

statistique descriptive

statistiques: données numériques

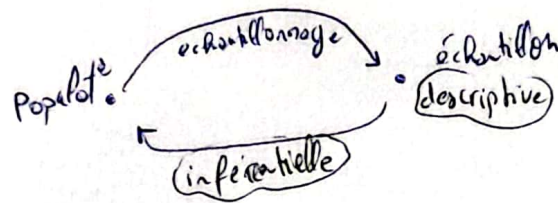
statistique = méthode d'analyse des données numériques

statistique descriptive

describ les caractéristiques d'un échantillon

statistique inférentielle

généraliser à partir d'un échantillon



Etude statistique (steps):

- ① problématique → ② collecte des données → ③ propriétés données
- ④ Analyse statistique (déductive, "globe", inductive, "généraliser") → ⑤ production résultat

collecte de Données

Population
 ↙ recensement (tous) "enquête exhaustive"
 ↘ échantillonnage (pas tous) "enquête partielle (sondage)"

méthode d'échantillonnage

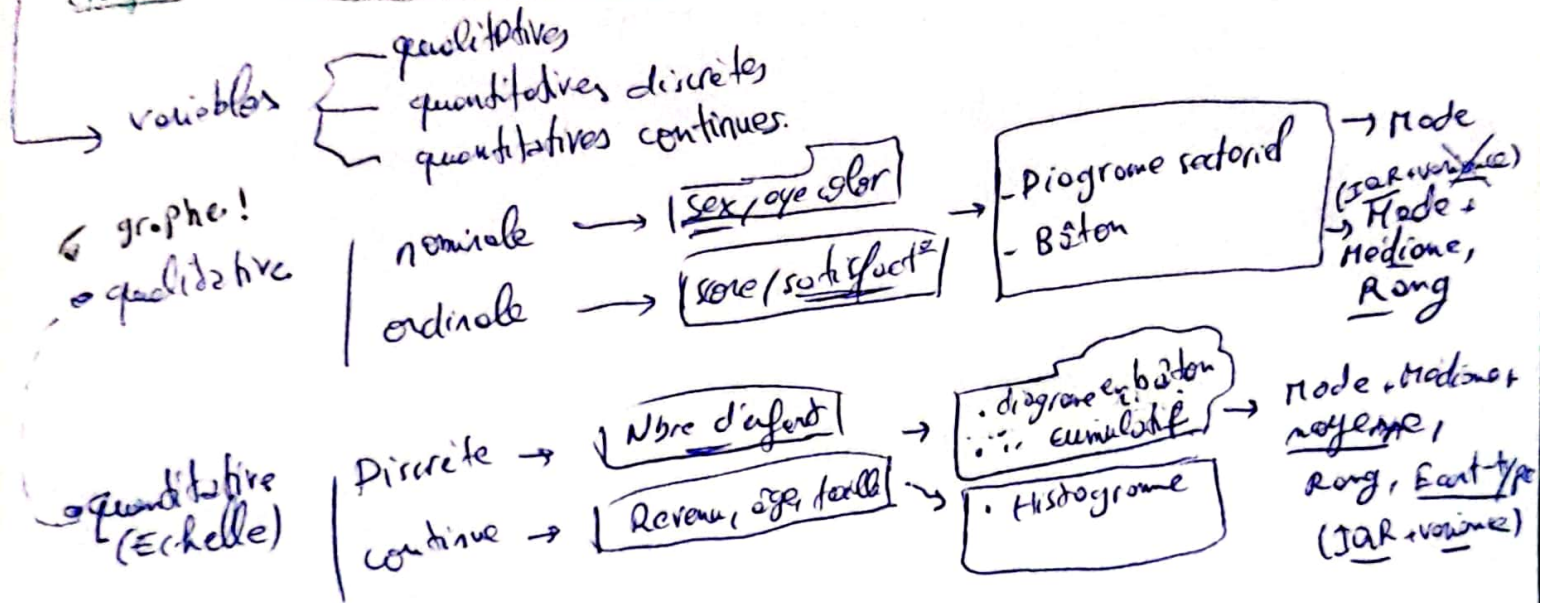
- Echantillonnage
 - aléatoire simple (avec / sans remise)
 "Sans remise est plus fréquent et plus précis"
 - stratifié ex par groupes
 "pour Population hétérogène"
 - à 2 ou plus niveaux
 "tirage au sort des familles"
 - méthode des quotas (quotas)
 "sondage d'opinion"

Précision de l'enquête
 dépend de l'échantillon (valeur, %)
 degré d'hétérogénéité (variabilité)

possibilité des observations
 (mesures / valeurs) — non-numeric (o qualitative (catégorielle)) — Nominales
 — numeric (o quantitative) — ordinales
 — continue
 — discrète
 (catégorie)
 (ordre de catégoris)
 (intervalle)
 (autres)
 Rang

Chap 1 Représentations graphiques des données :

Chap 2 Caractéristiques de positi^o, de dispersion et de forme :



effectifs cumulés

fréquences cumulées

indicateurs !

caractéristiques de 3 tendances centrales

moyenne arithmétique médiane mode

Dispersion

- Étendue
- écart type

indice de variabilité

N. # obs total = "effectif cumulé"

Forme

- quantiles
- indice de symétrie
- indice d'aplatissement

median (cas) valeurs de la variable

valeurs distinctes (N valeurs)

N impair

N pair

$$Q_2 = x_{(k)}$$

$$+9 \quad k = \frac{(N+1)}{2}$$

$$Q_2 = \frac{x_{(k)} + x_{(k+1)}}{2}$$

$$+9 \quad k = \frac{N}{2}$$

valeurs répétées

effectifs croissants

$\frac{\sum x_i}{N}$

Q_2 correspond

variable continue

go classe ouverte

$$Q_2 = \text{Barycentre} \left(\frac{N}{2} - F \right) / f_k$$

• quartiles Q_1 / Q_3 (1^{er} 3^{eu} quartile)

→ Q_1 $N_i < \frac{N}{4}$ Q_3 $N_i < \frac{3N}{4}$

→ rang $Q_1 = \frac{N}{4}$ → rang $Q_3 = \frac{3N}{4}$

Mode = valeur plus fréquente:

moyenne arithm.

discret (valeur)
 $\bar{X} = \frac{1}{N} \sum_{i=1}^N x_i$

continue "centre de classe"
 $\bar{X} = \frac{1}{N} \sum_{i=1}^n n_i x_i = \sum f_i x_i$

moyenne géométrique
 taux d'accroissement annuel

harmonique
 moyenne de %

quadratique
 mesure d'écart

Etendue

différence (max - min)

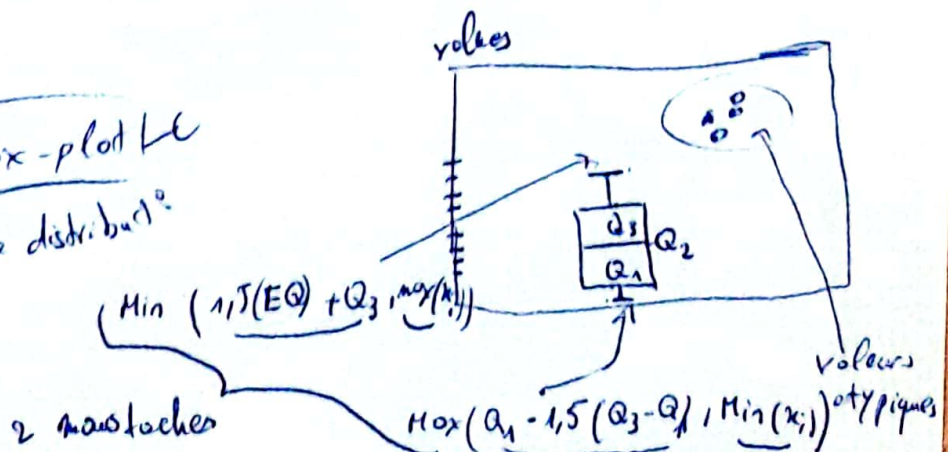
$s^2 = \frac{1}{N} \sum_{i=1}^n (x_i - \bar{x})^2$
 $= \sum_{i=1}^n f_i (x_i - \bar{x})^2$
 valeur & centre de classe

écart moyen
 sous-carré
 (comme normal)

Variances

Ecart interquartile $EQ = Q_3 - Q_1 \equiv \pm Q_R$ (range) $\equiv \pm q$
 (distributions asymétriques)

boîte à moustache **box-plot**
 représentation graphique d'une distribution



Ecart-type et variance sont des mesures de dispersion absolue

- ✓
 n. Coefficient de variat° \equiv indice de variabilité
 " dispersion relative pour comparer 2 distribut°s"

$$CV = \frac{S}{\bar{X}} (100) \quad \equiv \frac{\text{écart-type}}{\text{moyenne}} \times 100$$

→ dispersion relative plus grande que celle de la distribut°
 \Rightarrow moins homogène ...

Mesures de forme:
 Coefficient d'asymétrie (SKWENESS)

$$AS = \frac{m_3}{s^3}; \quad m_3 = \frac{1}{N} \sum_{i=1}^k n_i (x_i - \bar{x})^3 = E(X - E(X))^3$$

$AS = 0$
 distribut° symétrique

$AS > 0$
 distribut° présente une asymétrie à gauche

$AS < 0$
 distribut° présente une asymétrie à droite

Coeficient d'aplatissement : (KURTOSIS)

$$AP = \frac{m_4}{s^4}; \quad m_4 = \frac{1}{N} \sum_{i=1}^k n_i (x_i - \bar{x})^4$$

$AP = 3$
 Distribut° a un coef. AP similaire à la distribut° normale :
 distribut° mésocurtique

$AP > 3$
 Distribut° sera plus taillée que la distribut° normale avec des queues épaisses :
 distribut° leptocurtique

$AP < 3$
 Distribut° présente des queues plus fines que celle loi normale :
 distribut° platycurtique

- Distribut^o Normale
- chap 3: lien entre les variables

1. Distribut^o Normale : (Gauss)

"forme de cloche"

"distribut^o de référence"



→ propriétés :

- unimodale et symétrique (moyenne = mode = médiane)
- 150% > moyenne
- 50% < moyenne

$(\mu - \sigma, \mu + \sigma)$
68% populat^o

• écart type et moyenne

$(\mu - 2\sigma, \mu + 2\sigma)$

95.44% populat^o

$(\mu - 3\sigma, \mu + 3\sigma)$
99.74%

→ vérificat^o de la normalité d'une variable :

• cloche

• unimodale

• AS = 0

• AP = 3

• Programmes

PP Probability-Probability plot
QA Quantile-Quantile plot

Normale
+ unimodale
↓
moyenne
= mode
= médiane

1. Scores standardisés

1 score Z :

score centré et réduit à partir d'un score individuel.
donne une indicat^o précise de la posit^o du score de l'individu
au sein de la distribut^o.

$$Z = \frac{x - \mu}{s}$$

distribut^o normal centrée réduite (système de mesure) où :
→ mode = moyenne = médiane = \square
→ $s = 1$ (toujours) "écart type"

Chap 3

- ① liaison entre 2 variables qualitatives
- coefficient χ^2
 - Représentation graphique

- ② entre 2 quantitatives
- coefficient de corrélation
 - représentation graphique

- ③ entre 1 qualit. et 1 quantit.
- représentation graphique

① tableau de contingence

X \ Y	y ₁	y _j	y _L	total
x ₁				n ₁
⋮				
x _i		n _{ij}		
⋮				
x _k				
total		n _j		n

lignes
totales
des
individus

$$n_i = \sum_j n_{ij}$$

$$n = \sum_{ij} n_{ij}$$

$$= \sum_i n_i$$

$$= \sum_j n_j$$

$$\chi^2 = \sum_{i=1}^k \sum_{j=1}^L \left(n_{ij} - \frac{n_{i.} n_{.j}}{n} \right)^2 \left(\frac{n_{i.} n_{.j}}{n} \right)$$

(valeur théorique si indépendance)

0
indépendance
(profils identiques)

$\neq 0, \nearrow$ si profils différents

Diagrammes à bâtons superposés :

- ②
- tableau de données (forme de n couples (x_i, y_i) $i=1, \dots, n$)
 - Représentation graphique : Nuage de points M;
 - coefficient de corrélation, $r \in [-1, 1]$
 - " indicateur de variété de 2 variables quantitatives continues"
 - si sens corrélation positif
 - si sens corrélation négative

$$\boxed{r_{xy} = \frac{\text{cov}(x,y)}{s_x s_y}} ; \quad \text{cov}(x,y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

r → son unité, indépendant de l'origine choisie.

si il est proche de -1 ou 1 : il y a une relation presque affine entre x et y

⑧

si x qualitative à k modalités, l'ensemble des n individus peut être subdivisé en k groupes sur lesquels x est constante

si y est quantitative, les k groupes peuvent être représentés par une boîte à moustache de façon à pouvoir les comparer.

? (girl, boy) ?