



**HAL**  
open science

## Cours de statistique descriptive

Bardin Bahouayila

► **To cite this version:**

| Bardin Bahouayila. Cours de statistique descriptive. DEUG. Congo-Brazzaville. 2016. cel-01317598

**HAL Id: cel-01317598**

**<https://hal.science/cel-01317598>**

Submitted on 2 Aug 2016

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



REPUBLIQUE DU CONGO

Institut Africain de la Statistique

(IAS)

Option: HDTS

Année académique : 2015/2016

# STATISTIQUE DESCRIPTIVE

Rédigé par :

**BAHOUAYILA MILONGO Chancel Bardin<sup>1</sup>**

<sup>1</sup> E-mail : [bardinbahouayila@yahoo.fr](mailto:bardinbahouayila@yahoo.fr) / [bardin.bahouayila@facebook.com](https://www.facebook.com/bardin.bahouayila)

Tel : 05 075 33 71 / 06 837 81 85

# Sommaire

<b>INTRODUCTION</b> .....	1
<b>CHAPITRE 1</b> .....	2
I-1. LA POPULATION .....	2
I-2. L'UNITÉ STATISTIQUE OU L'INDIVIDU .....	2
I-3. L'ÉCHANTILLON .....	2
I-4. LE CARACTÈRE OU LA VARIABLE .....	2
I-5. LA MODALITÉ .....	3
I-6. LA DISCRÉTISATION .....	3
<b>CHAPITRE 2</b> .....	4
II-1. LES FREQUENCES ABSOLUE, RELATIVE ET CUMULEE .....	4
II-2. LA MOYENNE .....	5
II-3. LE MODE .....	7
II-4. LA MEDIANE .....	8
II-5. LES FRACTILES .....	8
<b>CHAPITRE 3</b> .....	9
III-1. L'ÉTENDU ET LE RAPPORT DE VARIATION .....	9
III-2. L'INTERVALLE INTERQUARTILE .....	11
III-3. LA VARIANCE ET L'ÉCART-TYPE .....	11
III-4. LE COEFFICIENT DE VARIATION .....	13
<b>CHAPITRE 4 :</b> .....	14
IV-1. LE CAS DES VARIABLES CONTINUES .....	14
IV-2. LE CAS DES VARIABLES DISCRETES .....	15
IV-3. LE CAS DES VARIABLES QUALITATIVES .....	15

# INTRODUCTION

En présence d'un ensemble de données chiffrées l'esprit a un besoin spontané de simplification. Selon les critères qui lui sont propre, il cherche d'une part à représenter et à classer ces données ; d'autre part, il souhaite résumer la multiplicité et la complexité des notations par des caractéristiques synthétiques. De ce fait, l'homme est conduit à déterminer **les caractéristiques centrales** (moyenne, médiane, etc.), à construire des **graphiques** (histogramme, camembert, etc.), à calculer des **caractéristiques de dispersion** (écart-type, rapport de variation, intervalle interquartile, etc.) et à comparer des « séries statistiques ». C'est en voulant tout cela qu'est née la notion de statistique descriptive. Le but de la **statistique descriptive** est donc de décrire des données en mettant de l'ordre et une certaine régularité; c'est comme si l'on faisait le résumé du livre : le résumé à l'avantage d'être plus court, plus facile à lire et comporte les éléments essentiels, mais le résumé néglige certains aspects pour faciliter la lecture. Ceci dit, en dehors de la statistique descriptive, il existe la **statistique inférentielle** qui permet de savoir à quel point l'on peut résumer sans perdre des informations essentielles et quel est le meilleur résumé avec le moins d'erreur. Cette branche des statistiques s'intéresse davantage à **extrapoler** des résultats issus d'échantillons en vue de caractériser une population mère inconnue, de faire des **prévisions** de comportements basées sur le calcul de probabilités. Malheureusement, dans ce cours, nous ne nous focaliserons que sur la statistique descriptive.

Ce cours est destiné en priorité à un public n'ayant aucune formation en statistique et cependant confronté de façon récurrente à la manipulation et à l'analyse de séries de données.

Aucun pré-requis en mathématique n'est exigé si ce n'est la connaissance des opérations mathématiques de base. Volonté, curiosité et ténacité permettront de maîtriser sans encombre les notions abordées qui, malgré leur complexité apparente, demeurent relativement simples. Cette formation se présente davantage comme une initiation à la rigueur que nécessite la manipulation d'ensembles de données afin d'utiliser à bon escient les méthodes appropriées pour éviter de faire parler faussement les chiffres.

Les concepts et méthodes statistiques seront abordés au travers de nombreux exemples.

Au final, il s'agira de se familiariser avec les données et de connaître la méthode statistique en général en vue de décrire, de résumer et d'analyser une population ou un ensemble de données.

# CHAPITRE 1

## PRÉSENTATION DES DONNÉES

La **statistique** est une méthode scientifique qui consiste à réunir des données chiffrées sur des ensembles nombreux, puis à analyser, à commenter et à critiquer ces données. En d'autres termes, c'est une science qui a pour objectif :

- ➔ la planification du projet ;
- ➔ la collecte, la codification, la saisie, le traitement et l'analyse des données ;
- ➔ la publication des résultats.

Il ne faut pas confondre **la statistique** qui est la science qui vient d'être définie et **une statistique** qui est un ensemble de données chiffrées sur un sujet précis.

Les premières statistiques correctement élaborées ont été celles des recensements démographiques. Ainsi le vocabulaire statistique est essentiellement celui de la démographie.

Les ensembles étudiés sont appelés **population**. Les éléments de la population sont appelés **individus ou unités statistiques**. La population est étudiée selon un ou plusieurs **caractères**.

### I-1. LA POPULATION

C'est l'ensemble des individus (ou unités statistiques) présentant un caractère commun. Pour une thématique donnée, la population regroupe toujours la totalité des individus relatifs à cette thématique (notion d'exhaustivité).

*Exemples* : la population congolaise, les pays de la CEMAC, les clients d'une entreprise.

La population est en général notée P

L'effectif total d'une population est noté N.

### I-2. L'UNITÉ STATISTIQUE OU L'INDIVIDU

C'est l'élément de base constitutif de la population à laquelle il appartient. Il est indivisible et peut être un pays, un végétal, un humain ou une entreprise.

### I-3. L'ÉCHANTILLON

C'est un sous-ensemble construit et représentatif d'une population donnée.

### I-4. LE CARACTÈRE OU LA VARIABLE

C'est la (les) caractéristique(s) de l'individu intégrant la population étudiée.

*Exemple* : la couleur, le sexe, le poids, la taille, la marque, le modèle, l'espèce, le prix, la surface, etc.

#### I-4.1 Variable qualitative

Une variable statistique est dite de nature qualitative si ses modalités ne sont pas mesurables.

Les modalités d'une variable qualitative sont les différentes catégories d'une nomenclature. Ces catégories doivent être exhaustives (chaque individu est affecté à une modalité) et incompatibles (un individu ne peut être affecté à plusieurs modalités) de façon à créer une partition.

Le sexe, la profession, l'état matrimonial sont quelques exemples de variables qualitatives.

Pour ses enquêtes auprès des ménages, l'Insee utilise la nomenclature des Professions et catégories socioprofessionnelles (PCS-2003). Les modalités d'une variable qualitative peuvent être classées sur deux types d'échelle : nominale ou ordinale. À ces deux types d'échelle correspondent deux types de variables qualitatives.

### ➔ Variable qualitative nominale

Une variable statistique qualitative est dite définie sur une échelle nominale si ses modalités **ne sont pas naturellement ordonnées**.

Exemples : Situation d'activité, statut matrimonial.

### ➔ Variable qualitative ordinale

Une variable statistique qualitative est dite ordinale si l'ensemble de ses modalités peut être doté d'**une relation d'ordre**.

Exemple : Niveau d'instruction.

## **I-4.2 Variable quantitative**

Toute variable qui n'est pas qualitative ne peut être que quantitative. Les différentes modalités d'une variable quantitative constituent l'ensemble des valeurs numériques que peut prendre la variable.

Une variable statistique est dite de nature quantitative si ses modalités sont mesurables. Les modalités d'une variable quantitative sont des nombres liés à l'unité choisie, qui doit toujours être précisée.

Il existe deux types de variables quantitatives : les variables discrètes et les variables continues.

Ces variables ont en commun des modalités clairement ordonnées, pour lesquelles l'écart entre les valeurs possède une signification, et sur lesquelles il est possible de réaliser des opérations mathématiques telles que des calculs de moyennes, etc. Néanmoins, elles ont des propriétés et des traitements spécifiques qui nécessitent une étude séparée.

### ➔ Variable quantitative discrète

Lorsque les modalités sont des valeurs numériques isolées, comme le nombre d'enfants par ménage, on parle de variable discrète.

Exemples : Âge, salaire, nombre de lit dans un hôpital.

### ➔ Variable quantitative continue

Lorsque la variable, par exemple la taille d'un individu, peut prendre toutes les valeurs d'un intervalle, ces valeurs peuvent alors être regroupées en classes, et on parle dans ce cas de variable continue.

Exemples : Poids, taux du sucre, taille, taux du sel.

## **I-5. LA MODALITÉ**

C'est la valeur qualitative ou quantitative que peut prendre le caractère précédemment défini.

*Exemple* : sexe féminin ou masculin, poids 45 kg, couleur verte, etc.

Attention, les modalités sont exhaustives et mutuellement exclusives. Chaque individu doit pouvoir être classé dans une et une seule modalité.

## **I-6. LA DISCRÉTISATION**

Lorsque les modalités sont des valeurs numériques isolées, comme le nombre d'enfants par ménage, on parle de variable discrète.

Ce découpage en classes pose de nombreuses questions : choix des amplitudes, amplitudes constantes ou variables, nombre de classes, etc. Nous ne rentrerons pas ici dans le détail de ces opérations.

## CHAPITRE 2

# CARACTÉRISTIQUES DE TENDANCE CENTRALE DES DONNÉES

Les paramètres de tendance centrale ou « mesures de tendance centrale » sont des grandeurs susceptibles de représenter au mieux un ensemble de données. L'appellation « tendance centrale » vient du fait que ces paramètres donnent une idée de ce qui se passe au centre d'une distribution, d'un ensemble de données.

On distingue trois mesures de tendance centrale :

- \* La moyenne ;
- \* Le mode ;
- \* La médiane.

Tous trois ne décrivent par la même chose et sont, de ce fait, complémentaires dans la description et l'analyse d'une distribution.

**Ces statistiques ne se calculent que dans le cas où nous avons à faire à des variables quantitatives.** Dans le cas où nous avons des variables qualitatives, on procède aux fréquences.

Avant d'analyser ces trois indicateurs de position, nous allons d'abord aborder la notion de la fréquence.

### II-1. LES FREQUENCES ABSOLUE, RELATIVE ET CUMULEE

A chaque modalité de variable X, peut correspondre un ou plusieurs individus dans l'échantillon de taille n.

On appelle effectif de la modalité  $x_i$ , le nombre  $n_i$ . Il est aussi appelé **fréquence absolue**.

**La fréquence relative** est le nombre  $f_i$  tel que  $f_i = \frac{n_i}{n}$

La fréquence cumulée croissante est cependant le nombre  $F_i$  tel que  $F_i = \sum_{p=1}^i f_p$

**Exemple** : Représentons la fréquence relative et la fréquence cumulée du tableau ci-dessous

$X_i$	$n_i$
1	8
2	18
3	14
4	10
Total	50

**Solution** :

$X_i$	$n_i$	$f_i$	FCC	FCD
1	8	$8/50=0,16$	0,16	1
2	18	$18/50=0,36$	$0,16+0,36=0,52$	$1-0,16=0,84$
3	14	$14/50=0,28$	$0,52+0,28=0,8$	$0,84-0,36=0,48$
4	10	$10/50=0,2$	$0,8+0,2=1$	$0,48-0,28=0,2$
Total	50	$50/50=1$		

## II-2. LA MOYENNE

La moyenne constitue l'un des paramètres fondamentaux de tendance centrale mais non suffisant pour caractériser une distribution. Complémentaire du mode et surtout de la médiane, la moyenne constitue à n'en point douter, la mesure la plus calculée et la plus utilisée lors de la description de séries statistiques.

Il existe plusieurs types de moyennes, chacun adapté à des situations précises :

DESIGNATION	NOTATION COURANTE
Moyenne arithmétique	$\bar{X}$
Moyenne géométrique	$\bar{G}$ ou $\bar{x}_G$
Moyenne harmonique	$\bar{H}$ ou $\bar{x}_H$
Moyenne quadratique	$\bar{Q}$ ou $\bar{x}_Q$

**Attention !** On ne peut pas calculer la moyenne sur des données qualitatives.

### II-2.1 La moyenne arithmétique

C'est la plus simple et la communément utilisée et ce, pas toujours à bon escient. Elle se note la plupart du temps par  $\bar{X}$ . Elle peut être simple ou pondérée.

#### → La moyenne arithmétique simple

Sa version simple correspond à une somme de résultats divisée par le nombre de résultats et s'écrit :

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

#### → La moyenne arithmétique pondérée

La moyenne arithmétique pondérée, autant le dire tout de suite, donne, dans son utilisation classique (c'est-à-dire lorsque tous les individus ont le même poids), le même résultat que la moyenne arithmétique simple. Sa formule est cependant différente puisqu'elle introduit la notion de poids via un terme supplémentaire qui peut s'avérer utile dans certaines situations, notamment lorsque justement les individus composant une population n'ont pas le même poids ou coefficient : certains individus, pour diverses raisons, ont davantage d'influence dans ladite population que les autres. Ce peut être le cas par exemple lorsque l'on a affaire à une série de notes dont le coefficient n'est pas le même. Cette moyenne s'écrit de la manière suivante :

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

- Exemples :**
- calculer la moyenne des valeurs suivantes : 800, 400, 200, 1000
  - exemple des notes de la classe.
  - exemple du calcul de la moyenne de 8, 9, 3, 5, 5, 4, 6, 4, 7, 9



## II-2.2 La moyenne géométrique

Sa définition purement mathématique est un peu rébarbative mais son utilité est grande comme nous allons le démontrer.

La moyenne géométrique de  $n$  valeurs positives  $x_i$  est la racine  $n^{\text{ième}}$  du produit de ces valeurs. Elle est notée  $\bar{G}$  et s'écrit :

$$\bar{G} = \sqrt[n]{\prod_{i=1}^n x_i}$$

La moyenne géométrique est un instrument permettant de calculer des taux moyens notamment des taux moyens annuels. Son utilisation n'a un sens que si les valeurs ont un caractère multiplicatif.

**Exemple :** a) les prix de l'immobilier ancien ont augmenté ces trois dernières années de la façon suivante : 2, 4, 8.

b) exemple du taux de **pauvreté** moyen.

## II-2.3 La moyenne harmonique

On utilise la moyenne harmonique lorsqu'on veut déterminer un rapport moyen dans des domaines où il existe des liens de proportionnalité inverse.

**Exemples :**

- Pour une distance donnée, le temps de trajet est d'autant plus court que la vitesse est élevée.
- Un loyer dans le parc privé est d'autant plus élevé que la taille ou la surface du logement est petite.

Cette moyenne s'écrit de la manière suivante :

$$\bar{H} = \frac{n}{\sum_{i=1}^n \frac{1}{x_i}}$$

## II-2.4 La moyenne quadratique

Une moyenne qui trouve des applications lorsque l'on a affaire à des phénomènes présentant un caractère sinusoïdal avec alternance de valeurs positives et de valeurs négatives. Elle est, de ce fait, très utilisée en électricité. Elle permet notamment de calculer la grandeur d'un ensemble de nombre. A titre d'information, elle s'écrit :

$$\bar{Q} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i^2)}$$

### II-3. LE MODE

Le mode,  $M_o$  d'une série statistique est la valeur du caractère la plus fréquente ou dominante dans l'échantillon. Autrement dit, c'est la valeur qui a la fréquence (absolue ou relative) la plus élevée. Lorsque la distribution a plus d'un mode, on parle d'une distribution « **multimodale** » (bimodale, trimodale, etc). Par contre, si l'on est en présence de données groupées en classes, le mode se rapportera à la classe comportant le plus grand nombre d'individus : on parlera alors de **classe modale**. Cependant, il peut y arriver que l'on s'intéresse à avoir la valeur approchée ou exacte de ce mode. Par conséquent, il est recommandé d'appliquer la démarche suivante :

- Pour avoir une valeur approximative du mode, on calcule la moyenne de la classe qui a la fréquence la plus élevée ;
- Pour avoir une valeur exacte, le mode se calcule de la manière suivante

$$M_o = x_m + \frac{i\Delta i}{\Delta s + \Delta i}$$

Avec

$X_m$  : limite inférieure de la classe modale ;

$i$  : amplitude de la classe modale ;

$\Delta i$  : écart d'effectif entre la classe modale et la classe inférieure la plus proche

$\Delta s$  : écart d'effectif entre la classe modale et la classe supérieure la plus proche

#### Exemples :

a) Donner le mode des séries des données suivantes

$$S_1 = \{9, 10, 9, 9, 11, 10, 11, 11, 11, 9\}$$

$$S_2 = \{2, 3, 5, 5, 6, 7, 7, 8, 8, 9, 2, 2, 1, 3, 3, 4, 2, 1, 1, 1\}$$

b) donner le mode des données se trouvant dans le tableau ci-après

Classes	Effectifs
[0-5[	10
[5-10[	15
[10-15[	12
[15-25[	25
[25-30[	10

## II-4. LA MEDIANE

Dans le calcul de la médiane, on distingue deux cas :

\* Si la variable est discrète

On désigne par  $n$  le nombre d'observations.

- ✓ Si  **$n$  est pair** : la médiane est alors égale à la moyenne des valeurs encadrant le milieu de la série.
- ✓ Si  **$n$  est impair** alors il est possible d'identifier simplement la valeur qui partage la population en deux effectifs égaux. Le rang central étant égal à  $[(n+1)/2]$ .

\* Si la variable est continue et qu'elle est groupée en classe.

On cherche la classe contenant le  $\frac{n}{2}$  individu de l'échantillon. Cette classe est appelée la **classe médiane**. En supposant que tous les individus de cette classe sont uniformément répartis à l'intérieur, la médiane se calcule de la façon suivante par interpolation linéaire :

$$M_e = x_m + a \left[ \frac{\frac{n}{2} - N_i}{n_i} \right]$$

$x_m$  : limite inférieure de la classe médiane ;

$a$  : amplitude de la classe médiane ;

$n_i$  : effectif de la classe médiane

$N_i$  : effectif cumulé inférieur à  $x_m$

$n$  : taille de l'échantillon

### Exemples :

a) Calculer la médiane des séries suivantes :

$S_1 = \{9, 10, 9, 9, 11, 10, 11, 11, 11, 9\}$

$S_2 = \{2, 3, 5, 5, 6, 7, 7, 8, 8, 9, 2, 2, 1, 3, 3, 4, 2, 1, 1\}$

b) Calculer la médiane des données se trouvant dans le tableau suivant :

$X_i$	$n_i$	$N_i$	$f_i$	FCC
0-10	48	48	0,24	0,24
10-15	40	88	0,2	0,44
15-20	56	144	0,28	0,72
20-30	32	176	0,16	0,88
30-50	24	200	0,12	1
Total	200		1	

## II-5. LES FRACTILES

Dans une distribution dont les individus ont été au préalable triés par ordre croissant, les fractiles correspondent aux valeurs qui partagent une population en sous-ensembles de même taille, c'est-à-dire d'effectifs égaux.

Il existe plusieurs fractiles à savoir: les **quartiles**, **déciles**, les **centiles**, etc.

Les quartiles sont respectivement  $Q_1$ ,  $Q_2$  et  $Q_3$ . Ils représentent respectivement 25 %, 50 % et 75 % des effectifs de la population.

De la même manière, et dans le but de préciser et d'affiner encore l'analyse de la dispersion d'une distribution, on peut faire appel aux notions de déciles et de centiles. Le principe demeure le même que pour les quartiles à la différence que la population est ici divisée respectivement en 10 et 100 sous-populations d'égal effectifs:

## CHAPITRE 3

Les indices de tendance centrale définissent le comportement général des données.

Mais les données peuvent varier beaucoup autour de cette tendance. On doit donc définir un indice qui caractérise la variabilité des données dans l'échantillon. Cet indice est appelé indice de dispersion parce qu'il renseigne sur la dispersion ou l'éparpillement des données autour notamment des paramètres de tendance centrale.

Nous étudierons quatre paramètres de dispersion parmi les principaux, en mettant plus particulièrement l'accent sur la variance et l'écart-type :

- ➔ l'étendue et le rapport de variation
- ➔ l'intervalle interquartile
- ➔ la variance et l'écart-type
- ➔ le coefficient de variation

### III-1. L'ETENDU ET LE RAPPORT DE VARIATION

Le Minimum et le maximum d'une série statistique correspondent respectivement, comme leur nom l'indique, aux valeurs minimale et maximale rencontrées dans une série. Ces deux paramètres ont une double utilité. Ils permettent de calculer:

- ➔ l'étendue de la distribution, également appelée Intervalle de Variation (IV), c'est-à-dire l'écart entre le minimum et le maximum.

$$\text{Etendu} = \text{Maximum} - \text{Minimum}$$

- ➔ le Rapport de Variation (RV), c'est-à-dire le rapport de la valeur maximale de la distribution à la valeur minimale de la même distribution.

$$\text{Rapport de variation} = \frac{\text{Maximum}}{\text{Minimum}}$$

**Exemple 4** : Les notes d'élèves de deux classes au même examen ont donné les résultats suivants.

Classe 1	Classe 2
8	3
11	12
13	16
5	5
8	3
14	7
6	10
12	7
5	19
10	16
16	5
7	11
12	13
13	11
8	9
13	13
8	9
7	10
13	12
13	8
9	15
17	15
10	8
13	
6	
13	
7	
14	

**Solution :**

	Classe 1	Classe 2
<b>Minimum</b>	<b>5</b>	<b>3</b>
<b>Maximum</b>	<b>17</b>	<b>19</b>
<b>Etendu</b>	<b><math>17-5=12</math></b>	<b><math>19-3=16</math></b>
<b>Rapport de variation</b>	<b><math>17/5=3,4</math></b>	<b><math>19/3=6,3</math></b>

Le rapport de variation nous apprend que dans la classe 1 la meilleure note est 3,4 fois plus élevée que la note la plus faible. Ce rapport est plus important dans la classe 2 pour laquelle il est 6,3.

### III-2. L'INTERVALLE INTERQUARTILE

C'est la différence entre le troisième quartile ( $Q_3$ ) et le premier quartile ( $Q_1$ ).

Noté IQ, il s'écrit :  $IQ = Q_3 - Q_1$

### III-3. LA VARIANCE ET L'ECART-TYPE

Considérons une distribution pour laquelle on a calculé les paramètres de tendance centrale comme la médiane ou la moyenne. Comme leurs noms l'indiquent, et comme mentionné plus haut, ces mesures caractérisent le centre de la distribution. Parmi celles-ci, considérons la moyenne comme une référence.

Que pensez-vous de l'écart entre chaque valeur de la distribution et cette moyenne?

$$(x_i - \bar{x})$$

Plus cet écart sera faible, plus la valeur  $x_i$  sera proche de la moyenne et donc du centre de la distribution. A contrario, plus l'écart sera important et plus  $x_i$  sera éloignée du centre de la distribution. La prise en compte de la somme de l'ensemble des écarts à la moyenne, c'est-à-dire de la somme de tous les écarts entre les  $x_i$  et la moyenne donne logiquement 0.

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

Si l'on veut tenir compte de l'ensemble des distances à la moyenne sans pâtir d'une somme nulle, résultat de la compensation entre écarts négatifs et écarts positifs, il est nécessaire d'élever au carré chaque écart de telle sorte que l'on est :

$$\sum_{i=1}^n (x_i - \bar{x})^2 \neq 0$$

Que pensez-vous alors de la moyenne calculée de ces écarts élevés au carré?

$$\sigma^2 = \frac{1}{n} * \sum_{i=1}^n (x_i - \bar{x})^2$$

Ce paramètre  $\sigma^2$  est la variance. La variance satisfait à toutes les exigences énoncées plus haut relativement à la mesure de la dispersion d'une distribution. La variance pose toutefois le problème de proposer un résultat en unité élevée au carré. Si les données  $x_i$  sont en euros, la moyenne sera en euros, de même que l'écart  $(x_i - \bar{x})$  alors que la variance sera en euros carrés.

Il faut noter que la valeur de  $\sigma^2$  est la variance de la population; la variance de l'échantillon est de ce fait :

$$s^2 = \frac{1}{n-1} * \sum_{i=1}^n (x_i - \bar{x})^2$$

Pour revenir à l'unité initiale, il faut extraire la racine carrée de la variance. C'est ce qui nous donne l'écart-type.

Comme ce fut le cas pour le calcul de la moyenne de données groupées, pour calculer la variance des données groupées, il faut prendre en compte le centre de chaque classe et considérer que les individus d'une même classe ont tous la même valeur, celle du centre de leur classe.

**Exemple** : Trouver la variance et l'écart-type de la distribution suivante :

Classes	$n_i$	Centre des classes ( $c_i$ )	$n_i * c_i$	$c_i$ -moyenne	$n_i * (c_i - \text{moyenne})^2$
[4-6[	2	$(4 + 6) / 2 = 5$	$2 * 5 = 10$	$5 - 15,5 = -10,5$	$2 * (-10,5)^2 = 220,5$
[6-10[	5	$(6 + 10) / 2 = 8$	$5 * 8 = 40$	$8 - 15,5 = -7,5$	$5 * (-7,5)^2 = 281,25$
[10-20[	8	$(10 + 20) / 2 = 15$	$8 * 15 = 120$	$15 - 15,5 = -0,5$	$8 * (-0,5)^2 = 2$
[20-30[	4	$(20 + 30) / 2 = 25$	$4 * 25 = 100$	$25 - 15,5 = 9,5$	$4 * (9,5)^2 = 361$
[30-50[	1	$(30 + 50) / 2 = 40$	$1 * 40 = 40$	$40 - 15,5 = 24,5$	$1 * (24,5)^2 = 600,25$
<b>Total</b>	<b>20</b>		<b>310</b>		<b>1465</b>

**Solution** :

La moyenne de ces données est :  $\bar{x} = \frac{310}{20} = 15,5$

La variance est  $\sigma^2 = \frac{1}{n} * \sum_{i=1}^n [n_i * (c_i - \bar{x})^2] = \frac{1465}{20} = 73,25$

L'écart-type est  $\sigma = \sqrt{\sigma^2} = \sqrt{73,25} = 8,56$

### III-4. LE COEFFICIENT DE VARIATION

L'écart-type, malgré sa pertinence dans la mesure de la dispersion d'une distribution, possède un inconvénient majeur : il est exprimé dans l'unité de la variable à laquelle il se rapporte. Il est alors impossible de comparer les dispersions de deux distributions ayant un lien entre elles (lien de causalité ou autre) et dont les valeurs s'expriment dans des unités différentes.

Pour comparer la dispersion de deux séries qui ne sont pas exprimées dans les mêmes unités, on utilise le coefficient de variation. Cette statistique est une mesure neutre qui s'exprime la plupart du temps en pourcentage. Il se calcule en divisant l'écart-type par la moyenne et s'écrit donc :

$$CV = \frac{\sigma}{\bar{x}}$$

Plus grand est le coefficient de variation, plus grande est la dispersion.

**Exemple** : Trouver le coefficient de variation des données suivantes :

N°	Age	Salaire
1	37	300
2	35	310
3	36	290
4	36	305
5	41	305
6	38	295
7	40	300
8	36	310
9	35	290
10	37	295

**Solution** :

	Age	Salaire
<b>Moyenne</b>	<b>37,1</b>	<b>300</b>
<b>Variance</b>	<b>3,69</b>	<b>50</b>
<b>Ecart-type</b>	<b>1,92</b>	<b>7,07</b>
<b>CV</b>	<b>0,05</b>	<b>0,02</b>

Nous constatons que l'écart-type de l'âge est 1,92 ans et l'écart-type du salaire est 7,07FCFA. On ne peut pas dire que les salaires sont plus dispersés que les âges car les deux variables n'ont pas les mêmes unités (ans pour l'âge et FCFA pour le salaire). Pour comparer la dispersion dans les deux variables, il faut faire l'analyse avec le coefficient de variation car cette statistique n'a pas d'unité. Il ressort de ce fait que la dispersion est plus énorme dans les salaires que dans les âges.



## CHAPITRE 4 :

# REPRÉSENTATION DES DONNÉES

Les graphiques sont les corollaires d'une bonne analyse et d'une interprétation la plus complète possible de séries statistiques ou de résultats sur des traitements de données. Ces modes de représentation de la donnée participent à la compréhension des phénomènes, au même titre que les tableaux simples ou élaborés, apportant une information certes agrégée, synthétique mais très visuelle et en cela plus facile à aborder et à interpréter que ne le ferait un tableau de chiffres.

Chaque type de graphique est adapté à une ou plusieurs situations ou façon de représenter l'information. Selon la nature des données, la nature de variable, le nombre de variables et ce que l'on souhaite montrer, il sera judicieux de choisir la représentation graphique la mieux adaptée.

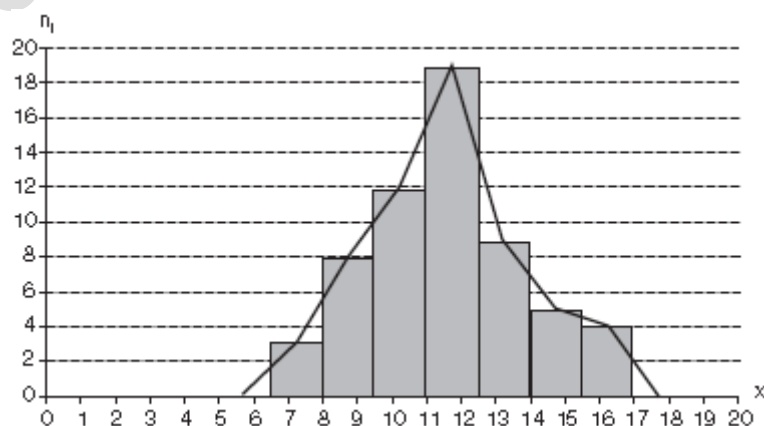
### IV-1. LE CAS DES VARIABLES CONTINUES

Quand on a une variable sous forme de classe (une variable discrétisée), on ne peut que faire l'histogramme. Un histogramme est un diagramme composé de rectangles contigus dont les aires sont proportionnelles aux effectifs (ou aux fréquences) et dont les bases sont déterminées par les intervalles de classes.

**Exemple :** Le responsable des ressources humaines d'une entreprise a relevé la distribution statistique suivante correspondant à l'ancienneté du personnel cadre dans l'entreprise, exprimée en années :

Classes	Effectifs
[6,5 ; 8[	3
[8 ; 9,5[	8
[9,5 ; 11[	12
[11 ; 12,5[	19
[12,5 ; 14[	9
[14 ; 15,5[	5
[15,5 ; 17[	4
Total	60

**Solution :**



## IV-2. LE CAS DES VARIABLES DISCRETES

Quand on a une variable discrète, on peut faire le diagramme en bâton, le diagramme en ligne ou le nuage de points.

### Diagramme en bâton

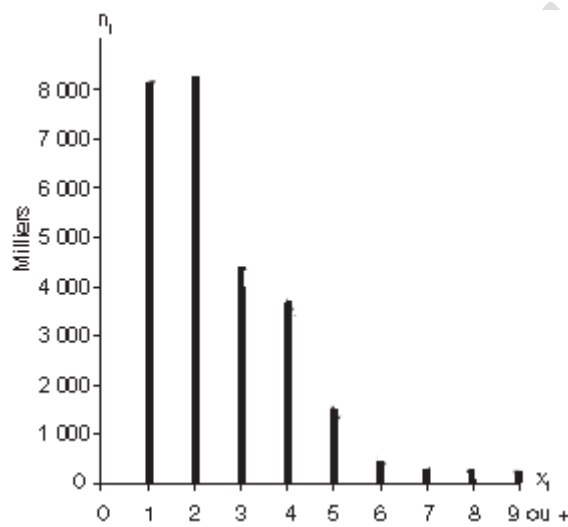
On appelle diagramme en bâtons un graphique qui à chaque modalité d'une variable quantitative discrète associe un segment (bâton) dont la hauteur est proportionnelle à l'effectif (ou à la fréquence).

Exemple :

Faites le diagramme en bâtons du nombre de personnes par ménage en France 1999 se trouvant dans le tableau ci-dessous.

$X_i$	1	2	3	4	5	6	7	8	9 ou plus
$n_i$	8000	8100	4500	3500	1500	500	300	200	300

Solution :



## IV-3. LE CAS DES VARIABLES QUALITATIVES

### → Diagramme circulaire (camembert)

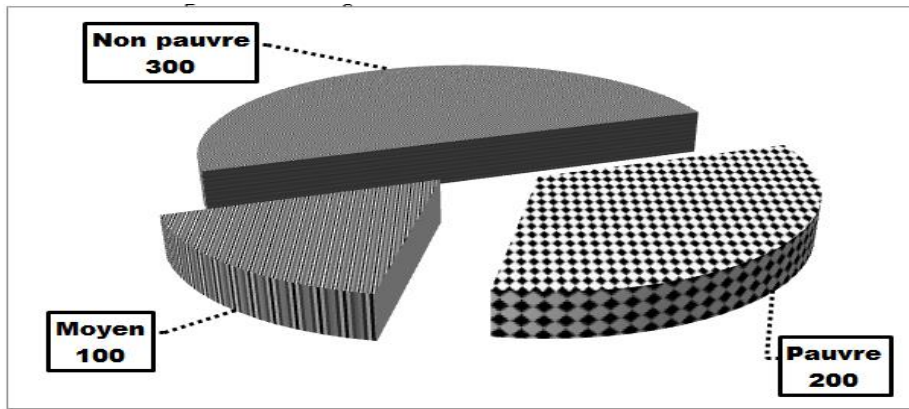
Un **diagramme circulaire** est un graphique constitué d'un cercle divisé en secteurs dont les angles au centre sont proportionnels aux effectifs (ou aux fréquences). De fait, les aires des secteurs sont proportionnelles aux effectifs. L'angle  $\alpha_i$  d'une modalité d'effectif  $n_i$  est donné en degrés par :

$$\alpha_i = \frac{n_i}{n} \times 360 = f_i \times 360$$

Exemple :

	$n_i$	$f_i$	$\alpha_i$
Pauvre	200	0,3	120
Moyen	100	0,2	60
Non pauvre	300	0,5	180
Total	600		

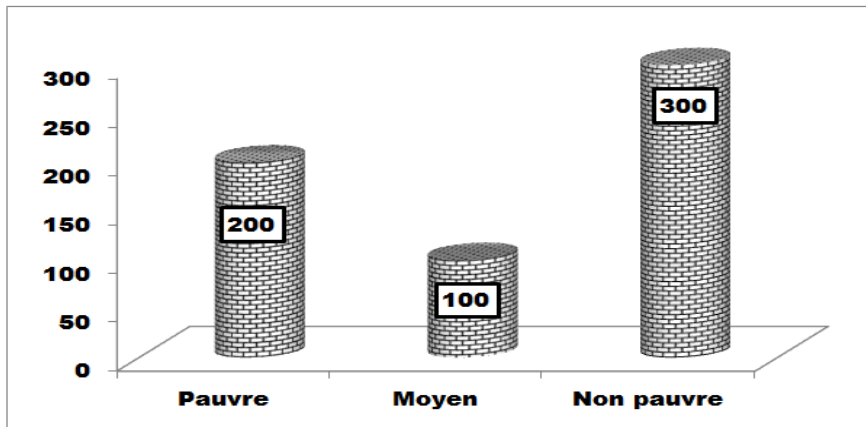
Solution :



➔ Diagramme en tuyau d'orgue

Un **diagramme en tuyaux d'orgue** est un graphique qui à chaque modalité d'une variable qualitative associe un rectangle de base constante dont la hauteur est proportionnelle à l'effectif (ou à la fréquence). De fait, les aires des secteurs sont proportionnelles aux effectifs. Les rectangles sont en général disjoints, verticaux ou horizontaux.

Exemple :

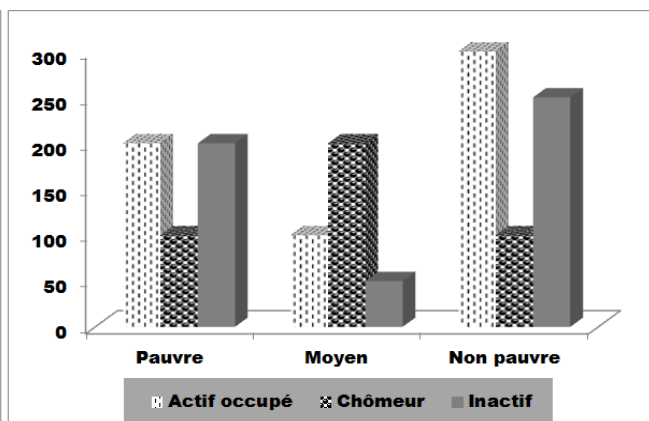
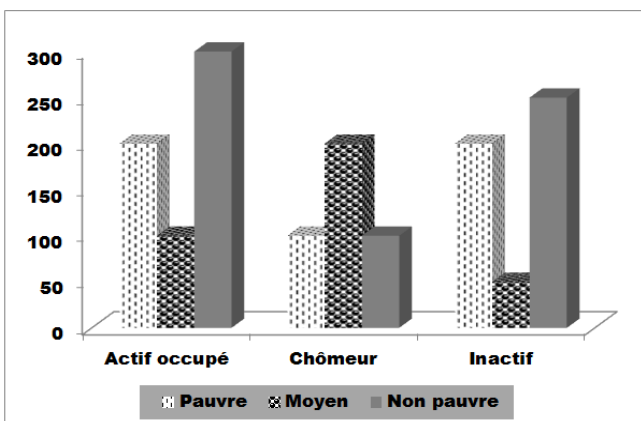


➔ Diagramme en barre multiple

Exemple :

	Actif occupé	Chômeur	Inactif	Total
Pauvre	200	100,0	200	500
Moyen	100	200,0	50	350
Non pauvre	300	100,0	250	650
Total	600	400	500	1500

Solution



SI VOUS AVEZ BESOIN DE LA CORRECTION DES EXERCICES PROPOSÉS, DE NOS EXERCICES DE TRAVAUX DIRIGÉS, DE NOS SUJETS DE DEVOIR OU D'EXAMEN AVEC SOLUTION, CONTACTEZ NOUS :

95, rue Malanda (Moukondo vers la Tsiémé, Ouenzé)

VOUS POUVEZ AUSSI APPELER OU ECRIRE A L'AUTEUR DE CE DOCUMENT.

E-mail : [bardinbahouayila@yahoo.fr](mailto:bardinbahouayila@yahoo.fr) / [bardinbahouayila@gmail.com](mailto:bardinbahouayila@gmail.com)

Tel : 05 075 33 71 / 06 837 81 85