



9

STATISTIQUES

STATISTIQUES DESCRIPTIVES UNIVARIÉES ET BIVARIÉES

INTRODUCTION...

La **statistique** et les **probabilités** sont les deux aspects complémentaires de l'étude des phénomènes aléatoires.

- Les probabilités peuvent être considérées comme une branche des mathématiques pures, basée sur des théories abstraites, déconnectée de la réalité. On cherche parfois à les appliquer en collant des modèles probabilistes sur des phénomènes aléatoires concrets.
- La statistique est la science dont l'objet est de recueillir, de traiter et d'analyser des données issues de l'observation de phénomènes aléatoires. La statistique se retrouve dans de nombreux domaines : assurance / finance (fixation des primes d'assurance, gestion de portefeuille, évaluation d'actifs financiers...), biologie / médecine (essais thérapeutiques, dynamique des populations...), sciences de la terre (prévisions météorologiques, exploration pétrolière...), sciences humaines (sondages, enquêtes...), sciences de l'information et de la communication (traitement des images, reconnaissance de formes et de parole...)... On en distingue deux classes :
 - * la **statistique descriptive** (analyse de données) : elle a pour but de résumer l'information contenue dans les données via des indicateurs numériques (moyenne, médiane, quartiles...), des tableaux, des graphiques...
 - * la **statistique inférentielle** : elle a pour but de faire des prévisions et prendre des décisions à partir du traitement des données recueillies. Ces prévisions sont basées sur d'importants résultats de probabilités ; nous en parlerons lors du chapitre 16.

POUR BIEN DÉMARRER...

1. Soit X une variable aléatoire telle que $X(\Omega) = \{x_1, x_2, \dots, x_n\}$.

$$\mathbb{E}(X) = \sum_{k=1}^n x_k \mathbb{P}([X = x_k])$$

$$\mathbb{V}(X) = \mathbb{E}\left(\left(X - \mathbb{E}(X)\right)^2\right) = \sum_{k=1}^n (x_k - \mathbb{E}(X))^2 \mathbb{P}([X = x_k])$$

$$\sigma(X) = \sqrt{\mathbb{V}(X)}$$

2. Soient X et Y deux variables aléatoires telles que $X(\Omega) = \{x_1, x_2, \dots, x_n\}$ et $Y(\Omega) = \{y_1, y_2, \dots, y_m\}$.

$$\text{Cov}(X, Y) = \mathbb{E}\left(\left(X - \mathbb{E}(X)\right)\left(Y - \mathbb{E}(Y)\right)\right) = \sum_{(i,j) \in \llbracket 1;n \rrbracket \times \llbracket 1;m \rrbracket} (x_i - \mathbb{E}(X))(y_j - \mathbb{E}(Y)) \mathbb{P}([X = x_i] \cap [Y = y_j])$$

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma(X)\sigma(Y)}$$

3. Que dire du coefficient de corrélation linéaire d'un couple (X, Y) de variables aléatoires discrètes ?

- $-1 \leq \rho(X, Y) \leq 1$.
- $\rho(X, Y) = 1$ si, et seulement si, l'une des variables aléatoires est presque-sûrement fonction affine strictement croissante de l'autre ;
- $\rho(X, Y) = -1$ si, et seulement si, l'une des variables aléatoires est presque-sûrement fonction affine strictement décroissante de l'autre ;

I VOCABULAIRE STATISTIQUE

On utilise les statistiques pour décrire un ensemble appelé **population** constituée d'**individus** auxquels est associé un **caractère**. On distingue :

- les statistiques **quantitatives** : le caractère est quantitatif, il peut prendre différentes **valeurs**,
- les statistiques **qualitatives** : le caractère est qualitatif, il peut prendre différentes **modalités**.

Lorsque l'on ne peut pas observer toutes les données relatives à la population, il convient alors d'observer seulement un sous-ensemble de cette population, appelé **échantillon**.

Le programme concerne des variables quantitatives... Sur lesquelles nous pourrions calculer des moyennes, médianes, quartiles... Nous ne pouvons pas faire ceci avec des variables qualitatives !

Nous en étudierons deux types :

- les variables quantitatives **discrètes** : les valeurs sont des réels isolés,
- les variables quantitatives **continues** : les valeurs sont des intervalles de nombres.

EXEMPLE 1

On souhaite étudier les étudiants d'ECG.

- Population : l'ensemble des étudiants d'ECG de France.
- Échantillons possibles : les ECG des Chartreux, ou seulement la classe d'ECG4 des Chartreux ; ou les ECG des Maristes...
- Caractères possibles :
 - * caractères quantitatifs : notes au dernier DS de mathématiques, âge, nombre de frères et sœurs, taille, revenus des parents, nombre de cafés consommés chaque jour, temps passé au CDI chaque semaine,...
 - * caractères qualitatifs : couleur de cheveux / d'yeux, sexe, catégorie socio-professionnelle des parents, admission en fin de 2A,...

Les éléments d'analyse (indicateurs, représentations graphiques) diffèrent selon le type de variable étudiée : quantitative discrète, quantitative continue, qualitative. Dans la suite, nous nous limiterons aux statistiques quantitatives discrètes.

II SÉRIES STATISTIQUES QUANTITATIVES DISCRÈTES UNIVARIÉES

Dans cette partie, on s'intéresse à un certain échantillon constitué de n individus d'une population. On étudie un caractère quantitatif discret dont le n -uplet des valeurs observées (la série statistique) est noté $x = (x_1, x_2, \dots, x_n)$. On note z_1, z_2, \dots, z_p (avec $p \in \mathbb{N}^*$) les p valeurs possibles *distinctes* de ce caractère.

DÉFINITIONS 1

EFFECTIF, FRÉQUENCE

Avec les notations précédentes...

D1 Pour tout $i \in \llbracket 1; p \rrbracket$, l'**effectif** de z_i , noté e_i , est le nombre de fois où la valeur z_i est prise dans l'échantillon.

D2 L'**effectif total** est le nombre d'individus de l'échantillon, ici n . On a : $n = \sum_{i=1}^p e_i$.

D3 L'**effectif cumulé croissant** de z_i , noté ECC_i , est le nombre d'individus dont le caractère a une valeur inférieure ou égale à celle-ci : $ECC_i = \sum_{k=1}^i e_k$.

D4 La **fréquence** de z_i , notée f_i , est le réel défini par : $f_i = \frac{e_i}{n}$.

D5 La **fréquence cumulée croissante** de z_i , notée FCC_i , est le réel défini par : $FCC_i = \frac{ECC_i}{n}$.

EXEMPLE 2

On s'intéresse au nombre de frères et sœurs d'un échantillon de 20 étudiants. La série statistique observée est résumée dans les deux premières lignes du tableau ci-dessous, que nous complétons :

valeurs	0	1	2	3	4	5
effectifs	4	5	6	3	1	1
ECC	4	9	15	18	19	20
fréquence	$\frac{1}{5}$	$\frac{1}{4}$	$\frac{3}{10}$	$\frac{3}{20}$	$\frac{1}{20}$	$\frac{1}{20}$
FCC	$\frac{1}{5}$	$\frac{9}{20}$	$\frac{3}{4}$	$\frac{9}{10}$	$\frac{19}{20}$	1

Vocabulaire

On parle aussi de **variable quantitative** ou **variable qualitative**.

Remarque

Les valeurs x_1, x_2, \dots, x_n peuvent être égales !

Pour info...

L'analyse des fréquences (analyse fréquentielle) est utilisée en cryptanalyse pour décrypter des codes !

Remarque

$$FCC_i = \sum_{k=1}^i f_k.$$

Remarque

En laissant toutes les fractions sur 20, les calculs sont encore plus rapides !

II.1 INDICATEURS

Dans l'étude d'une série statistique quantitative univariée, on peut utiliser deux types d'indicateurs :

- des **indicateurs de position**,
- des **indicateurs de dispersion**.

Commençons par des indicateurs de position :

DÉFINITIONS 2

MODE, MOYENNE, MÉDIANE

- D1** Le (ou les) **mode(s)** d'une série statistique est la (ou les) valeur(s) de plus grand effectif.
- D2** En reprenant les notations des définitions précédentes, la **moyenne (empirique)** de la série x , notée \bar{x} , est le réel défini par :
- $$\bar{x} = \frac{1}{n} \sum_{i=1}^p e_i z_i$$
- D3** Une **médiane**, notée \tilde{x} , est un réel tel que la moitié de la série statistique soit inférieure ou égale à \tilde{x} et la moitié soit supérieure ou égale à \tilde{x} .

Important !

- $\bar{x} = \sum_{i=1}^p f_i z_i$
- si $x = (x_1, x_2, \dots, x_n)$ est la série complète (les valeurs peuvent être égales), alors : $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$.

♣ **MÉTHODE 1** ♣ Pour déterminer une médiane d'une série statistique d'effectif total n , on commence par ranger les valeurs de la série par ordre croissant. Puis :

- si l'effectif total n est impair, la seule médiane alors possible est la valeur de rang central dans la série ordonnée ;
- si l'effectif total n est pair, alors, on choisit comme médiane "habituelle" la moyenne des valeurs des deux rangs centraux dans la série ordonnée.

EXEMPLE 3

Reprenons l'exemple précédent.

- Le mode est 2.
- L'effectif de l'échantillon étant pair, toutes les valeurs comprises entre les deux valeurs centrales de la série ordonnée sont des médianes possibles.
Ici, les deux valeurs centrales sont la 10^{ème} et la 11^{ème}, toutes deux égales à 2. Ainsi :

$$\tilde{x} = 2$$

Et maintenant, des indicateurs de dispersion :

DÉFINITIONS 3

VARIANCE ET ÉCART-TYPE

Avec les notations des définitions précédentes...

- D1** La **variance (empirique)** de la série statistique x , notée v_x , est le réel défini par :

$$v_x = \frac{1}{n} \sum_{i=1}^p e_i (z_i - \bar{x})^2$$

- D2** L'**écart-type (empirique)** de la série statistique x , noté s_x , est le réel défini par :

$$s_x = \sqrt{v_x}$$

Important !

- $v_x = \sum_{i=1}^p f_i (z_i - \bar{x})^2$
- si $x = (x_1, x_2, \dots, x_n)$ est la série complète (les valeurs peuvent être égales), alors : $v_x = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$
- v_x est une somme de termes positifs, donc $v_x \geq 0$, ce qui justifie la définition de l'écart-type.
- On note souvent juste s_x^2 la variance et s_x l'écart-type.

PROPRIÉTÉS 1

- P1** La variance d'une série statistique est nulle si, et seulement si, la série statistique est constante.

P2 $s_x^2 = \left(\frac{1}{n} \sum_{i=1}^p e_i z_i^2 \right) - \bar{x}^2$ (formule de Koenig-Huygens)

- P3** Soit x une série statistique de moyenne \bar{x} , de médiane \tilde{x} et de variance s_x^2 . Soient $a, b \in \mathbb{R}$ et y la série statistique obtenue par transformation affine de x de la sorte : $y = ax + b$. On a :

- la moyenne de y est donnée par : $\bar{y} = a\bar{x} + b$
- la médiane de y est donnée par : $\tilde{y} = a\tilde{x} + b$
- la variance de y est donnée par : $s_y^2 = a^2 s_x^2$

Important !

- Si $x = (x_1, x_2, \dots, x_n)$ est la série complète (les valeurs peuvent être égales), alors : $s_x^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$.

★ DÉMONSTRATION :

P1. Avec les notations précédentes :

$$v_x = 0 \iff \frac{1}{n} \sum_{i=1}^p e_i(z_i - \bar{x})^2 = 0$$

↙ une somme de termes positifs est nulle si, et seulement si, tous ses termes sont nuls

$$\iff \forall i \in \llbracket 1; p \rrbracket, e_i(z_i - \bar{x})^2 = 0$$

↙ par définition, pour tout $i \in \llbracket 1; p \rrbracket, e_i \neq 0$

$$\iff \forall i \in \llbracket 1; p \rrbracket, z_i = \bar{x}$$

D'où le résultat.

P2. Immédiat...

P3. Assez immédiat...

★

DÉFINITIONS 4

ÉTENDUE, QUANTILES

- D1** L'**étendue** d'une série statistique est la différence entre sa plus grande et sa plus petite valeur.
- D2** Soit $p \in]0; 1[$. Le **quantile d'ordre p** est la plus petite valeur de la série statistique pour laquelle la FCC atteint ou dépasse p . En particulier :
- les **quartiles** sont les quantiles d'ordre 0, 25 ; 0, 5 et 0, 75 :
 - * le premier quartile, noté Q_1 , est la plus petite valeur de la série telle qu'au moins 25% des valeurs prises lui soient inférieures ou égales ;
 - * le troisième quartile, noté Q_3 , est la plus petite valeur de la série telle qu'au moins 75% des valeurs prises lui soient inférieures ou égales.
 - les **déciles** sont les quantiles d'ordre 0, 1 ; 0, 2 ; ... ; 0, 9.
- D3** L'**intervalle interquartile** est l'intervalle $[Q_1; Q_3]$. L'**écart interquartile** est $Q_3 - Q_1$: c'est l'amplitude de l'intervalle interquartile.

Remarque

Quand on souhaite comparer deux séries statistiques, on peut utiliser :

- le couple (moyenne, écart-type),
- le triplet (médiane, écart-interquartile, étendue).

Important !

Par définition, l'intervalle interquartile contient au moins 50% des valeurs de la série.

EXEMPLE 4

Toujours avec les données de l'Exemple 2 :

- L'amplitude est égale à 5.
- $Q_1 = 1$, car 1 est la plus petite valeur de la série pour laquelle la FCC atteint ou dépasse 0, 25.
- $Q_3 = 2$, car 2 est la plus petite valeur de la série pour laquelle la FCC atteint ou dépasse 0, 75.
- L'intervalle interquartile est $[1; 2]$, il contient ici 55% des valeurs de la série.
- L'écart interquartile est égal à 1.
- Le premier décile vaut 0, car 0 est la plus petite valeur de la série pour laquelle la FCC atteint ou dépasse 0, 1.

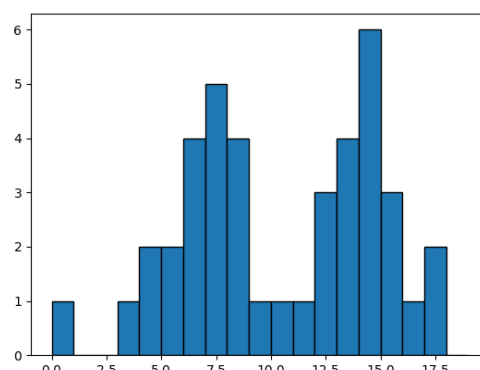
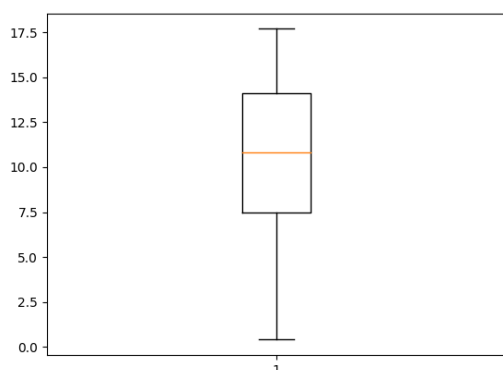
II.2 REPRÉSENTATIONS GRAPHIQUES

Les trois représentations graphiques les plus fréquentes pour des séries statistiques quantitatives discrètes univariées :

- le diagramme en barres (d'effectifs ou de fréquences) ou en bâtons,
- l'histogramme,
- le diagramme de Tukey, ou boîte à moustaches.

EXEMPLE 5

A gauche, le diagramme de Tukey des notes de l'épreuve 1 du CB, à droite l'histogramme d'effectifs avec des intervalles de la forme $[k; k + 1[$, où $k \in \llbracket 0; 19 \rrbracket$.



Remarque

Ne pas confondre diagramme en barres et histogramme ! L'histogramme est utilisé lorsque les valeurs obtenues sont regroupées en intervalles...

Un peu d'histoire

John Tukey (1915-2000, américain) est un statisticien à qui on doit des méthodes d'analyse de données ainsi que l'algorithme de transformation de Fourier rapide (FFT en anglais), sur lequel il a travaillé avec James Cooley.

Remarque

Le diagramme en barres n'aurait pas beaucoup d'intérêt ici, car trop peu de notes sont identiques ! Il aurait en revanche un intérêt pour les données de l'exemple 2.

III SÉRIES STATISTIQUES QUANTITATIVES DISCRÈTES BIVARIÉES

Comme nous avons pu le constater dans la partie précédente, le principe de l'étude d'une série statistique quantitative discrète univariée est de voir comment les effectifs de chaque valeur peuvent influencer la série et ensuite d'analyser cette série avec des indicateurs tenant compte de ces effectifs.

Dans ce qui va suivre, nous ne nous intéresserons plus aux effectifs ; mais au comportement d'une série de données par rapport à une autre...

On considère à présent une population sur laquelle nous mesurons deux caractères quantitatifs discrets notés X et Y . Sur un certain échantillon de taille n de cette population, nous obtenons ainsi deux séries statistiques quantitatives univariées $x = (x_1, x_2, \dots, x_n)$ et $y = (y_1, y_2, \dots, y_n)$ **chaque couple (x_i, y_i) étant associé à un même individu de l'échantillon.**

Plutôt que d'étudier chaque caractère (et donc chaque série statistique) séparément, intéressons-nous plutôt à la façon dont X et Y se comportent l'une par rapport à l'autre : existe-t-il une relation entre X et Y ?

Remarque

En fait, on peut voir X et Y comme des variables aléatoires et x et y comme des n -uplets de réalisations de ces variables aléatoires...

III.1 NUAGE DE POINTS ET POINT MOYEN

DÉFINITIONS 5

Soient $x = (x_1, x_2, \dots, x_n)$ et $y = (y_1, y_2, \dots, y_n)$ deux séries statistiques.

D1 Le **nuage de points** du couple (x, y) est l'ensemble des points du plan de coordonnées (x_i, y_i) pour $i \in \llbracket 1; n \rrbracket$.

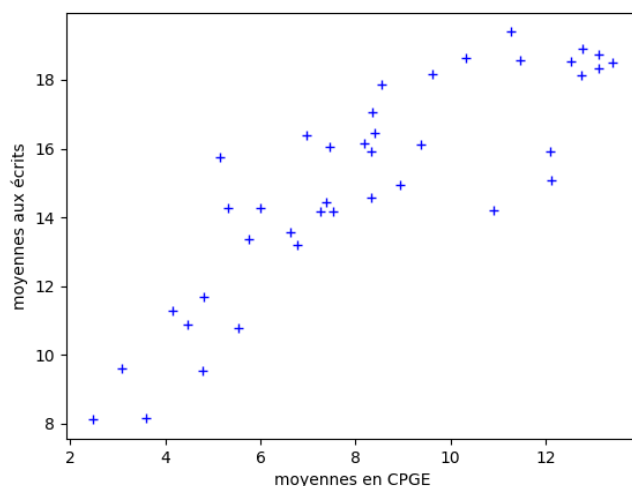
D2 Le **point moyen** du nuage est le point de coordonnées (\bar{x}, \bar{y}) .

EXEMPLE 6

On considère la promo 2025 de maths appli dont on a récupéré la moyenne en mathématiques en deuxième année de CPGE ECG ainsi que la moyenne aux écrits de mathématiques.

Les commandes **Python** ci-dessous ont permis d'obtenir le nuage de points qui suit.

```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3
4 notes=pd.read_excel('notes_promo_2025.xlsx')
5 x=notes["AnneeTotal"]
6 y=notes["EcritsTotal"]
7 plt.plot(x,y,"b+")
8 plt.xlabel("moyennes en CPGE")
9 plt.ylabel("moyennes aux écrits")
10 plt.show()
```



Quelles remarques pouvons-nous faire ?

- La corrélation entre la moyenne annuelle et la moyenne aux écrits semble assez forte : plus la moyenne annuelle est élevée, plus la moyenne aux écrits l'est également.
- Le nuage de points est *relativement proche d'une droite*.

III.2 COVARIANCE ET COEFFICIENT DE CORRÉLATION LINÉAIRE

DÉFINITIONS 6

COVARIANCE ET COEFFICIENT DE CORRÉLATION LINÉAIRE

Soient $x = (x_1, x_2, \dots, x_n)$ et $y = (y_1, y_2, \dots, y_n)$ deux séries statistiques de moyennes et écart-types respectifs \bar{x}, \bar{y} et s_x, s_y .

D1 La **covariance** du couple (x, y) est le réel, noté $\text{Cov}(x, y)$, défini par :

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

D2 Si σ_x et σ_y sont non nuls, le **coefficient de corrélation linéaire** du couple (x, y) est le réel, noté $\rho(x, y)$, défini par :

$$\rho(x, y) = \frac{\text{Cov}(x, y)}{s_x s_y}$$

Comme dans le cas des couples de variables aléatoires discrètes, on retrouve :

PROPRIÉTÉS 2

Avec les notations précédentes...

P1 $\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}$ (formule de Koenig-Huygens)

P2 $\text{Cov}(x, y)^2 \leq s_x^2 s_y^2$ (inégalité de Cauchy-Schwarz)

P3 $-1 \leq \rho(x, y) \leq 1$

P4 $|\rho(x, y)| = 1 \iff \exists (a, b) \in \mathbb{R}^* \times \mathbb{R} / y = ax + b$
Et :

- si $\rho(x, y) = 1$, alors $a > 0$
- si $\rho(x, y) = -1$, alors $a < 0$

* DÉMONSTRATION : Analogues à celles dans le cas des couples de variables aléatoires discrètes... mais plus simples, car toutes les sommes en jeu sont des sommes finies !

*

III.3 RÉGRESSION LINÉAIRE

La propriété précédente (P4), permet de savoir quand y est fonction affine de x ; mais on ne sait en revanche pas encore comment trouver les coefficients a, b de cette fonction affine. Voyons un peu cela...

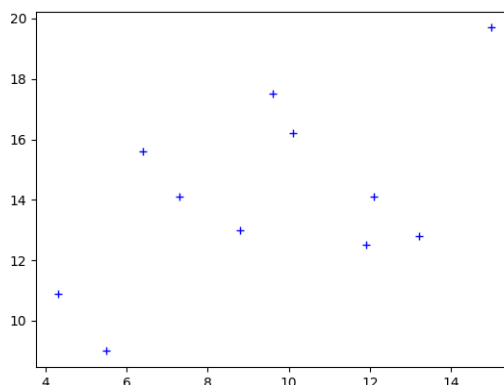
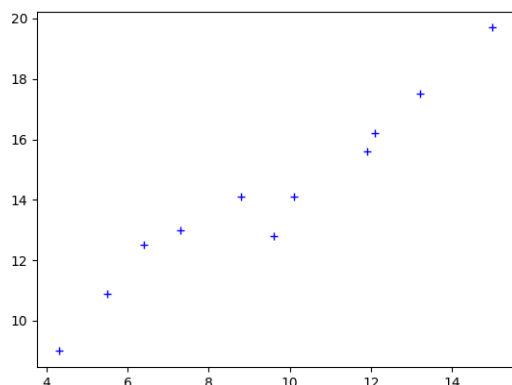
Chercher un modèle de régression consiste à chercher s'il existe une fonction f et une variable aléatoire ε telles que :

$$Y = f(X) + \varepsilon$$

où f est alors appelée **fonction de régression** et ε **erreur d'ajustement** (ou résidu). Dans ce qui suit, nous nous limiterons au cas de la régression linéaire : cas où f est affine.

✗ Attention !
Régression **linéaire** et pourtant fonction **affine** !

Certains nuages de points ont des formes pouvant penser qu'il existe un ajustement affine entre X et Y ...



La forme du nuage de gauche nous laisse penser qu'un ajustement affine est possible entre X et Y ; celle du nuage de droite beaucoup moins.

On cherche à identifier une droite qui ajusterait assez bien le nuage de points, selon des critères à définir. Si l'on suppose qu'il existe des réels a, b tels que pour les séries x et y , on a $y = ax + b$, alors pour tout $i \in \llbracket 1; n \rrbracket$ l'erreur commise en utilisant la valeur $ax_i + b$ pour estimer y_i est égale à $y_i - (ax_i + b)$.

Pour déterminer les coefficients a et b , on va utiliser ici le **principe des moindres carrés** qui consiste à chercher les valeurs de a et b qui minimisent la somme des carrés de ces erreurs : $\sum_{i=1}^n (y_i - ax_i - b)^2$.

THÉORÈME 1

Soient $x = (x_1, x_2, \dots, x_n)$ et $y = (y_1, y_2, \dots, y_n)$ deux séries statistiques de moyennes et écart-types respectifs \bar{x}, \bar{y} et s_x, s_y , avec s_x et s_y non nuls, et de covariance $\text{Cov}(x, y)$.

L'unique couple (a, b) rendant minimale la quantité $\sum_{i=1}^n (y_i - ax_i - b)^2$ est défini par :

$$a = \frac{\text{Cov}(x, y)}{s_x^2} ; \quad b = \bar{y} - a\bar{x}$$

Vocabulaire

- Pour ce tel couple, la droite d'équation $y = ax + b$ est appelée **droite des moindres carrés**.
- X est la **variable explicative** et Y la **variable à expliquer**.

À retenir...

Puisque $\rho(x, y)$ et $\text{Cov}(x, y)$ ont même signe, $\rho(x, y)$ et le coefficient directeur de la droite des moindres carrés ont également le même signe.

★ DÉMONSTRATION : Soit $(a, b) \in \mathbb{R}^2$. Considérons la fonction $f : (a, b) \mapsto \sum_{i=1}^n (y_i - ax_i - b)^2$.

1. **Cas particulier.** Supposons que $\bar{x} = \bar{y} = 0$.

1.a. Montrons alors :

$$\forall (a, b) \in \mathbb{R}^2, f(a, b) \geq a^2 \sum_{i=1}^n x_i^2 - 2a \sum_{i=1}^n x_i y_i + \sum_{i=1}^n y_i^2$$

avec égalité si, et seulement si $b = 0$.

Soit $(a, b) \in \mathbb{R}^2$. On a :

$$\begin{aligned} f(a, b) &= \sum_{i=1}^n (y_i - ax_i - b)^2 \\ &= \sum_{i=1}^n (y_i^2 + a^2 x_i^2 + b^2 - 2ax_i y_i - 2by_i + 2abx_i) \\ &= \sum_{i=1}^n y_i^2 + a^2 \sum_{i=1}^n x_i^2 + nb^2 - 2a \sum_{i=1}^n x_i y_i - 2b \sum_{i=1}^n y_i + 2ab \sum_{i=1}^n x_i \\ &= \sum_{i=1}^n y_i^2 + a^2 \sum_{i=1}^n x_i^2 + nb^2 - 2a \sum_{i=1}^n x_i y_i \end{aligned} \quad \leftarrow \bar{x} = 0 \text{ et } \bar{y} = 0$$

Puisque $nb^2 \geq 0$, on a bien :

$$f(a, b) \geq a^2 \sum_{i=1}^n x_i^2 - 2a \sum_{i=1}^n x_i y_i + \sum_{i=1}^n y_i^2$$

avec égalité si, et seulement si $nb^2 = 0$, autrement dit, si et seulement si, $b = 0$.

1.b. Dédudons-en que $f(a, b)$ est minimale lorsque $a = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$ et $b = 0$.

D'après ce qui précède, $f(a, b)$ est minimale si, et seulement si, $b = 0$ et que la quantité $a^2 \sum_{i=1}^n x_i^2 - 2a \sum_{i=1}^n x_i y_i + \sum_{i=1}^n y_i^2$ est minimale.

Or $a^2 \sum_{i=1}^n x_i^2 - 2a \sum_{i=1}^n x_i y_i + \sum_{i=1}^n y_i^2$ est une expression polynomiale de degré 2 (en effet, les x_i ne sont pas tous nuls, sinon, on aurait $s_x = 0$...) en a à coefficient dominant positif ; elle est donc minimale si, et seulement si,

$$a = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}.$$

Conclusion : $f(a, b)$ est minimale lorsque $a = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$ et $b = 0$.

Rappels...

- Si $\alpha > 0$, alors la fonction $x \mapsto ax^2 + \beta x + \gamma$ possède un minimum, atteint en $-\frac{\beta}{2\alpha}$.
- Si $\alpha < 0$, alors la fonction $x \mapsto ax^2 + \beta x + \gamma$ possède un maximum, atteint en $-\frac{\beta}{2\alpha}$.

2. **Cas général.** Considérons deux nouvelles séries statistiques x' et y' définies par :

$$x' = x - \bar{x} \quad ; \quad y' = y - \bar{y}$$

2.a. Que dire de $\bar{x'}$ et $\bar{y'}$?

On a alors :

$$\bar{x'} = 0 \quad ; \quad \bar{y'} = 0$$

2.b. Exprimons $f(a, b)$ en fonction des x'_i et y'_i , en posant $b' = b - \bar{y} + a\bar{x}$.

Pour tout $(a, b) \in \mathbb{R}^2$:

$$\begin{aligned} f(a, b) &= \sum_{i=1}^n (y_i - ax_i - b)^2 && \left. \begin{array}{l} \nearrow \\ \searrow \end{array} \right\} x' = x - \bar{x} \text{ et } y' = y - \bar{y} \\ &= \sum_{i=1}^n (y'_i + \bar{y} - a(x'_i + \bar{x}) - b)^2 \\ &= \sum_{i=1}^n (y'_i - ax'_i - (b - \bar{y} + a\bar{x}))^2 \\ &= \sum_{i=1}^n (y'_i - ax'_i - b')^2 \end{aligned}$$

2.c. Concluons !

D'après le cas particulier (licite car $\bar{x'} = \bar{y'} = 0$), la quantité $\sum_{i=1}^n (y'_i - ax'_i - b')^2$ est minimale lorsque $b' = 0$

$$\text{et } a = \frac{\sum_{i=1}^n x'_i y'_i}{\sum_{i=1}^n x'^2_i}.$$

Autrement dit, d'après la question précédente, $f(a, b)$ est minimale lorsque $b = \bar{y} - a\bar{x}$ et $a = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$.

Mais :

$$\begin{aligned} \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} &= \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \frac{\text{Cov}(x, y)}{s_x^2} \end{aligned}$$

Conclusion : la quantité $\sum_{i=1}^n (y_i - ax_i - b)^2$ minimale lorsque

$$a = \frac{\text{Cov}(x, y)}{s_x^2} \quad ; \quad b = \bar{y} - a\bar{x}$$

Remarque

Une autre démonstration pourra être vue dans le chapitre 13 sur les fonctions de deux variables.

PROPRIÉTÉ 3

La droite de régression linéaire obtenue par la méthode des moindres carrés passe par le point moyen du nuage de points.

★ DÉMONSTRATION : Notons a le coefficient directeur de la droite des moindres carrés et b son ordonnée à l'origine. D'après le théorème précédent :

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

Autrement dit :

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

Le point moyen, de coordonnées (\bar{x}, \bar{y}) , appartient donc à la droite des moindres carrés, puisque ses coordonnées vérifient l'équation réduite de cette droite.

PROPRIÉTÉ 4

INTERPRÉTATION DE LA VALEUR DE $\rho(x, y)$

Plus $|\rho(x, y)|$ est proche de 1, meilleur est l'ajustement affine par la droite des moindres carrés.

Remarque

Cette propriété va dans le sens de Propriétés 3 - P4 !

* DÉMONSTRATION :

Notons a le coefficient directeur de la droite des moindres carrés et b son ordonnée à l'origine. Montrons que, plus $|\rho(x, y)|$ est proche de 1, plus l'erreur commise en approchant y par $ax + b$ est faible. On a, d'après le théorème 1 :

$$\begin{aligned} \sum_{i=1}^n (y_i - ax_i - b)^2 &= \sum_{i=1}^n (y_i - ax_i - \bar{y} + a\bar{x})^2 \\ &= \sum_{i=1}^n (y_i - \bar{y} - a(x_i - \bar{x}))^2 \\ &= \sum_{i=1}^n (y_i - \bar{y})^2 - 2a \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) + a^2 \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= ns_y^2 - 2an \text{Cov}(x, y) + a^2 ns_x^2 \\ &= n \left(s_y^2 - 2 \frac{\text{Cov}(x, y)}{s_x^2} + \frac{\text{Cov}(x, y)^2}{s_x^2} \right) \quad \leftarrow \text{théorème 1} \\ &= n \left(s_y^2 - \frac{\text{Cov}(x, y)^2}{s_x^2} \right) \\ &= n \left(s_y^2 - \frac{\rho(x, y)^2 s_x^2 s_y^2}{s_x^2} \right) \quad \leftarrow \text{définition de } \rho(x, y) \\ &= ns_y^2 (1 - \rho(x, y)^2) \end{aligned}$$

Par conséquent, plus $\rho(x, y)^2$ est proche de 1, plus $\sum_{i=1}^n (y_i - ax_i - b)^2$ est proche de 0 et donc meilleur est l'ajustement affine par la droite des moindres carrés.

Important !

La quantité $\sum_{i=1}^n (y_i - ax_i - b)^2$ est la somme des carrés des écarts entre la valeur exacte et la valeur approchée obtenue par la méthode des moindres carrés. Plus cette quantité est faible, meilleure est l'approximation.

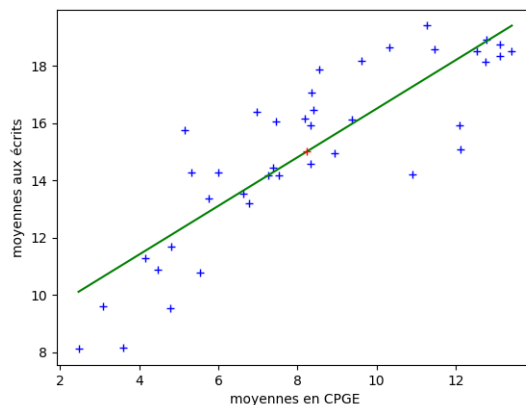
Remarque

De la dernière relation, on retrouve également :

- $-1 \leq \rho(x, y) \leq 1$
- $|\rho(x, y)| = 1$ ssi $\sum_{i=1}^n (y_i - ax_i - b)^2 = 0$ ssi (une somme de termes positifs est nulle ssi tous ses termes sont nuls) : $\forall i \in \llbracket 1; n \rrbracket, y_i = ax_i + b$: tous les points du nuage sont sur la droite des moindres carrés : l'ajustement affine est parfait !

EXEMPLE 7

Dans l'exemple étudié, on obtient graphiquement :



La droite des moindres carrés est la droite d'équation $y = ax + b$, avec $a \simeq 0,84$ et $b \simeq 8$. On trouve : $\rho(x, y) \simeq 0,85$. L'ajustement affine est bon.

Pour info...

Parfois, on n'a pas directement $y = ax + b$, mais on peut par exemple avoir :

- $y = ax^2 + b$ et l'ajustement affine sera entre x^2 et y
- $y = Ce^{ax}$ et l'ajustement affine sera entre $\ln(y)$ et x ...

III.4 MISE EN GARDE !

Il ne faut surtout pas confondre corrélation et causalité ! Deux variables augmentant en fonction d'un même troisième facteur seront assez fortement corrélées, et pourtant il n'y a pas forcément de lien de causalité entre ces deux variables. Il est même possible que la corrélation soit le fruit du hasard !

Pour rire un peu...

Ici : <https://www.tylervigen.com/spurious-correlations>

Lors d'une analyse statistique révélant un fort coefficient de corrélation linéaire, on peut se poser trois questions :

- une variable a-t-elle un lien de cause à effet avec l'autre ?
- existe-t-il une troisième variable ayant un lien de causalité avec les deux précédentes ?
- la corrélation est-elle fortuite ?