

# Data Exploration

## objectif:

Introduire des méthodes qui nous permet d'extraire de l'information pertinente d'un jeu de données.

A l'issue de cette formation, vous serez capable de:

→ Faire apparaitre des comportements particuliers des objets observés (détecter les individus ayant un comportement atypique, trouver des ensembles d'individus ayant un comportement similaire)

→ Trouver des liens entre les variables étudiées

→ utiliser R.

## Corpus / jeu de données:

un jeu de données est un tableau avec:

- en ligne  
unité statistique)

les variables (attributs étudiés)  
- en colonne

Exp:

①

	couleur yeux	couleur cheveux	taille
Karine	Noir	Blond	164
Olivier	Bleu	Noir	174
Sarah	Maron	Noir	158
Paul	Noir	Brun	180

Dans notre exemple:

+ les individus sont les étudiants

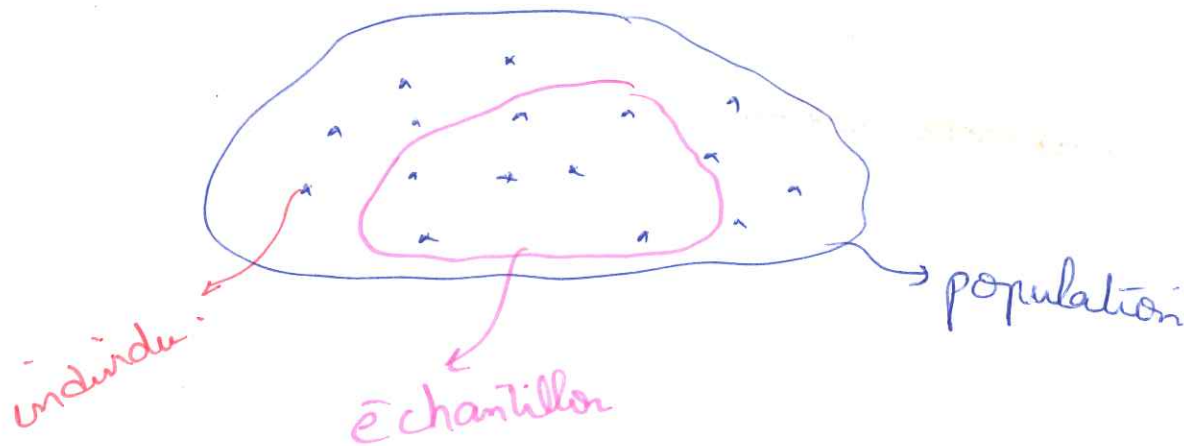
+ les variables sont: couleur de yeux, couleur de cheveux, taille,

→ A partir de ce tableau, tirer les informations importantes pour qu'il soit compréhensible pour tout le monde.  
donc on fait un résumé numérique (Moyenne, Médiane...) ou une représentation graphique (diagramme en barre, histogramme ---) présentat. synthétique.

## Population et échantillon:

→ population = Ensemble des individus

→ Échantillon = Sous ensemble de la population.



→ Recensement: Etude de tous les individus

→ Sondage: sous ens de la population.



Si la population infinie  
alors le Recensement est impossible  
→ on fait un sondage.

# 2 types d'études statistiques

## Analyse exploratoire (Statistique descriptive)

- (tableau de contingence)
- Résumé numérique
- Représentation graphique
- Recherche de sous groupes homogènes

## Statistique inférentielle

- Trouver des estimateurs sans biais et efficace pour passer de l'échantillon à la population.

Statistique inférentielle

Echantillonnage aléatoire

Statistique descriptive

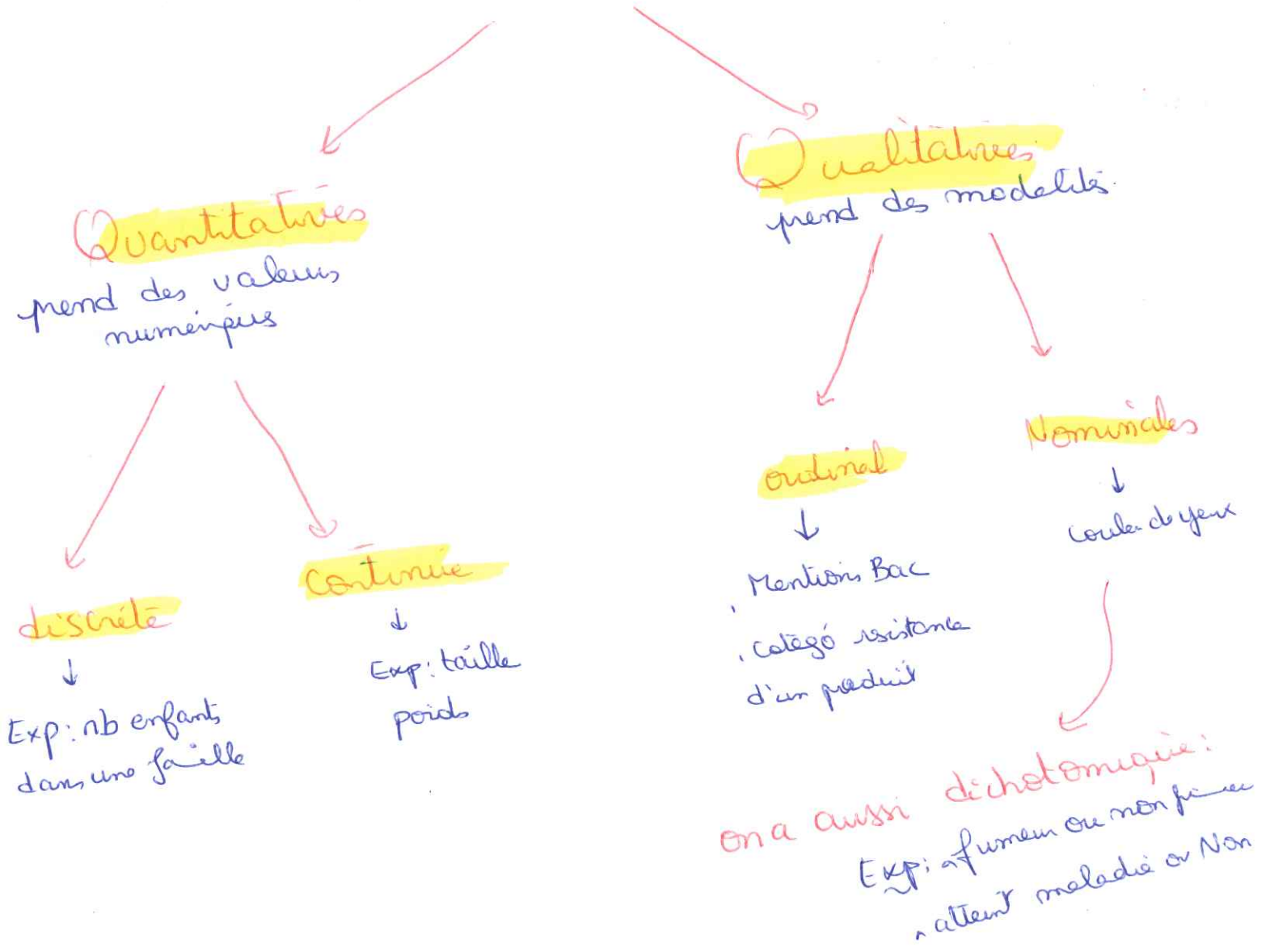
Analyse exploratoire

Caractéristiques



→ Chaque individu est décrit par un ensemble de caractéristiques appelé "Variables" ou attributs.

## Nature Variables



# Types d'analyse.

## Statistique uni variées



On s'intéresse à une seule variable

## Statistique Bi variées



On s'intéresse à deux variables simultanément et on fait étude selon la nature de ces deux variables.

## Statistique multi variées.



On s'intéresse à plusieurs ou p variables.

Dans la suite on s'intéresse à :

## Statistique uni variée



Représentat synthétique

Représentat graphique

Représentat numérique.

# Représentation Synthétique

Exp:  
21

étudiant	Mention Bac	Couleur Yeux	Couleur Cheveux
A	Bien	Vert	Marron
B	Bien	Noir	Noir
C	Bien	Bleu	Noir
D	T.B	Vert	Noir
E	AB	Noir	Noir
F	Bien	Noir	Marron
G	AB	Marron	Noir
H	T.B	Noir	Marron
I	Bien	Bleu	Blond
J	T.B	Noir	Noir

→ jeu de données

À partir de "Jeu de données", on peut créer pour chaque variable un "tableau de contingence".

Exp: Variable "Mention en Bac":  
 \* les modalités sont: AB - Bien - TB  
 \* On compte pour chaque modalité, le nombre des individus (ici étudiants) ont ce modalité.

Exp: on a: 5 étudiants ont mention Bien en Bac.  
 ou 5 c'est l'effectif de la modalité Bien.

\* On fait pareil pour les 2 autres modalités "AB" et "TB" et on stocke tout dans un tableau →

(5)

Mention	effectif	freq
$i=1 \leftarrow$ AB	2 $\sim n_1=2$	$\frac{2}{10} = 0.2$ ou 20%
$i=2 \leftarrow$ B	5 $\sim n_2=5$	$\frac{5}{10} = 0.5$ ou 50%
$i=3 \leftarrow$ TB	3 $\sim n_3=3$	$\frac{3}{10} = 0.3$ ou 30%
total	10	

car on a 10 étudiants

Fréquence =  $f_i = \frac{n_i}{n}$  eff chq modalité / total  
 ou  $f_i$  en pourcentage =  $f_i \times 100$

Interprétation; Il ya 50% des étudiants qui ont mention "Bien" en Bac.

## Représentation graphique

### ① Variables qualitatives:

① Variables nominales : diag en barre ou diagramme circulaire.

② Variables ordinales : Diag en bâtons



## ② Variables ordinales: Diagramme en bâton

### Variables quantitatives:

Discrete

diagramme en bâton

nb états attribués à 10 hôtels parisiens

Exp: nb états

nb état	Effectif
1	5
2	3
3	1
4	1

total 10

eff (ou fréq)

Diagramme en bâton.



Effectif (ou fréquence)  
augmente  
⇓

Hauteur du rectangle aussi

Continue

histogramme.

Exp: Ici les valeurs sont regroupées par des intervalles.

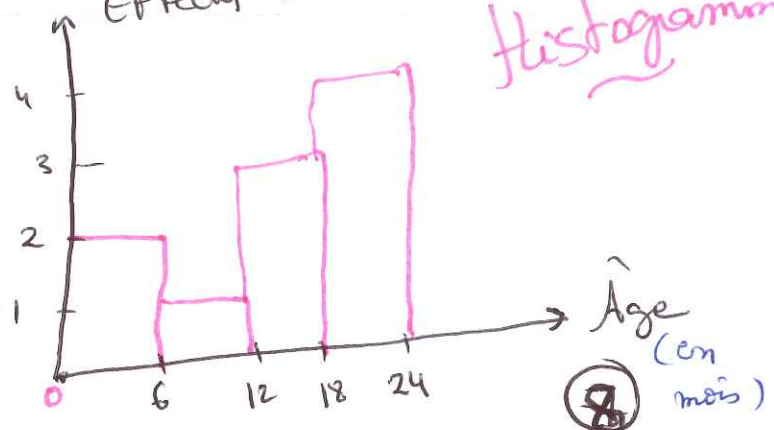
Exp: Âge enfants crèche municipale de 10 enfants (en mois)

centre classe

âge	Effectif
[0, 6[	2
[6, 12[	1
[12, 18[	3
[18, 24[	4
total	10

Effectif (ou fréq)

histogramme





# Résumés Numériques.



Que pour les variables quantitatives.



on distingue deux types de Résumé numérique

Indicateur de position

(Moyenne, mode, médiane, quantiles...),

Ils positionnent la série de valeurs observées autour d'une tendance centrale.

Exp: Moy. Âge des étudiants

Indicateur de dispersion

(Variance, écart type, étendue, interquartile...)

Indiquent la fluctuation des valeurs de la série autour d'une tendance centrale en générale.

## ★ Indicateurs de position:

① Mode: La valeur observée qui a le plus grand effectif.

Exp: nb étoile : Mode est 1 (qui a eff le plus haut = 5)

• Si données sont quantitatives continues, on parle alors du classé modal.

Exp: note, classe modale est: [18, 24[



⚠ • "La mode" est le seul indicateur où on peut calculer pour les variables qualitatives aussi:

Exp: Mention, le Mode est "Bien"

## ② Moyenne:

① Discrete:

$$\bar{x} = \frac{1}{n} \sum_i n_i \cdot x_i$$

Exp: (Etoile):

$$\begin{aligned} \bar{x} &= \frac{1}{10} [(1 \times 5) + (2 \times 3) + (3 \times 1) + (1 \times 4)] \\ &= 1.8 \end{aligned}$$

② Continue:

$$\bar{x} = \frac{1}{n} \sum_i n_i \cdot x_i$$

Ici le  $\underline{x_i}$  est le centre du classe.

⚠ classe  $[a, b] \Rightarrow \text{centre}(x_i) = \frac{a+b}{2}$

Exp: (Crèche):

$$\begin{aligned} \bar{x} &= \frac{1}{10} [(3 \times 2) + (9 \times 1) + (15 \times 3) + (4 \times 2)] \\ &= 14,4 \text{ mois} \end{aligned}$$

### ③ Quantiles:

• Médiane: c'est la valeur qui sépare une population en deux groupes d'effectifs égaux. Elle n'a pas de sens que si les données sont rangées par ordre croissant.  
Ex: Médiane des notes en Maths est 13; c'-à-d:

50% des étudiants ont une note inférieure à 13.  
et 50% (le moitié) ont une note supérieure à 13.

### • Quantiles:

① 1<sup>er</sup> quantile  $Q_1$ : 25% des valeurs sont inf à  $Q_1$  (ou les 3/4 des valeurs sont sup à  $Q_1$ ); donc le  $\frac{3}{4}$  restant sont sup à  $Q_1$

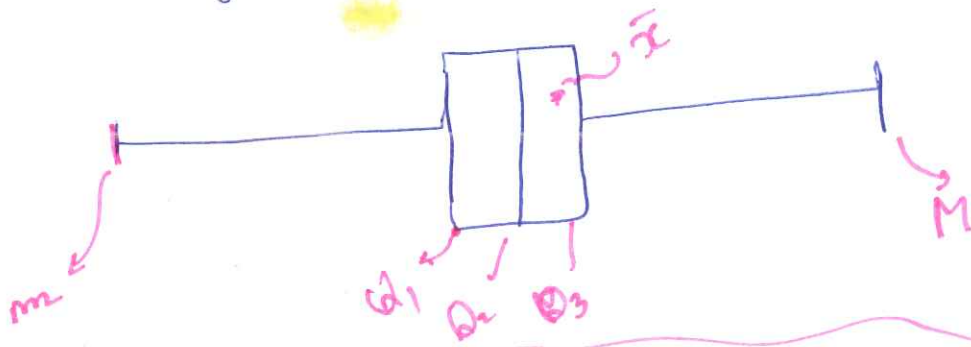
② 2<sup>e</sup> quantile  $Q_2$ : 50% des valeurs sont inf à  $Q_2 \Rightarrow$  Médiane

③ 3<sup>e</sup> quantile  $Q_3$ : 75% des valeurs sont inf à  $Q_3$ : [ou  $\frac{3}{4}$  des valeurs inf à  $Q_3$ ]  
alors le  $\frac{1}{4}$  restant sont sup à  $Q_3$ .



# Boîte de Tukey (Boîte à Moustaches, Boxplot).

Représente les indicateurs de position :



$m$  et  $M$  sont les moustaches:

$$\begin{cases} m = Q_1 - 1,5 \cdot (Q_3 - Q_1) \\ M = Q_3 + 1,5 \cdot (Q_3 - Q_1) \end{cases}$$

Les valeurs de la série en dehors des moustaches est considérée comme "atypique" ou valeur aberrante.



Boxplot nous permet d'avoir une aperçue graphique rapide de la distribution des valeurs de la série

Exp:

Note	1	2	3	4	5	6	7
Elève L	9	10	8	7	10	9	11
Elève P	14	2	16	5	6	5	16

Série ordonnée :

					Me		
Elève L:	7	8	9	9	10	10	11
Elève P:	2	5	5	6	14	16	16
					Me		

Elève L: Moyenne  $\bar{x} = \frac{1}{7} (7+8+\dots+11) = 9,1$

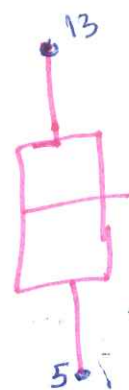
• Médiane  $Me = 9$

•  $Q_1 = 8$

•  $Q_3 = 10$

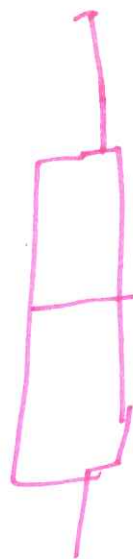
$$m = Q_1 - 1,5(Q_3 - Q_1) = 8 - 1,5(10 - 8) = 5$$

$$M = Q_3 + 1,5(Q_3 - Q_1) = 10 + 1,5(10 - 8) = 13$$



⚠ pas de valeurs aberrantes car tous les notes de l'élève sont entre 5 et 13.

Elève P:



$$\left\{ \begin{array}{l} \bar{x} = 9,1 \\ Me = 6 \\ Q_1 = 5 \\ Q_3 = 14 \\ m = 0 \\ M = 27,5 \end{array} \right.$$

13

## \* Indicateurs de dispersion:

① Variance: Mesure l'écart au carré entre les valeurs de la série et leur moyenne.

$$V = S^2 = \frac{1}{n} \sum n_i (x_i - \bar{x})^2$$

⚠ tout le temps positif.

Exp: (étoile)

$$S^2 = \frac{1}{10} \left[ 5(1-1,8)^2 + 3(2-1,8)^2 + 1(3-1,8)^2 + 1(4-1,8)^2 \right]$$

• Ecart-type:

$$S = \sqrt{S^2}$$

• Ecart-Médiane: Mesure l'écart entre les valeurs de la série et leur médiane;

$$e_m = \frac{1}{n} \sum n_i |x_i - \text{med}|$$

• Etendue:  $\max \{x_i\} - \min \{x_i\}$

Exp (étoile):  $4 - 1 = 3$

• Ecart-inter-quartiles:

$$Q_3 - Q_1$$



## Variables centrées - réduites.

On définit la série centrée réduite de la façon suivante:

$$\tilde{x}_i = \left( \frac{x_i - \bar{x}}{s_x} \right)$$

→ Moyenne

→ écart-type

→ La série est dite :

- . Centrée car elle a Moyenne nulle.
- . Réduite car de variance 1.

