

III Statistique descriptive. (En Analyse des données)

01 Échantillons et caractères

1 Échantillon (ou population).

C'est un ensemble de personnes ou d'objets sur lequel on fera une étude. Les elts de l'ensemble s'appellent individus.

Remarque

En Analyse des données, le mot "échantillon" n'a pas le même sens en statistique, où les individus sont supposés provenir d'une population plus vaste.

Exemple :

Personnel d'entreprise.

Ensemble des étudiants en 1^{ère} année de PG à l'Institut de Mathématique.

2. Les différents types de caractères

Définition

Un caractère des individus d'une population est une application qui, à chaque individu de la population, fait correspondre un elt d'un ensemble appelé ensemble des modalités de caractères.

Exemple

Population : ensemble des étudiants en 1^{ère} année de PG à USTHB.

Caractère : couleur des cheveux.

ensemble des modalités : Noire, Blanc, Brun, Roux.

a) Caractères qualitatifs :

Un caractère qualitatif, est un caractère dont les modalités échappent à la mesure, elles peuvent seulement être constatées.

Exemple :

Couleur des cheveux,
sexe d'une personne

b) Caractères quantitatifs.

On dit qu'un caractère est quantitatif, lorsqu'il est mesurable ou repérable. A chaque individu correspond alors un nombre qui est la mesure ou la valeur du caractère. Ace nbre on donne le nom de variable.

Exemples -

Taille, poids ou âge d'une personne.
Durée de vie d'une lampe

I Echantillon d'une variable.

I = ensemble de n individus

x = variable quantitative.

$$x : I \longrightarrow \mathbb{R}$$

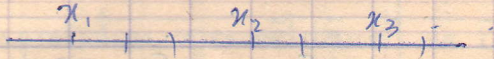
$$i \in I \longrightarrow x(i) = x_i$$

On appelle échantillon de la variable x , l'ensemble $e = \{x_i, i \in I\}$.

On suppose qu'à tout individu i est associé un poids $p_i > 0$
tg $\sum_{i=1}^n p_i = 1$.

1 | Deux représentations de \mathcal{E} .
On peut regarder \mathcal{E} soit

- comme un sous ensemble de \mathbb{R} .



- comme un vecteur de \mathbb{R}^n .

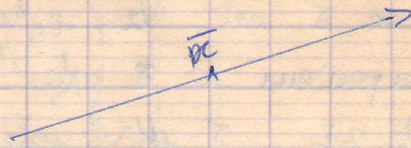
$\underline{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$ = vecteur de \mathbb{R}^n associé à l'échantillon \mathcal{E}

2) (Trois) caractéristiques de valeurs centrale

a) La moyenne

$$\bar{x} = \sum_{i=1}^n p_i \cdot x_i$$

(Propre)



Son image sur la droite \mathbb{R} est le centre de gravité des pts x_i munis de poids p_i

Propriétés

$$\min_{1 \leq i \leq n} x_i \leq \bar{x} \leq \max_{1 \leq i \leq n} x_i$$

b) La médiane

L'idée générale de la définition d'une médiane est de partager le domaine de la variable en deux parties. De manière que chaque (chaque) des parties, contiennent, si possible, la moitié de l'effectif total, ou, si ce n'est pas possible de se rapprocher au mieux de (celle) est objectif

Exemple.

Une population N a subi un examen noté de 0 à 5
La distribution des notes est donnée par le tableau suivant

valeur des notes	0	1	2	3	4	5
effectifs	5	15	30	25	20	5
Frequence	0,05	0,15	0,30	0,25	0,20	0,05

Soit δ un nbre compris entre 2 et 3 ($2 < \delta < 3$)
 la proportion de la population A qui a obtenu une note inférieure à δ est $1/2$, la proportion qui a obtenu une note supérieure à δ est $1/2$.

$\forall \delta \in]2, 3[$, δ est une médiane

En duo que l'intervalle $]2, 3[$ est un intervalle médiane

c) Moyenne des valeurs extrêmes.

$$\frac{1}{2} \left(\min_{1 \leq i \leq n} x_i + \max_{1 \leq i \leq n} x_i \right) = ME(x)$$

3/ Trois caractéristiques de dispersion

a) variance et écart-type.

$$\begin{aligned} \text{Var}(x) &= \sum_{i=1}^n p_i (x_i - \bar{x})^2 \\ &= \sum_{i=1}^n p_i x_i^2 - \bar{x}^2 \end{aligned}$$

$$\sigma_x = \sqrt{\text{Var}(x)} = \sqrt{\sum_{i=1}^n p_i (x_i - \bar{x})^2}$$

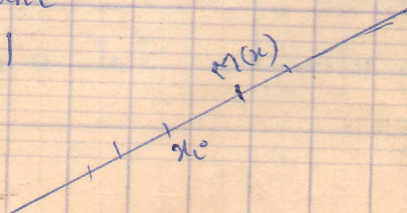
Remarques.

- si $\text{var}(x) = 0$ E est formé de n valeurs égales
- L'écart-type est "homogène à x " (il s'exprime avec la même unité que x)

b) écart moyen

$M(x) =$ médiane.

$$ECM(x) = \sum_{i=1}^n p_i |x_i - M(x)|$$



c) Etendue

$$W(x) = \max_{1 \leq i \leq n} x_i - \min_{1 \leq i \leq n} x_i$$

4) Interprétation géométrique de la moyenne et de la variance de \mathbb{R}^n
 Exercice 4 de la série

$$D = \begin{pmatrix} p_1 & 0 \\ 0 & p_n \end{pmatrix} \quad \sum_{i=1}^n p_i = 1, \quad p_i \geq 0$$

$$2) d(x, y) = N(x - y) = \sqrt{(x - y)^T D (x - y)}$$

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$D(x - y) = \begin{pmatrix} p_1(x_1 - y_1) \\ \vdots \\ p_n(x_n - y_n) \end{pmatrix}$$

$$\begin{aligned} (x - y)^T D (x - y) &= (x_1 - y_1, \dots, x_n - y_n) \begin{pmatrix} p_1(x_1 - y_1) \\ \vdots \\ p_n(x_n - y_n) \end{pmatrix} \\ &= \sum_{i=1}^n p_i (x_i - y_i)^2 \end{aligned}$$

$$\Rightarrow d(x, y) = \sqrt{\sum_{i=1}^n p_i (x_i - y_i)^2}$$

$$3) u = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \quad \Delta u = \langle u \rangle$$

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

any = αu projection D-orth de x sur Δu

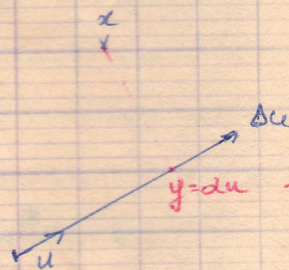
$$x - y \perp_D u \quad (\Leftrightarrow) \quad \langle x - y, u \rangle_D = 0$$

$$(\Leftrightarrow) \quad \langle x, u \rangle_D = \alpha \langle u, u \rangle_D$$

$$\langle u, u \rangle_D = {}^t u D u = (1, \dots, 1) \begin{pmatrix} p_1 & 0 \\ 0 & p_n \end{pmatrix} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = 1$$

$$\alpha = \langle x, u \rangle_D = {}^t x D u = \sum_{i=1}^n p_i x_i = \bar{x}$$

Conclusion: l'abscisse de la projection orthogonale de x sur Δu



est égale à \bar{x} , $y = \bar{x}u$ // null

$$b. d(x, \Delta u) = \inf_{y \in \Delta u} d(x, y) = d(x, y)$$

$$d(x, y) = d(x, \alpha u) = \sqrt{t(x - \alpha u) D(x - \alpha u)} \quad \alpha = \bar{x}$$

$$\begin{pmatrix} p_1 & & & \\ & 0 & & \\ & & \ddots & \\ & & & p_n \end{pmatrix} \begin{pmatrix} x_1 - \bar{x} \\ \vdots \\ x_n - \bar{x} \end{pmatrix} = \begin{pmatrix} p_1(x_1 - \bar{x}), \dots, p_n(x_n - \bar{x}) \end{pmatrix}$$

$$t(x - \bar{x}u) D(x - \bar{x}u) = \sum_{i=1}^n p_i (x_i - \bar{x})^2$$

$$d(x, \Delta u) = \sqrt{\sum_{i=1}^n p_i (x_i - \bar{x})^2} = \sigma_x$$

II. Echantillon d'un couple de variables.

I = ensemble de n individus.

x et y = deux variables quantitatives

$$I \xrightarrow{x} \mathbb{R}$$

$$i^\circ \quad v \longrightarrow x(i^\circ) = x_i$$

$$I \xrightarrow{y} \mathbb{R}$$

$$i^\circ \quad v \longrightarrow y(i^\circ) = y_i$$

On appelle échantillon du couple (x, y) , l'ensemble $E = \{(x_i, y_i) \mid 1 \leq i \leq n\}$.

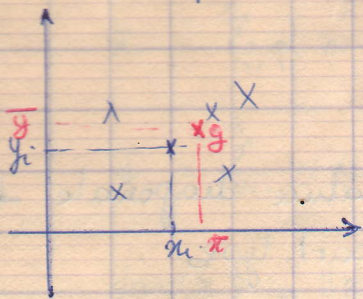
On suppose que chaque individu est muni d'un poids $p_i > 0$ ($\sum_{i=1}^n p_i = 1$)

A l'échantillon E correspondent les échantillons

$$E_x = \{x_i \mid 1 \leq i \leq n\} \quad E_y = \{y_i \mid 1 \leq i \leq n\}$$

On peut décrire E_x et E_y séparément, mais c'est insuffisant car il est bon de mettre en évidence les liens éventuels entre les variables x et y .

1. Représentation de l'échantillon dans \mathbb{R}^2 .



On représente E dans \mathbb{R}^2 par le nuage de n pts (x_i, y_i)
 $g = (\bar{x}, \bar{y})$ est le centre de gravité du nuage
 $\bar{x} = \frac{1}{n} \sum_{i=1}^n p_i x_i$
 $\bar{y} = \frac{1}{n} \sum_{i=1}^n p_i y_i$

2. Covariance et coefficient de corrélation linéaire.

• Covariance de variable x et y .

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

• Coefficient de corrélation linéaire des variables x et y

$$\text{correl}(x, y) = r(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y}$$

• Matrice de variance-covariance associée à l'échantillon E .

$$V = \begin{pmatrix} \text{Var } x & \text{Cov}(x, y) \\ \text{Cov}(x, y) & \text{Var } y \end{pmatrix}$$

• Matrice de corrélation associée à l'échantillon E

$$R = \begin{pmatrix} 1 & r(x, y) \\ r(x, y) & 1 \end{pmatrix}$$

Remarques

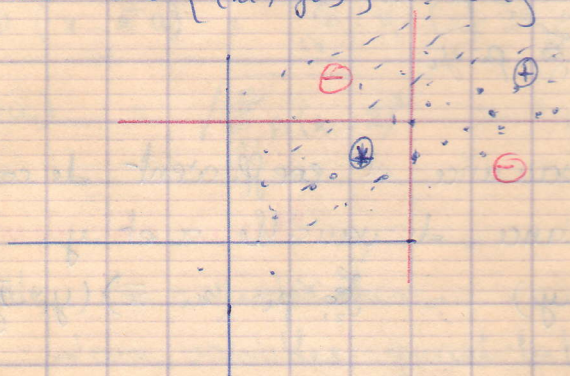
* La covariance est sensible aux changements d'échelle sur x et y , le coefficient de corrélation ne l'est pas

$$r(ax, by) = r(x, y) \quad \forall a, b \in \mathbb{R}$$

$$* R = D_{1/\sigma} \cdot V \cdot D_{1/\sigma}$$

où $D_{1/\sigma} = \begin{pmatrix} 1/\sigma_x & 0 \\ 0 & 1/\sigma_y \end{pmatrix}$ matrice diagonale des inverse, des écart-type.

+ signe de la covariance et du coefficient de corrélation
 Reprenons le nuage des n pts (x_i, y_i) qui visualise de \mathbb{R}^2
 l'échantillon $E = \{(x_i, y_i); 1 \leq i \leq n\}$



si le nuage a une allure \nearrow , $\text{cov}(x, y) > 0$ $r(x, y) > 0$

si le nuage a une allure \searrow , $\text{cov}(x, y) < 0$

3. Représentation de l'échantillon dans \mathbb{R}^n .

Signification du coefficient de corrélation.

ds \mathbb{R}^n , l'échantillon $E = \{(x_i, y_i); 1 \leq i \leq n\}$ a pour image le couple de pts (x, y) où

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \quad \text{et} \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

\mathbb{R}^n est muni de la métrique D_p .

$$D_p = \begin{pmatrix} p_1 & & 0 \\ & \ddots & \\ 0 & & p_n \end{pmatrix}$$

Nous avons vu que

$$\bar{x} = \sum_{i=1}^n p_i x_i = {}^t x D_p u$$

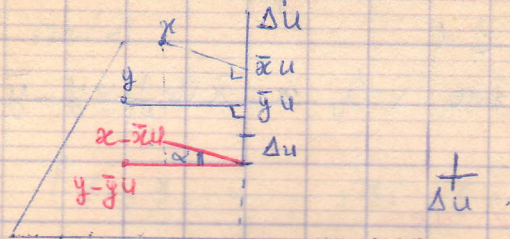
$$u = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, u \in \mathbb{R}^n.$$

$$\bar{y} = \sum_{i=1}^n p_i y_i = {}^t y D_p u$$

$$\text{Var}(x) = \|x - \bar{x}u\|_{D_p}^2$$

$$\text{Var}(y) = \|y - \bar{y}u\|_{D_p}^2$$

Si on considère la décomposition en somme directe $\mathbb{R}^n = \Delta u \oplus \Delta u^\perp$



où Δu désigne le sous-espace vectoriel D_p -orthogonal à Δu . En cas, $x - \bar{x}u \in \Delta u^\perp$ et $y - \bar{y}u \in \Delta u^\perp$

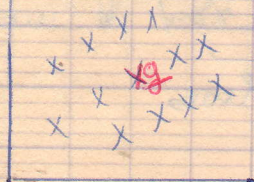
Comme tout vecteur de \mathbb{R}^n , $x - \bar{x}u$ et $y - \bar{y}u$ symbolisent les caractères, on dit que ce sont les caractères et y centrés, en effet leur moyenne

$$\sum_{i=1}^n p_i (x_i - \bar{x}) = 0 \quad \sum_{i=1}^n p_i (y_i - \bar{y}) = 0$$

Les caractères centrés, $x - \bar{x}u$ et $y - \bar{y}u$ sont les projections D_p -orthogonales sur Δu^\perp des vecteurs x et y .

Remarque.

ds \mathbb{R}^n , cette revient à mettre l'origine au centre de gravité \bar{g} du nuage de n p^t $\bar{g} = \begin{pmatrix} \bar{x} \\ \bar{y} \end{pmatrix}$



Interpretation du coefficient de corrélation
 Soit α , angle que forment les vecteurs x et y centrés

$$\cos \alpha = \frac{\langle x - \bar{x}u, y - \bar{y}u \rangle}{\|x - \bar{x}u\|_{D_p} \cdot \|y - \bar{y}u\|_{D_p}} = r(x, y)$$

$\langle x, \bar{x}u; y, \bar{y}u \rangle_{D_p} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \text{cov}(x, y)$
 Le coefficient de corrélation linéaire entre x et y est égal au $\cos \alpha$, où α est le D_p angle formé de $\bar{x}u$ par les caractères centrés

Remarque.

$$-1 \leq r(x, y) \leq 1$$

$$r(x, y) = 0 \Leftrightarrow x - \bar{x}u \perp_{D_p} y - \bar{y}u$$

III) Echantillon de p variables.

I est un ensemble de n individus ou chacun des quels on a mesuré p caractères x^j , $1 \leq j \leq p$

$$I \xrightarrow{x^j} \mathbb{R}$$

$$i \mapsto x^j(i) = x_{ij}$$

On appelle échantillon les p variables $\{x^1, \dots, x^p\}$, l'ensemble $E = \{(x_{i1}, x_{i2}, \dots, x_{ip}), i = 1, \dots, n\}$.

Soit le tableau suivant.

$$X = \begin{pmatrix} 1 & x_{11}^1 & x_{11}^2 & \dots & x_{11}^p & x_{11}^p \\ \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\ i & x_{i1}^1 & x_{i1}^2 & \dots & x_{i1}^p & x_{i1}^p \\ \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\ n & x_{n1}^1 & x_{n1}^2 & \dots & x_{n1}^p & x_{n1}^p \end{pmatrix} \left. \vphantom{\begin{pmatrix} 1 \\ \vdots \\ i \\ \vdots \\ n \end{pmatrix}} \right\} n \text{ lignes}$$

p (Lignes) colonnes

La i^{eme} ligne de X symbolise l'individu " i "

La j^{eme} colonne de X symbolise le caractère α^j

Le vecteur $\underline{x}_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ij} \\ \vdots \\ x_{ip} \end{pmatrix} \in \mathbb{R}^p$ symbolise l'individu " i "

Le vecteur $\underline{x}_j = \begin{pmatrix} x_{1j} \\ \vdots \\ x_{ij} \\ \vdots \\ x_{nj} \end{pmatrix} \in \mathbb{R}^n$ symbolise le caractère α^j

\mathbb{R}^p est l'espace des "individus", l'échantillon E, y est représenté par un nuage de n pts

\mathbb{R}^n est l'espace des "caractères", l'échantillon E, y est représenté par un nuage de p pts

On pourra noter

$$X = (\underline{x}_1, \underline{x}_2, \dots, \underline{x}_i, \dots, \underline{x}_n)$$

$${}^t X = (\underline{x}_1, \underline{x}_2, \dots, \underline{x}_i, \dots, \underline{x}_n)$$

Matrice de variance, et matrice de corrélation, tableau de description

Les elts de l'ensemble I sont munis des poids, p_i ($p_i > 0$)

$$\sum_{i=1}^n p_i = 1$$

Soit D_p la matrice diagonale

$$D_p = \begin{pmatrix} p_1 & & 0 \\ & \ddots & \\ 0 & & p_n \end{pmatrix}$$

D_p est symétrique définie positive

D_p est matrice sur \mathbb{R}^n , espace des caractères

DS \mathbb{R}^p , le i^{th} moyen qui réalise E est :

$$g = \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_j \\ \vdots \\ \bar{x}_p \end{pmatrix} = \begin{pmatrix} \langle \underline{x}_1, u \rangle_{D_p} \\ \langle \underline{x}_2, u \rangle_{D_p} \\ \vdots \\ \langle \underline{x}_j, u \rangle_{D_p} \\ \vdots \\ \langle \underline{x}_p, u \rangle_{D_p} \end{pmatrix} \quad u \in \mathbb{R}^n \quad u = \begin{pmatrix} u_1 \\ \vdots \\ u_j \\ \vdots \\ u_p \end{pmatrix}$$

mettre l'origine de \mathbb{R}^p en g , c'est centrer les caractères.
Le tableau X est alors dit "centré".

Les matrices de covariance et de corrélation associées à l'échantillon E sont les matrices symétriques suivantes

$$V = \begin{pmatrix} V_{11} & V_{12} & \dots & V_{1p} \\ V_{21} & V_{22} & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots \\ V_{p1} & \dots & \dots & V_{pp} \end{pmatrix}$$

ou $V_{ij} = \text{cov}(x^i, x^j)$

$$R = \begin{pmatrix} 1 & r_{12} & \dots & r_{1p} \\ r_{21} & 1 & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots \\ r_{p1} & \dots & \dots & 1 \end{pmatrix}$$

ou $r_{ij} = r(x^i, x^j)$