

Statistiques

Chapitre 1 . Séries statistique – Généralité

I) Définition

- Population : ensemble d'éléments faisant l'objet d'une étude statistique.
- Echantillon : partie de la population.
- Unité statistique : Un élément de la population étudié dont on veut examiner un ou plusieurs caractères.
- Caractère : -qualitatif (nominale, ordinale) : non mesurable : (bon, mauvais; lieu d'habitation)
-quantitatif (discrète, continue) : mesurable à l'aide d'une variable statistique.
- Variable statistique : -discrètes : prend des valeurs isolées.
-continues : prend toutes les valeurs sur un intervalle.
- Séries statistique : ensemble des valeurs prises par une variable statistique sur l'ensemble de la population ou sur un échantillon.

II) Séries statistique d'un caractère quantitatif discret

Soit un échantillon de taille n (x_1, x_2, x_3, \dots) sont les valeurs positives du caractère x , avec p le nombre de valeurs possible des x .

- Séries statistique : les valeurs prise par les n membres de l'échantillon.
- Effectif total : n
- Effectif partiel de x_i (fréquence absolue) : c'est le nombre n_i de fois qu'apparaît la valeur x_i ($n_1 + n_2 + \dots + n_p = n$).
- Fréquence relative de x_i : $f_i = \frac{n_i}{n}$ ($f_1 + f_2 + \dots + f_p = 1$)
- Etude de la série : écart entre la plus grande et la plus faible des valeurs de x_i
- Exemple :

On considère 100 familles de 4 enfants.

On étudie le caractère : nombre de garçons.

Série statistique : 3, 1, 3, 2, 4, 0, 1, 2, ...

Effectif total : 100 étendu : $0 - 4 = 4$

On compte les valeurs :

| | | | | | | |
|----|---|----|----|----|---|-------|
| xi | 0 | 1 | 2 | 3 | 4 | Total |
| ni | 7 | 20 | 43 | 25 | 5 | 100 |

Fréquence relative de famille avec 2 garçons : $43/100 : 43\%$

III) Séries statistique d'un caractère quantitatif continu

Le nombre p de valeurs possible des x est infini (en réalité c'est fini à cause de la précision des mesures → les f_i sont presque nulles).

On va constituer les classes en divisant l'étendu de la série dans un certain nombre d'intervalles

Avec pour la classe i : un centre de classe x_i et un effectif u_i .

- Exemple : des nourrissons

Pesées effectuée à 10 g près. Donnée le poids entre 2,24 et 4,45.

| | | | | | | | | | |
|----------|---------|---------|----------|---------|---------|---------|---------|---------|--------|
| | 2,2-2,5 | 2,5-2,8 | 2,8 -3,1 | 3,1-3,4 | 3,4-3,7 | 3,7-4,0 | 4,0-4,3 | 4,3-4,6 | total |
| classe | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | |
| centre | 2,35 | 2,65 | 2,95 | 3,25 | 3,55 | 3,85 | 4,15 | 4,45 | |
| effectif | 5 | 11 | 24 | 40 | 42 | 20 | 13 | 6 | 161 |
| f_i | 3,1 | 6,8 | 14,9 | 24,8 | 26,1 | 12,4 | 8,1 | 3,7 | 99,9 % |

IV) Séries statistiques d'un caractère qualitatif

On groupe les résultats en autant de classes qu'il existe de modalités de caractères

- Variable cardinale : si les modalités peuvent être ordonnées (taille vestimentaire (+ ou - grande)).
- Variable nominale : si elles ne peuvent être ordonnées (couleurs)
- Variable dichotomiques : avec deux modalités.

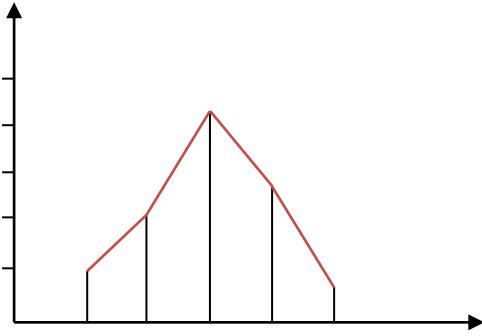
V) Représentation graphique des série statistiques

1) Cas des variables discrètes

- Diagramme en bâtons : (des effectifs et des fréquences absolues)

On trace $n_i = f(x_i)$ ou $f_i = f(x_i)$

Exemple des familles de 4 enfants : (x_i = nombre de garçons)



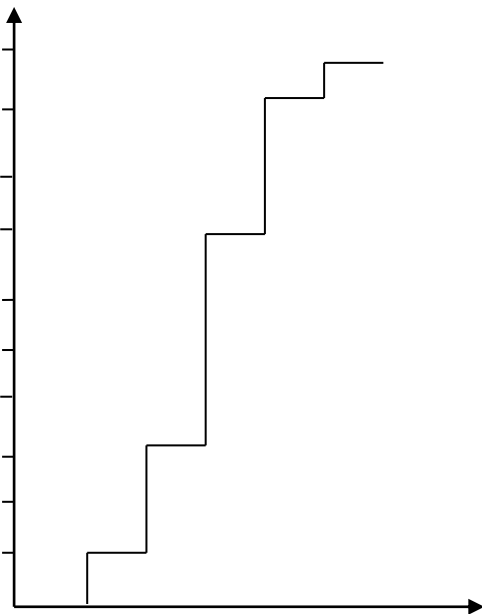
- Polygone : (des effectifs et de fréquences absolues)

On joint l'extrémité des bâtons par des segments (voir en rouge sur le graphique précédent)

- Diagramme cumulatif : (des effectifs et des fréquences)

Effectif cumulé jusqu'à la $i^{\text{ème}}$ valeur du caractère $n_1 + n_2 + \dots + n_i$

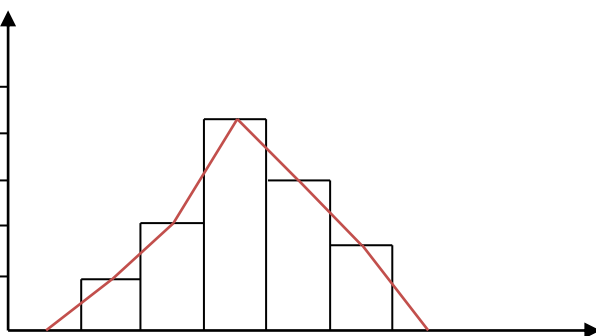
Fréquence relative cumulée jusqu'à la $i^{\text{ème}}$ valeur du caractère $f_1 + f_2 + \dots + f_i$



2) Caractère continu

- Histogramme :

On va partager le domaine des valeurs de x_i en classes : intervalles.



Si les classes sont inégales on les divise pour que toutes classes est la même 'largeur'.

- Polygones : (des effectifs est des fréquences)

Il s'agit d'une ligne brisée joignant les milieux des sommets des rectangles. (en rouge)

Chapitre 2 . Parammètre de position et de dispersion

I) Paramètre de position

1) Moyenne arithmétique

$$\text{Moyenne des } x = \frac{x_1+x_2+\dots+x_n}{n} = \frac{1}{n} \sum_{(i=1)}^n x_i$$

- Séries à caractère discret :

$$\text{Moyenne des } x = \frac{n_1x_1+n_2x_2+\dots+n_px_p}{n} = \frac{1}{n} \sum_{(i=1)}^p n_i x_i$$

- Séries à caractère continu :

$$\text{Moyenne des } x = \frac{n_1x_1+n_2x_2+\dots+n_px_p}{n} = \frac{1}{n} \sum_{(i=1)}^p n_i x_i \quad (x_i \text{ est le centre de classe})$$

On peut faire une moyenne sur le 'numéro de la classe' : (z_i : n° de classe ; M : minimum des x_i ; a : largeur de classe)

$$\text{si } x_i = M + a z_i \quad (z_i : \text{n}^\circ \text{ de classe}) \quad \text{moy}(x) = M + a * \text{moy}(z)$$

2) Médiane

On ordonne les valeurs de la série statistique et on prend la valeur du milieu.

Série discrète : série 1 : 3, 3, 4, **5**, 7, 7, 9 : médiane : Me = 5

Série 2 : 3, 3, 4, **5, 6**, 7, 7, 9 médiane : Me = (5 +6)/2 = 5,5

Série continu : On va considérer que la série est déjà ordonnée par classes et que la Me va tomber dans une classe qui correspond à la moitié de l'effectif total.

Me est, par approximation, dans la classe $[l_1, l_2[$ avec F1 effectif cummulé en l_1 et F2 en l_2 .

$$\text{Me} = l_1 + \left(\frac{n}{2} - F_1\right) \frac{l_2 - l_1}{F_2 - F_1}$$

3) Quartiles

Les quartiles divisent les effectifs cumulés en 4 avec Q1, Q2 = Me, Q3.

Même principe que la médiane sauf que pour la formule des séries continues il faut changer le $n/2$ par $n/4$ et $3n/4$.

4) Modes

Le mode est la valeur la plus fréquente (effectif le plus important)

Pour les séries continues, on a une classe dominante, mais on peut aussi trouver la valeur dominante.

$$\text{Mode : } Mo = l_1 + \frac{\Delta_1(l_2 - l_1)}{\Delta_1 + \Delta_2}$$

Où Δ_1 : excédant d'effectif de la classe modale par rapport à la classe inférieure et Δ_2 : excédant d'effectif de la classe modale par rapport à la classe supérieure

II) Paramètre de dispersion

1) Etendu

Ecart entre la plus grande et la plus petite valeur des x de la série.

2) Ecart moyen

Ecart de x_i : $e_i = |x_i - \text{moy}(x)|$

Ecart moyen de la série : $\text{moy}(e) = \frac{1}{N} \sum_{i=1}^n e_i$

Cas discret : $\text{moy}(e) = \frac{1}{n} \sum_{i=1}^p n_i \cdot e_i = \sum_{i=1}^p f_i \cdot e_i$

Cas continu : p classes de centre x_i

3) Variance et écart-type

Variance : moyenne arithmétique des carrés des écarts des x_i

Variance : $\sigma^2_x = \frac{1}{N} \sum_{i=1}^n (x_i - \text{moy}(x))^2$

Ecart-type : racine de la variance

Ecart-type : $\sigma_x = \sqrt{\frac{1}{N} \sum_{i=1}^n (x_i - \text{moy}(x))^2}$

Méthode rapide :

(Cas continu x_i : centre des classes)

$$\sigma^2x = \left(\frac{1}{N} \sum_{i=1}^n xi^2 \right) - \text{moy}(x)^2 \quad \rightarrow \quad \sigma^2x = \left(\frac{1}{N} \sum_{i=1}^p ni \cdot xi^2 \right) - \text{moy}(x)^2$$

remarque : si $xi = M + azi$ alors $\sigma x = a \sigma z$

4) Ecart-interquartile

L'écart-interquartile est représenté par $Q3 - Q1$ et il caractérise 50 % de la valeur de la série

5) Coefficient de variation

$$\text{Coefficient de variation} = \frac{\sigma x}{\text{moy}(x)}$$

Chapitre 3 . Ajustement de séries statistique

I) Méthode de lissage

1) Méthode des points médians

On trace l'enveloppe supérieur et inférieur avec les maximums/minimums locaux. Puis on trace un segment vertical pour chaque xi et on trace la nouvelle courbe qui passe par les milieux de chaque segments

2) Méthode de la moyenne mobile

On remplace y_i par $y'_i = \frac{y_{i-1} + y_i + y_{i+1}}{3}$ (on effectue une moyenne sur trois point).

Sauf pour le premier et le dernier point (moyenne sur deux points).

Remarque : Possibilité de pondérer les moyennes.

II) Méthode d'ajustement

1) Méthode graphique

On trace une droite ayant autant de point au dessus qu'au dessous. Ensuite si on connaît deux points de cette droite, on peut en déduire son équation.

2) Méthode des moindres carrés

On cherche $y'_i = f(xi)$ pour minimiser $\sum_{i=1}^p (y_i - y'_i)^2$

P = nombre de points.

Pour mesurer la qualité de l'ajustement on utilise le coefficient de corrélation tel que :

$$-1 \leq R \leq 1$$

$$R = \frac{\sum_i (x_i \cdot y_i) - \frac{1}{p} \sum_i x_i \sum_i y_i}{\sqrt{\left(\sum_i (x_i^2) - \frac{1}{p} (\sum_i x_i)^2\right) \left(\sum_i (y_i^2) - \frac{1}{p} (\sum_i y_i)^2\right)}}$$

a) Ajustement à l'aide d'une droite

Droite de la forme : $y = ax + b$ où a et b sont donnés par :

$$a = \frac{p \sum x_i y_i - \sum x_i \sum y_i}{p \sum x_i^2 - (\sum x_i)^2} \quad \text{et } b = \text{moy}(y) - a \cdot \text{moy}(x)$$

b) Ajustement à l'aide d'une parabole

Parabole la forme : $y = ax^2 + bx + c$ où a , b et c sont donnés par :

$$a \sum x_i^2 + b \sum x_i + cp = \sum y_i$$

$$a \sum x_i^3 + b \sum x_i^2 + c \sum x_i = \sum x_i y_i$$

$$a \sum x_i^4 + b \sum x_i^3 + c \sum x_i^2 = \sum x_i^2 y_i$$

c) Ajustement à l'aide d'une exponentielle

Exponentielle de la forme : $y = b \cdot a^x$

Donc $\rightarrow \log y = \log b + \log a^x$

$$Y = Ax + B \quad \text{où } y = \log y, B = \log b \text{ et } A = \log a$$

$$A = \frac{p \sum x_i Y_i - \sum x_i \sum Y_i}{p \sum x_i^2 - (\sum x_i)^2} \quad \text{et } B = \text{moy}(Y) - A \cdot \text{moy}(x)$$

d) Ajustement à l'aide d'une puissance

puissance de la forme : $y = b \cdot x^a$

Donc $\rightarrow \log y = \log b + \log x^a$

$$Y = aX + B \quad \text{où } y = \log y, X = \log x \text{ et } B = \log b$$

$$a = \frac{p \sum X_i Y_i - \sum X_i \sum Y_i}{p \sum X_i^2 - (\sum X_i)^2} \quad \text{et } B = \text{moy}(Y) - a \cdot \text{moy}(X)$$