_i = \frac{n_i}{N}$ est un nombre compris entre 0 et 1, que l'on peut écrire sous forme de pourcentage.

Exemple:

On peut représenter la **série A** par un tableau d'effectifs, et le compléter par la distribution des fréquences :

Notes	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
Eff.	0	1	2	1	1	2	3	5	6	2	3	0	2	0	1	1	0	0	0
Fréq. en %	0	3	7	3	3	7	10	17	20	7	10	0	7	0	3	3	0	0	0

! Remarque

On peut vérifier que la somme des fréquences est égale à 1 (ou à 100 si on les exprime en pourcentages).

On peut aussi faire un regroupement par classe, ce qui rend l'étude moins précise, mais qui permet d'avoir une vision plus globale.

Exemple:

Toujours pour la **série A**, si on regroupe les données par classes d'amplitude 5 points, on obtient :

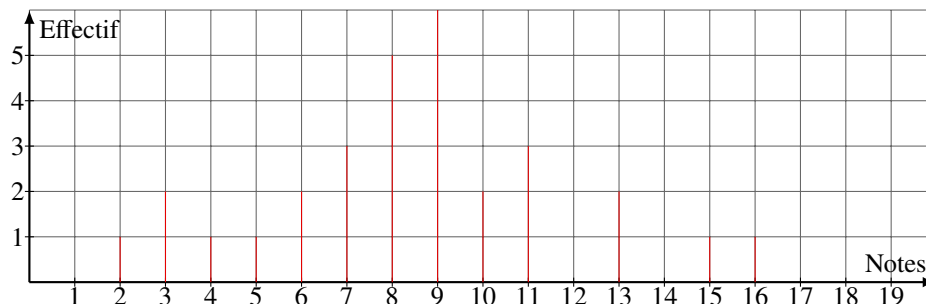
Notes	[0 ; 5 [[5 ; 10 [[10 ; 15 [[15 ; 20 [total
Effectif	4	17	7	2	30
Fréquence	0,13	0,57	0,23	0,07	1

1.1.3 Graphiques

Lorsque le caractère étudié est **quantitatif et discret**, on peut représenter la série statistique étudiée par un **diagramme en bâtons**: la hauteur de chaque bâton est alors proportionnelle à l'effectif (ou à la fréquence) associé à chaque valeur.

Exemple:

Voici le diagramme en bâtons représentant la série des notes de la **série A** :

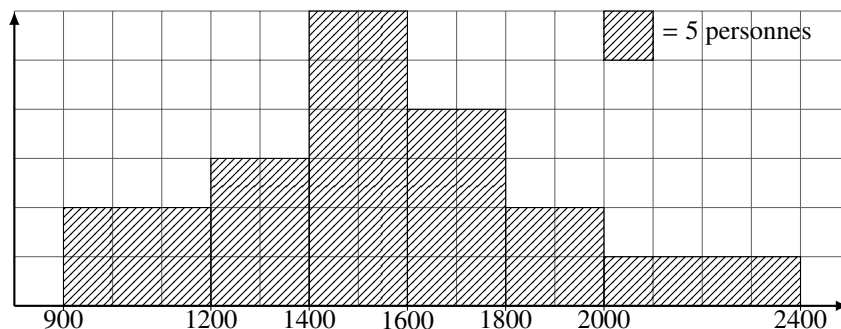


Lorsque le caractère étudié est **quantitatif et continu**, et lorsque les modalités sont regroupées en classes, on peut représenter la série par un **histogramme**: l'aire de chaque rectangle est alors proportionnelle à l'effectif (ou à la fréquence) associée à chaque classe.

Lorsque les classes ont la même **amplitude**, c'est la hauteur qui est proportionnelle à l'effectif.

Exemple:

Pour la **série B**, on obtient par exemple l'histogramme suivant :



Enfin, lorsque le caractère est **qualitatif**, on peut représenter la série par :

- **Un diagramme circulaire** ("camemberts") :

La mesure de chaque secteur angulaire est proportionnelle à l'effectif associé.

- **Un diagramme en tuyaux d'orgue :**

Chaque classe est représentée par un rectangle de même largeur et de longueur proportionnelle à l'effectif, donc à la fréquence.

- **Un diagramme en bandes :**

Chaque classe est représentée par un rectangle de même largeur et de longueur proportionnelle à l'effectif, donc à la fréquence.

Exemple:

Diagrammes de la **série C**.

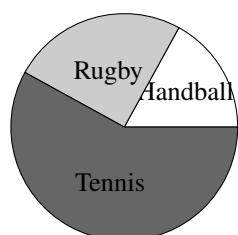


Diagramme circulaire

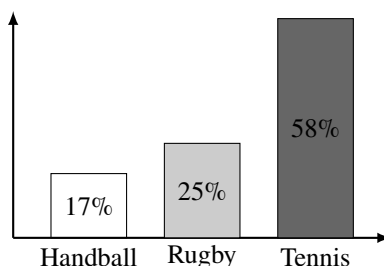


Diagramme en tuyau d'orgue

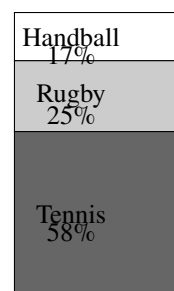


Diagramme en bande

1.2 Caractéristiques de positions

1.2.1 Moyenne

Définition:

Si une série statistique $x_1 ; \dots ; x_p$ prend n_i fois la valeur x_i , pour $i \in \{1, \dots, p\}$, alors on a la moyenne pondérée qui vaut

$$\bar{x} = \frac{n_1 \times x_1 + \dots + n_p \times x_p}{n_1 + \dots + n_p}$$

Exemple:

- Dans la **série A**, la moyenne du contrôle est égale à $\bar{m} = \frac{254}{30} \approx 8,47$.
- Dans la **série B**, une estimation du salaire moyen est donné par : $\bar{S} = \frac{460500}{280} \approx 1644,64$.

! Remarque

On peut aussi calculer une moyenne à partir de la distribution de fréquences : $\bar{x} = f_1 x_1 + f_2 x_2 + \dots + f_p x_p$

Propriété: linéarité de la moyenne

- Si on ajoute (ou soustrait) un même nombre k à toutes les valeurs d'une série, alors la moyenne de cette série se trouve augmentée (resp. diminuée) de k .
- Si on multiplie (ou divise) par un même nombre non nul k toutes les valeurs d'une série, alors la moyenne de cette série se trouve multipliée (resp. divisée) par k .

Démonstration:

$$\overline{x+k} = \frac{n_1 \times (x_1 + k) + \dots + n_p \times (x_p + k)}{n_1 + \dots + n_p} = \frac{n_1 \times x_1 + \dots + n_p \times x_p}{n_1 + \dots + n_p} + \frac{n_1 \times k + \dots + n_p \times k}{n_1 + \dots + n_p} = \bar{x} + k$$

Exemple:

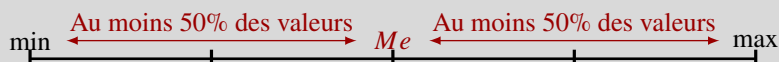
- Si on ajoute 1,5 points à chaque note du contrôle, alors la moyenne de classe devient $\bar{m} = 8,47 + 1,5 = 9,97$.
- Si on augmente chaque note de 10%, cela revient à multiplier chaque note par 1,1, ce qui donne $\bar{m} = 8,47 \times 1,1 = 9,32$.

1.2.2 Médiane et quartiles

Définition:

On appelle médiane Me d'une série statistique à n éléments, la plus petite valeur de la série qui est supérieure ou égale à 50% des valeurs de la série.

- Si n est impair alors la médiane est la valeur centrale.
- Si n est pair alors la médiane est la moyenne des deux valeurs centrales.



! Remarque

La médiane peut ne pas être égale à une des valeurs de la série statistique.

Exemple:

On relève les pointures de 10 personnes et on obtient :

45 ; 39 ; 42 ; 48 ; 41 ; 42 ; 43 ; 42 ; 38 ; 49

On réorganise ces données dans l'ordre croissant et on a donc la médiane : 42.

! Remarque

La médiane n'est pas affectée par les valeurs extrêmes.

Si on remplace le 49 par 54 dans l'exemple précédent le résultat ne change pas.

Exemple:

On souhaite calculer la médiane de la **série A**.

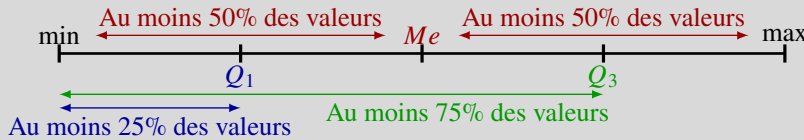
- Pour cela, on commence par remplir le tableau des effectifs cumulés croissants :

Notes	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15	16	17	18	19
Eff.	0	1	2	1	1	2	3	5	6	2	3	0	2	0	1	1	0	0	0
ECC.	0	1	3	4	5	7	10	15	21	23	26	26	28	28	29	30	30	30	30

- Ensuite, l'effectif étant de 30, on choisit la moyenne entre la 15^{ème} et la 16^{ème} note.
On obtient $Med = \frac{8+9}{2} = 8,5$.
- Ce qui signifie que la moitié des notes est inférieure ou égale à 8,5, et que l'autre moitié des notes est supérieure ou égale à 8,5.

Définition:

On appelle 1^{er} quartile Q_1 (respectivement 3^{ème} quartile Q_3) la plus petite valeur de la série qui est supérieure ou égale à 25% (respectivement 75%) des valeurs de la série.



Exemple:

On reprends l'exemple avec les pointures des 10 personnes

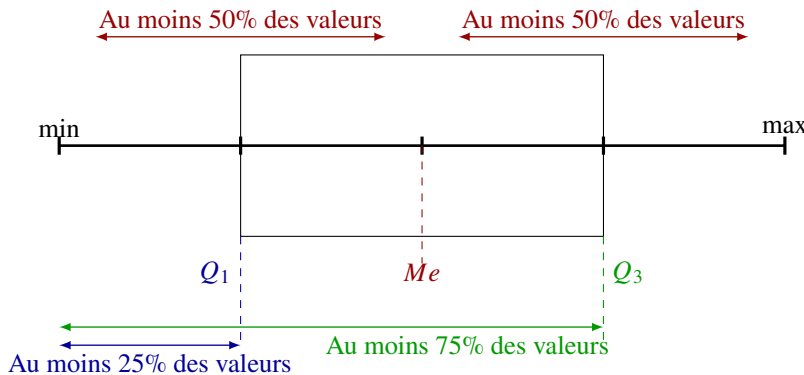
On obtient donc $Q_1 = 40$ et $Q_3 = 46,5$.

Exemple:

- Pour la **série A**, la calculatrice nous donne $Q_1 = 7$, $Med = 8,5$ et $Q_3 = 10$.
- Pour la **série B**, on trouve $Q_1 = 1500$, $Med = 1700$ et $Q_3 = 1900$.

Il est possible de représenter graphiquement les résultats de médiane et de quartiles d'une série statistique par un diagramme en boîte.

Les diagrammes en boîte, ou "boîte à moustaches", permettent de visualiser rapidement des caractéristiques de position (voir figure ci-dessous).



1.3 Caractéristiques de dispersion

1.3.1 Ecart interquartiles

Définition:

On appelle écart interquartile la donnée $EQ = Q_3 - Q_1$.

Exemple:

On reprends l'exemple avec les pointures des 10 personnes

On obtient donc $EQ = Q_3 - Q_1 = 4$.

! Remarque

Plus EQ est petit, plus la médiane est représentative de la série statistique.

1.3.2 Ecart-type d'une série statistique

Définition:

Si une série statistique $x_1 ; \dots ; x_p$ prend n_i fois la valeur x_i , pour $i \in \{1, \dots, p\}$, alors on définit respectivement la variance et l'écart-type par:

$$V = \frac{n_1(x_1 - \bar{x})^2 + \dots + n_p(x_p - \bar{x})^2}{n_1 + \dots + n_p} \text{ et } \sigma = \sqrt{V}$$

! Remarque

L'écart type représente la distance moyenne entre les valeurs de la série et la moyenne.

Au plus il est petit, au plus la moyenne est représentative de la série statistique.

Exercice:

Calculer la variance et l'écart type de la série statistique des pointures des 10 personnes.

2 Série statistique à deux variables

2.1 Ajustement affine

Pour étudier une population selon deux caractères quantitatifs, on considère une série statistique à deux variables x et y composée de n observations $(x_1; y_1), (x_2; y_2), \dots, (x_n; y_n)$, chaque x_i correspondant au caractère x et chaque y_i correspondant au caractère y .

Exemple:

On interroge 8 personnes en leur demandant chacune leur taille et leur poids. On a donc une série statistique x relative à la taille et une série statistique y relative au poids.

Personne	A	B	C	D	E	F	G	H
Taille	165	167	169	171	173	174	175	178
Poids	68	73	71	72	70	75	82	85

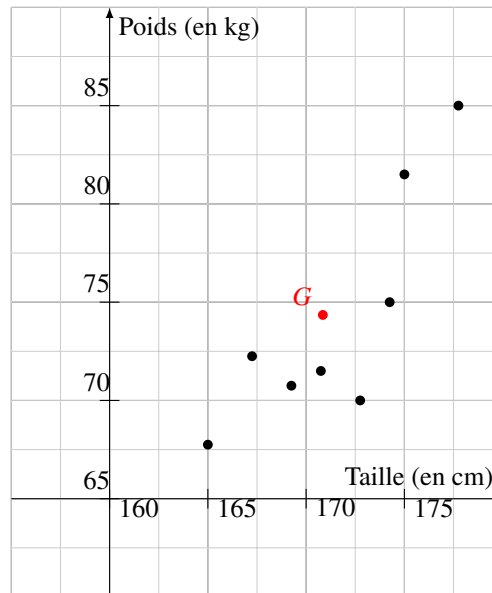
Définition:

A une série statistique on associe :

- Un nuage de points, qui est l'ensemble des n points $M_i(x_i; y_i)$ dans un repère du plan.
- Un point moyen dans le même repère : $G(\bar{x}; \bar{y})$ où \bar{x} et \bar{y} sont respectivement les moyennes des séries x et y .

Exemple:

En reprenant l'exemple précédent on a donc le nuage de point suivant :



En calculant les moyennes des séries statistiques taille (x) et poids (y), on obtient respectivement $\bar{x} = 171,5$ et $\bar{y} = 74,5$. On a donc le point moyen $G(171,5; 74,5)$.

Définition:

Le principe de l'ajustement est de chercher un lien éventuel et simple entre x et y .

Dans le cadre d'un ajustement affine, on cherche à lier x et y par une relation de la forme $y = ax + b$. On obtient alors une droite d'ajustement censée représenter le nuage de points.

Méthode:

Il y a plusieurs méthodes pour tracer la droite d'ajustement :

- "Au jugé", c'est-à-dire de la faire passer "au milieu" du nuage de points.
- En utilisant deux points données, appartenant ou non au nuage de points ou dont l'un des deux est le point G .
- A partir de l'équation de la droite d'ajustement, si elle est donnée.
- Si le coefficient directeur ou l'ordonnée à l'origine est donné et la droite d'ajustement passe par G ou un point du nuage.

! Remarque

Si les points ne sont pas à peu près alignés, un ajustement affine n'a pas d'intérêt. Cependant, on peut faire un changement de variable (c'est-à-dire transformer les valeurs de x ou y) pour obtenir un nuage de points à peu près alignés.

Exemple:

En reprenant l'exemple précédent on a donc diverses méthodes pour obtenir l'ajustement affine de cette série statistique.

- **Méthode 1:** On suppose que la droite passe deux points G et $A(159; 60, 5)$.
On a alors son coefficient directeur $m = \frac{y_G - y_A}{x_G - x_A} = 1,12$.
On a son ordonnée à l'origine donné par $y_G = 1,12x_G + b \iff b = -117,58$.
- **Méthode 2:** On suppose que la droite passe par G en connaissant son ordonnée à l'origine.
On a $y_G = ax_G - 117,58 \iff a \approx 1,12$.
- **Méthode 3:** On suppose que la droite passe par G en connaissant son coefficient directeur.
On a $y_G = 1,12x_G + b \iff b \approx -117,58$.

! Remarque

L'objectif d'un ajustement est de faire une estimation à partir de la relation trouvée entre x et y .

2.2 Méthode des moindres carrés

Il s'agit d'obtenir une droite équidistante des points situés de part et d'autre d'elle-même.

Pour réaliser ceci, on cherche à minimiser la somme des distances des points à la droite au carré.

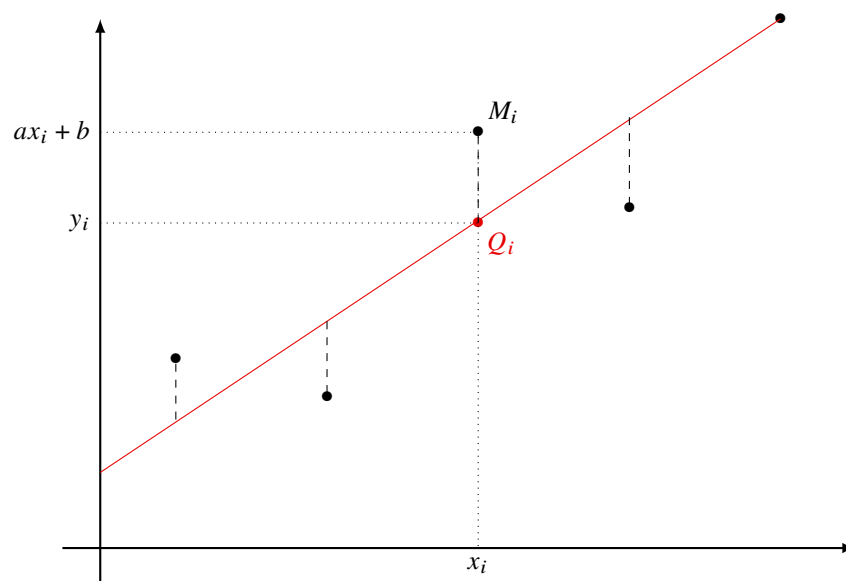
On considère une série statistique à deux variables représentée par un nuage justifiant un ajustement affine.

Définition:

Dans le plan muni d'un repère orthonormal, on considère un nuage de n points de coordonnées $(x_i; y_i)$.

La droite \mathcal{D} d'équation $y = ax + b$ est appelée droite de régression de y en x de la série statistique si et seulement si la quantité suivante est minimale :

$$\sum_{i=1}^n (M_i Q_i)^2 = \sum_{i=1}^n [y_i - (ax_i + b)]^2$$



Définition:

On appelle covariance de la série statistique double de variables x et y le nombre réel

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

! Remarque

- On a : $\text{Cov}(x, x) = \mathbb{V}(x) = \sigma_x^2$.
- Pour les calculs, on pourra aussi utiliser : $\sigma_{xy} = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}$.

Propriété:

La droite de régression \mathcal{D} de y en x a pour équation $y = ax + b$ où

$$\begin{cases} a = \frac{\sigma_{xy}}{[\sigma(x)]^2} \\ b \text{ vérifie } \bar{y} = a\bar{x} + b. \end{cases}$$

Démonstration:

On cherche à minimiser la fonction $F(a, b) = \sum_{i=1}^n [y_i - (ax_i + b)]^2$. On a donc :

$$\begin{aligned} \frac{\partial F}{\partial a}(a, b) &= -2 \sum_{i=1}^n x_i (y_i - ax_i - b) \\ \frac{\partial F}{\partial b}(a, b) &= -2 \sum_{i=1}^n (y_i - ax_i - b) \end{aligned} \tag{1}$$

Afin de minimiser F , on cherche à résoudre le système :

$$\begin{cases} \frac{\partial F}{\partial a}(a, b) = -2 \sum_{i=1}^n x_i (y_i - ax_i - b) = 0 \\ \frac{\partial F}{\partial b}(a, b) = -2 \sum_{i=1}^n (y_i - ax_i - b) = 0 \end{cases}$$

On réécrit la somme en sachant que $\sigma_{xy} = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}$; $\sigma_x^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$ et $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$.

On obtient ainsi le système :

$$\begin{cases} (\sigma_{xy} + \bar{x} \bar{y}) - a(\sigma_x^2 + \bar{x}^2) - b\bar{x} = 0 \\ \bar{y} - a\bar{x} - b = 0 \end{cases}$$

On obtient bien le résultat souhaité après résolution de ce système (par exemple par combinaison en additionnant la première ligne par la deuxième multipliée par \bar{x}).

! Remarque

Les réels a et b sont donnés par la calculatrice.

T.I.

- Touche **STAT**
- Menu **CALC**
- Item **LinReg**
- LinReg L_1, L_2

Casio

- Menu **STAT**
- Item **CALC**
- Règler les paramètres avec **set**
- Item **REG**
- Choisir **X**

Propriété:

Le point moyen G du nuage appartient toujours à la droite de régression de y en x .

Démonstration:

D'après la propriété précédente.

Exercice:

Déterminer une équation de la droite d'ajustement de l'exemple précédente de y en x obtenue par la méthode des moindres carrés.

2.3 Coefficient de corrélation

Définition:

Le coefficient de corrélation linéaire d'une série statistique de variables x et y est le nombre r défini par :

$$r = \frac{\sigma_{xy}}{\sigma_x \times \sigma_y}$$

Ce coefficient sert à mesurer la qualité d'un ajustement affine.

! Remarque

Plus le coefficient de régression linéaire est proche de 1 en valeur absolue, meilleur est l'ajustement linéaire.
Lorsque $r = \pm 1$, la droite de régression passe par tous les points du nuage, qui sont donc alignés.
