

# Méthodes quantitatives

Paul Hobeika

2022-01-28



# Contents

<b>À propos de ce document</b>	<b>5</b>
<b>1 Données et vocabulaire de la statistique</b>	<b>7</b>
1.1 Les sources statistiques en sociologie . . . . .	7
1.2 Le vocabulaire de la statistique . . . . .	9
1.3 Variables . . . . .	12
1.4 Mesures de tendance centrale . . . . .	14
<b>2 Statistique descriptive univariée</b>	<b>15</b>
2.1 Variables qualitatives . . . . .	15
2.2 Variables quantitatives . . . . .	17
2.3 La loi normale : une distribution importante . . . . .	27
<b>Références</b>	<b>31</b>



# À propos de ce document

Cette page accueille les notes de cours de méthodes quantitatives du M1 de Science Po Strasbourg pour l'année 2021-2022. Il s'agit d'une introduction aux statistiques destinée à des étudiants de master de sociologie politique. Elle ne requiert pas de bagage préalable en statistique ou en mathématique. Il a été généré par l'extension `bookdown` de Yihui Xie, et le code source est disponible sur [GitHub](#).



# Chapter 1

## Données et vocabulaire de la statistique

### 1.1 Les sources statistiques en sociologie

Nous avons évoqué la semaine dernière l'importance de la connaissance des sources statistiques pour la production de savoirs quantitatifs en sciences sociale. Il en existe différents types, qu'il est important de savoir identifier.

#### 1.1.1 Les enquêtes par questionnaire produites par les chercheur-es

C'est par exemple le cas des données exploitées dans *La distinction* [Bourdieu, 1979] dont nous avons parlé au premier semestre. À partir d'une problématique de départ parfois abstraite (dans le sens pas directement quantifiable), l'élaboration d'un questionnaire a souvent pour objectif de trouver des éléments empiriques concrets qui permettent de rendre opérationnelles certaines notions ou concepts. Par exemple, dans *La distinction*, le questionnaire porte sur les pratiques culturelles et permet d'opérationnaliser empiriquement la notion de *capital culturel*.

Remarque : si vous souhaitez produire vous-même des données dans le cadre de votre TER et de la validation du cours c'est tout à fait possible, mais nous n'aborderons pas la méthodologie du questionnaire dans ce cours. De bons manuels sont toutefois disponibles, je vous recommande par exemple celui de Bugeja-Bloch and Couto [2021], chapitres 3 et 4.

### 1.1.2 Les autres source de “première main”

En réalité, les chercheur-es peuvent effectuer des traitement quantitatifs sur d’autres types de sources que les données issues d’un questionnaire. Pour cette raison, Fanny Bugeja-Bloch et Marie-Paule Couto font une distinction entre les **techniques d’enquête** et les **techniques d’analyse** des données.

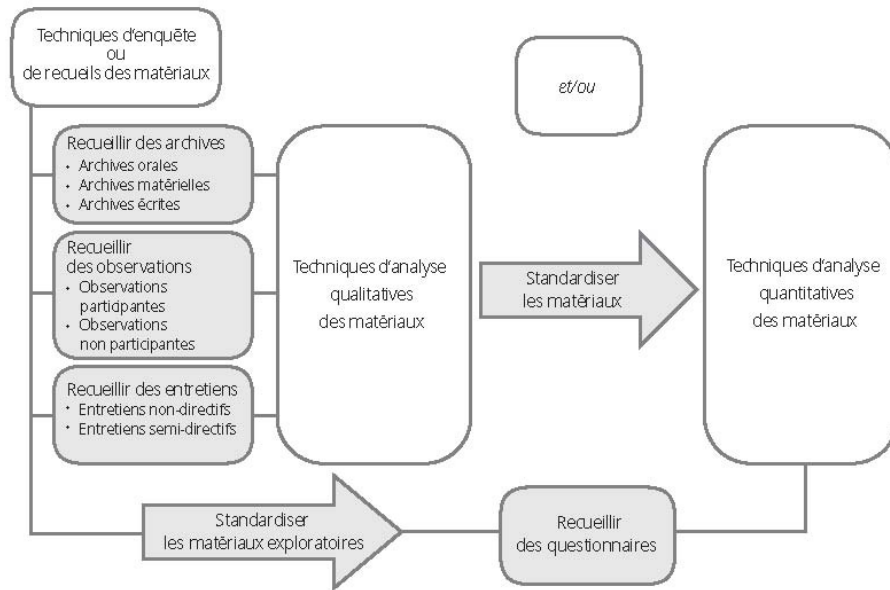


Figure 1.1: Techniques d’enquête et techniques d’analyse [Bugeja-Bloch and Couto, 2021]

Les **techniques d’enquête** désignent les différents modes de recueil des données : données d’archives, d’entretien ou encore issues d’observation. Les matériaux ainsi produits peuvent ensuite se prêter à différentes formes d’analyse. C’est seulement à ce niveau que l’on peut distinguer méthodes qualitatives et quantitatives. Une analyse qui se fondera sur le commentaire d’un ou quelques extrait d’entretien par exemple sera qualifiée de *qualitative*. Mais ces mêmes matériaux, lorsqu’ils sont *standardisés* et *mis en série* peuvent également être l’objet de techniques d’analyse quantitative. On peut produire des statistiques à partir d’archives [Lemerrier and Zalc, 2008], à partir d’entretiens (le questionnaire en est un cas particulier) ou encore à partir d’observations <sup>1</sup>.

<sup>1</sup>Un exemple tiré de la sociologie du travail est celui de l’enquête de Jean Peneff sur les urgences. Effectuant une enquête par observation participante en tant que brancardier dans un service d’urgence, il fait un certain nombre de comptages dans l’objectif d’objectiver certaines dimensions du travail aux urgences [Peneff, 1992].



### 1.1.3 L'analyse secondaire des données

Dans de nombreux cas, ce ne sont pas les sociologues ou politistes qui produisent les données qu'ils ou elles exploitent. On parle alors d'**analyse secondaire des données**. C'est le cas lorsqu'on travaille sur des données de l'Insee ou n'importe quelle base de donnée produite par une administration.

Quelques liens pour accéder aux données de la statistique publique française :

- le site de l'Adisp (Archives de données issues de la statistique publique) , qui rassemble les données de l'Insee et des directions statistiques ministérielles (santé, travail, culture, etc.)
- les données de l'Ined (Institut national de la recherche démographique)

### 1.1.4 Données d'enquête et données de gestion

Parmi l'ensemble des données accessibles produites par la statistique publique, on distingue en général deux grandes catégories [Desrosières, 2005] . D'un côté les bases de données produites via une **enquête par questionnaire** comme évoqué plus haut : elles sont réalisées à partir d'un échantillonnage au sein d'une population plus large (voir plus loin pour des définitions de ces termes), et comportent un grand nombre de variables, qui correspondent en général à des questions qui sont posées directement par des enquêteurs ou enquêtrices. De l'autre côté, certaines bases de données sont le **résultat du travail de gestion de certaines administrations** : par exemple, les employeurs effectuent chaque année ce qu'on appelle une "déclaration annuelle de données sociales", dans laquelle ils renseignent une série d'informations sur leurs différents salariés (parmi lesquelles leur salaire et leur profession). Ces "DADS" constituent un exemple de base de données administrative. Ils sont largement utilisés pour étudier les salaires. Ces bases de données sont intéressantes mais en général moins riches que les données d'enquête, car elles ne sont pas réalisées dans le but de produire de la connaissance. Je vous conseille d'éviter de choisir une telle base de données pour votre rendu du semestre, car il est souvent plus difficile d'en tirer des résultats intéressants à moins de savoir exactement ce qu'on cherche.

## 1.2 Le vocabulaire de la statistique

### 1.2.1 Bases de données

Il est temps d'expliquer plus précisément ce qu'on entend par "base de données". En voilà un premier exemple, issu du package R `titanic` :

Une **base de données** se présente sous la forme d'un tableau. Les lignes décrivent les **individus** : ici ce sont des passagers, mais gardez en tête que la nature des individus peut être à peu près n'importe quoi (ça peut être des ménages, des villes, des bactéries, n'importe quoi). Chaque colonne apporte des

Table 1.1: Extrait de la base de données des passagers du Titanic

PassengerId	Survived	Age	Name
1	0	22	Braund, Mr. Owen Harris
2	1	38	Cumings, Mrs. John Bradley (Florence Briggs Thayer)
3	1	26	Heikkinen, Miss. Laina
4	1	35	Futrelle, Mrs. Jacques Heath (Lily May Peel)
5	0	35	Allen, Mr. William Henry
6	0	NA	Moran, Mr. James

Table 1.2: Extrait de la base USArrests

	Murder	Assault	UrbanPop	Rape
Alabama	13.2	236	58	21.2
Alaska	10.0	263	48	44.5
Arizona	8.1	294	80	31.0
Arkansas	8.8	190	50	19.5
California	9.0	276	91	40.6
Colorado	7.9	204	78	38.7

éléments permettant de caractériser les individus (leur nom, leur âge, etc.). On appelle ces caractéristiques des **variables**.

### 1.2.2 Un autre exemple

Dans cet exemple, les individus sont des États des États-Unis. Les variables correspondent à des taux d'arrestation par la police pour meurtre, agression et viol pour 10000 habitants en 1973, ainsi que le pourcentage de la population urbaine.

### 1.2.3 Données “tidy”

Dans R, on qualifie certaines base de données de “tidy” [Wickham, 2014]. C’est la structure qu’on souhaite avoir en général. Ces bases de données ont un individu par ligne, une variable par colonne. Dans chaque case, on trouve la modalité d’une variable correspondant à l’individu décrit dans la ligne. L’exemple présenté sur la figure 1.2 en bas à droite n’est pas “tidy”, car il existe une variable (dont on ne connaît pas le nom) dont les modalités sont réparties dans deux colonnes différentes, qui représentent les années 1999 et 2000, c’est-à-dire les modalités d’une autre variable qui indique l’année. En remplaçant l’année et la variable observée chacune dans une colonne, on obtient un tableau ‘tidy’ (en bas à gauche).

country	year	cases	population
Afghanistan	1999	745	19557071
Afghanistan	2000	2666	20395360
Brazil	1999	37737	172006362
Brazil	2000	80488	174304898
China	1999	212258	1272315272
China	2000	213766	128028583

variables

country	year	cases	population
Afghanistan	1999	745	19557071
Afghanistan	2000	2666	20395360
Brazil	1999	37737	172006362
Brazil	2000	80488	174304898
China	1999	212258	1272315272
China	2000	213766	128028583

observations

country	year	cases
Afghanistan	1999	745
Afghanistan	2000	2666
Brazil	1999	37737
Brazil	2000	80488
China	1999	212258
China	2000	213766

country	1999	2000
Afghanistan	745	2666
Brazil	37737	80488
China	212258	213766

table4

Figure 1.2: Tidy data sets

Table 1.3: Extrait de la base de données beaver1

day	time	temp	activ
346	840	36.33	0
346	850	36.34	0
346	900	36.35	0
346	910	36.42	0
346	920	36.55	0
346	930	36.69	0

### 1.2.4 Séries temporelles

On parle parfois de **série temporelle** lorsque une base de donnée concerne un même individu statistique à différents instants. On parle dans ce cas là plutôt d'**observations** que d'individus. En voilà un exemple : la base de données **beaver1** accessible dans R présentent la température corporelle d'un castor en fonction du temps.

### 1.2.5 Autres types de bases de données

- Des données qui décrivent différents individus à un moment donné sont parfois qualifiées de **données en coupe** (ou *cross-sectional dataset*). Les données du titanic ou de USArrest en sont des exemples.
- Certaines bases de données décrivent un même groupe d'individus statistique de manière répétée dans le temps. On nomme ce genre de données des **données de panel** (exemple : enquête Emploi en continu).

## 1.3 Variables

### 1.3.1 Définition

Les variables sont les éléments qui permettent de décrire les individus présents dans la base de données. Lorsque les données sont issues d'un questionnaire, chaque question correspond en général à une variable.

Exemple :

- le sexe
- la catégorie socioprofessionnelle
- le niveau de diplôme
- le revenu

On appelle **modalités** les différentes valeurs que peuvent prendre une variable.

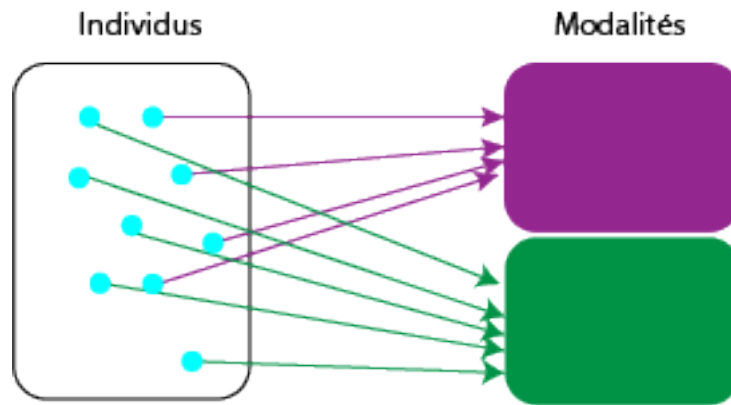


Figure 1.3: Les variables associent les individus à leur modalités

### 1.3.2 Variables qualitatives et variables quantitatives

On distingue les variables en fonction des opérations statistiques qu'on peut effectuer à partir de leurs modalités. Les deux grandes catégories de variables sont les **variables quantitatives**, dont les modalités sont des nombres, et les **variables qualitatives**, qui sont les autres.

### 1.3.3 Variables qualitatives

Parmi les variables qualitatives, on distingue encore :

- Les variables **qualitatives ordinales**, qui sont celles pour lesquelles on peut ordonner les modalités (exemple : niveau de diplôme)
- Les variables **qualitatives nominales** dont les modalités ne sont pas hiérarchisables (exemple : sexe, catégorie socioprofessionnelle)

### 1.3.4 Variables quantitatives

Les modalités des variables quantitatives (ou numériques) sont des nombres qui ont une signification (par exemple, le code postal n'est pas une variable quantitative). Parmi elles, on distingue :

- Les variables **continues**, qui peuvent prendre toutes les valeurs réelles dans un intervalle donné
- Les variables **discrètes**, qui ne peuvent prendre qu'un certain nombre de valeurs

**Pourquoi toutes ces catégories ?** À ces différents types de variables, on associe différentes méthodes statistiques. Il est donc important de comprendre et mémoriser ces définitions, car lorsque vous souhaitez étudier une variable, la première chose à faire sera d'identifier son type pour ensuite utiliser les méthodes statistiques appropriées.

## 1.4 Mesures de tendance centrale

Ce sont des manières de résumer l'information contenue dans une variable. En fonction du type de variable, il existe plusieurs indicateurs.

- La **moyenne** : lorsqu'on parle de moyenne, on fait généralement référence à la moyenne arithmétique d'un ensemble de valeurs numériques (par opposition à la moyenne géométrique ou harmonique). C'est une mesure très utilisée car elle fournit un premier résumé de la distribution statistique. Elle existe uniquement pour les **variables quantitatives**.

$$\bar{X} = \frac{1}{N} \sum_i^N x_i$$

- La **médiane** est la modalité d'une variable qui permet de séparer la population en deux parts égales. C'est-à-dire que 50% des individus auront une modalité supérieure ou égale à la médiane, et 50% une modalité inférieure ou égale à la médiane.
- Les **quantiles** sont une généralisation de la médiane : si vous voulez diviser votre population en groupe de 10%, vous pouvez utiliser les **déciles**. On appelle la médiane le **quantile d'ordre 2**, tandis que les déciles sont les **quantiles d'ordre 10**. Les quantiles les plus utilisés sont la médiane, les quartiles, les déciles et les centiles.
- Les quantiles n'existent que pour les variables dont les modalités peuvent être hiérarchisées : toutes les variables quantitatives et les variables qualitatives ordinales.
- Le **mode** indique la modalité la plus fréquente d'une variable. Par exemple, la plupart des passagers du Titanic sont décédés dans le naufrage, donc le mode de la variable "Survived" est 0. Le mode existe pour tous les types de variables.

## Chapter 2

# Statistique descriptive univariée

Le cours précédent était consacré à vous présenter différents types de variables. Celui de cette semaine présente les premiers éléments de **statistique descriptive univariée**, les outils permettant la description d'une unique variable. Ces outils dépendent de la nature de la variable étudiée.

### 2.1 Variables qualitatives

#### 2.1.1 Tris à plat

Pour décrire ce genre de variable, le principal traitement statistique est de compter le nombre d'individus correspondant à chaque modalité de la variable. C'est ce qu'on appelle un **tri à plat** (par opposition aux tris croisés qui font intervenir plusieurs variables). Un exemple issu des données du titanic (voir section 1.1).

	n	%
1	216	24.2
2	184	20.7
3	491	55.1
Total	891	100.0

À partir d'un tableau comportant une ligne par passager, on produit donc un tableau qui comporte seulement une ligne par classe de passagers. La colonne d'effectif montre le nombre de passagers par classe, tandis que la colonne de pourcentage indique le pourcentage de passagers des différentes classes parmi l'ensemble de passagers du Titanic. On peut lire le tableau de cette manière : parmi les 891 passagers du Titanic, 216 voyageaient en première classe. On cal-

cule le pourcentage de chaque catégorie en divisant l'effectif de chaque catégorie par l'effectif total, puis en multipliant par 100 ( $\frac{216}{891} * 100 = 24,2\%$ ).

Dans le cas où la variable est qualitative ordinale (c'est-à-dire qu'on peut ordonner ses modalités de manière hiérarchique, comme c'est le cas pour la variable de classe), on peut présenter dans ce tableau les **pourcentages cumulés**.

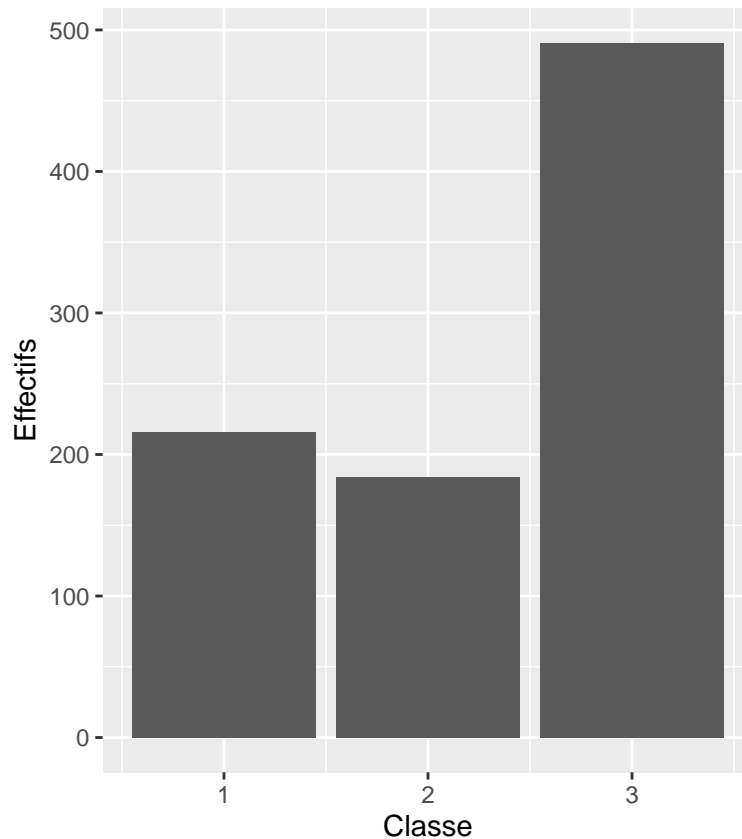
	n	%	%cum
1	216	24.2	24.2
2	184	20.7	44.9
3	491	55.1	100.0
Total	891	100.0	100.0

Ici, le chiffre 44,9% représente le pourcentage des passagers qui voyageaient **au moins en seconde classe**. Il s'agit simplement de la somme des pourcentages des passagers des première et seconde classe. La ligne suivante indique le pourcentage de passagers qui voyageaient au moins en troisième classe, ce chiffre est donc logiquement égal à 100%.

### 2.1.2 Diagrammes en barre

La représentation graphique associée à ce décompte est ce qu'on appelle généralement un **diagramme en barre** (vous pouvez aussi trouver "diagramme en bâtons ou diagramme en tuyaux d'orgue" qui désignent la même chose). Sur ce diagramme, on trace des barres verticales dont les hauteurs sont proportionnelles aux effectifs du tri à plat. Seule la hauteur des barres à une signification, la largeur est totalement arbitraire.





## 2.2 Variables quantitatives

Si la statistique univariée est très simple pour une variable qualitative, elle peut faire l'objet d'analyses plus approfondies lorsqu'on dispose de variables quantitatives.

### 2.2.1 Mesures de dispersion

La semaine dernière, je vous ai présenté quelques **mesures de tendance centrale**. Elles donnent des renseignements importants pour décrire une variable, mais n'en résument qu'une dimension. Deux séries statistiques peuvent avoir la même moyenne tout en étant très différentes.

Comparez par exemple ces deux séries de chiffres, qui représentent des profits (fictifs) en dollars de deux agriculteurs de deux régions A et B :

- A: 14, 16, 18, 20, et 22
- B: 2, 8, 18, 29, et 33

La somme de ces deux série est la même, 90 dollars, mais il apparaît rapidement que l'une des séries est beaucoup plus **dispersée** que l'autre, c'est-à-dire que les écarts par rapport à la moyenne sont en général beaucoup plus grands (la série B). Notre vision des risques et des profits liés à l'agriculture est informée par cette différence, et nous devrions en inclure des indices dans toute description statistique de cette variable.

Pour faire cela, nous avons besoin de mesures permettant de décrire la dispersion des modalités de la variable autour de sa moyenne.

### 2.2.1.1 L'étendue

C'est la différence entre la plus grande et la plus petite valeur de la série :

$$R = X_{max} - X_{min}$$

C'est une mesure de dispersion assez basique. Son défaut est assez évident : elle dépend uniquement des valeurs extrêmes, et aucunement du reste de la distribution.

### 2.2.1.2 L'écart interquartile

Pour prendre en compte plus que les deux valeurs extrêmes, on peut calculer la différence entre deux quantiles, des quartiles par exemple (voir définition dans le premier cours)

$$Q_d = Q_3 - Q_1$$

C'est une mesure un peu meilleure que l'étendue, parce que le maximum et le minimum sont des valeurs qui donnent généralement peu d'information sur la distribution en général. Cet écart représente l'étendue de la moitié de la distribution, moitié obtenue après avoir enlevé les 25% des valeurs les plus faibles et 25% des valeurs les plus hautes. L'écart interquartile est moins sensible aux valeurs extrêmes que l'étendue (puisque'on les a supprimées), mais résume tout de même l'ensemble de données sans prendre en compte la variabilité des données entre le premier et le 3ème quartile. Les mesures suivantes n'ont pas ces défauts.

### 2.2.1.3 La variance

Ce qu'on voudrait, c'est l'équivalent de la moyenne, mais pour mesurer la dispersion. On pourrait donc se dire qu'il suffirait de faire la **moyenne des écarts à la moyenne** de cette manière <sup>1</sup>:

---

<sup>1</sup>Les deux côtés de l'équation représentent la même chose, il s'agit juste d'une différence de notation. À gauche, on utilise ... pour indiquer la série de termes supplémentaires qu'on va inclure dans l'addition mais qu'on n'écrira pas. À droite, l'opérateur  $\sum_{i=1}^n$  est la somme pour i allant de 1 à n de l'expression qui est à droite. Cela signifie qu'il faut remplacer i par

$$\frac{1}{N}((X_1 - X_m) + (X_2 - X_m) + \dots + (X_n - X_m)) = \frac{1}{n} \sum_{i=1}^n (X_i - X_m)$$

Le problème, en faisant ça, c'est que compte tenu de la définition de la moyenne, les écarts à la moyenne vont se compenser terme à terme. On peut le voir facilement si l'on sépare la somme en deux :

$$\frac{1}{n} \sum_{i=1}^n (X_i - X_m) = \frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{n} \sum_{i=1}^n X_m$$

Du côté droit de l'équation, le terme de gauche ( $\frac{1}{n} \sum_{i=1}^n X_i$ ) est la définition de la moyenne (c'est la somme des termes  $X_1 + X_2 + \dots + X_n$  divisé par l'effectif total  $n$ ), tandis qu'à droite ( $\frac{1}{n} \sum_{i=1}^n X_m$ ) on ajoute  $n$  fois la moyenne  $X_m$  puis on la divise par  $n$ , donc on obtient encore la moyenne. Au final, cette somme est toujours égale à 0.

Pour éviter ce problème, on définit la variance comme la somme des écarts à la moyenne **au carré**.

$$Var = \frac{1}{n} \sum_{i=1}^n (X_i - X_m)^2$$

Cette définition a l'avantage de donner un résultat non nul, excepté dans le cas où la variable est une constante (qui est alors toujours égale à sa moyenne). Surtout, les écarts à la moyenne s'ajoutent, qu'ils soient générés par des valeurs supérieures ou inférieures à la moyenne.

#### 2.2.1.4 L'écart type

La variance a beaucoup de propriétés intéressantes et on l'utilise très largement en statistique. Malgré tout, elle a un dernier inconvénient, c'est de s'exprimer comme un carré de l'unité dans laquelle est mesurée la variable. Par exemple, si l'on mesure la taille des étudiant-es de la classe puis qu'on calcule la variance, on aura un résultat en centimètres ou en mètres au carré. Dans notre exemple, on obtient une mesure en "dollars au carré", dont le sens est difficile à interpréter.

On résout ce problème de manière simple en calculant la racine carrée de la variance. Cette opération nous permet d'obtenir notre dernière mesure de dispersion, **l'écart-type**.

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - X_m)^2}$$

---

1, puis par 2, 3, etc. jusqu'à  $n$  et faire la somme de tous les éléments ainsi obtenu. Ce qui revient exactement à ce qui est écrit de manière plus longue de l'autre côté de l'équation.

L'écart-type est la mesure de dispersion la plus utile et la plus fréquente. Elle est meilleure que la variance car elle se mesure dans la même unité que la variable en question. Par exemple, on peut dire que dans la région A, la moyenne des revenus agricoles est de 18 dollars, avec un écart-type de 2,8 dollars.

## 2.2.2 Représentations graphiques

Il existe plusieurs manières de représenter graphiquement la distribution d'une variable quantitative.

### 2.2.2.1 Histogramme

Les histogrammes sont l'équivalent des diagrammes en barres pour les variables quantitatives. Chaque barre (ou rectangle) qui compose l'histogramme a une aire qui est proportionnelle au nombre d'observation dont les valeurs sont dans l'intervalle sur lequel s'étend le rectangle.

Comme il s'agit de variables quantitatives, on peut choisir le nombre de rectangles comme on le souhaite (contrairement aux variables qualitative dont les modalités sont définies une fois pour toutes). Ici, on représente la même variable avec moins de rectangles.

Ici avec un plus grand nombre de rectangle (largeur = 1 an)

### 2.2.2.2 Densité

On représente parfois les variables quantitatives par une courbe que l'on appelle une 'densité'. C'est la courbe qu'on pourrait obtenir si on avait un très grand nombre de passagers et qu'on représentait l'histogramme avec des rectangles très fins.

On peut également produire une estimation de cette courbe à partir d'une transformation effectuée sur les histogrammes. Si l'on trace une ligne qui passe au milieu de chacun des segments supérieurs des rectangles qui composent l'histogramme, on obtient alors un graph qu'on nomme un **polygone de fréquence**.

La forme de ce polygone peut être être "lissée" à l'aide de techniques mathématiques, pour donner la courbe de densité recherchée. Elle donne une idée de la forme du polygone de fréquence si l'on avait un très grand nombre d'individu dans notre échantillon.

Ces représentations graphiques sont un bon moyen de visualiser la **forme** de la distribution d'une variable quantitative continue, et spécifiquement de la manière dont les données sont réparties autour de leur valeur "centrale".

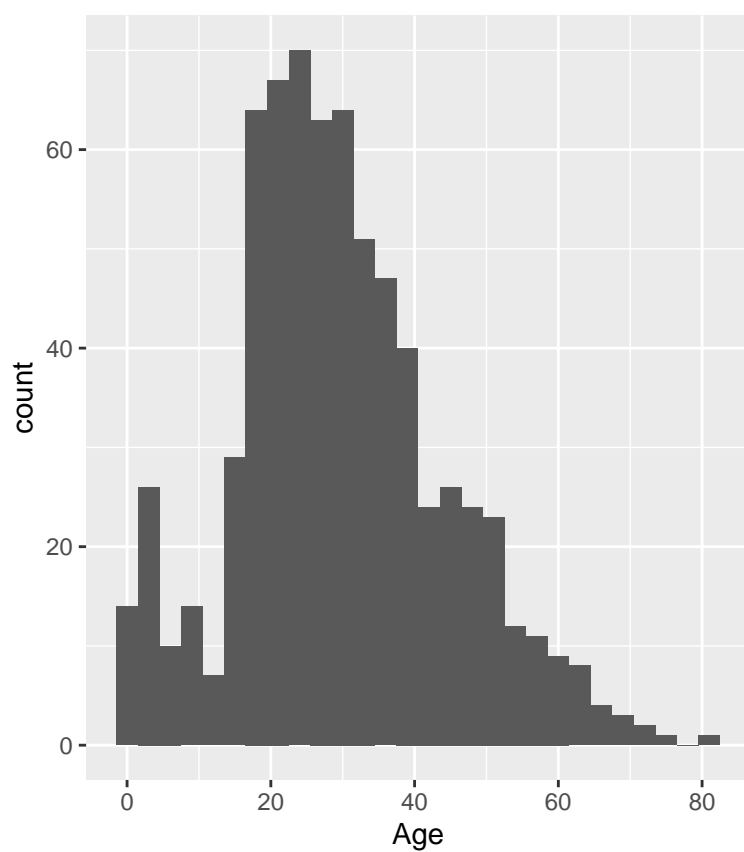


Figure 2.1: Distribution de l'âge des passagers du Titanic. Chaque rectangle a une largeur de 3 ans

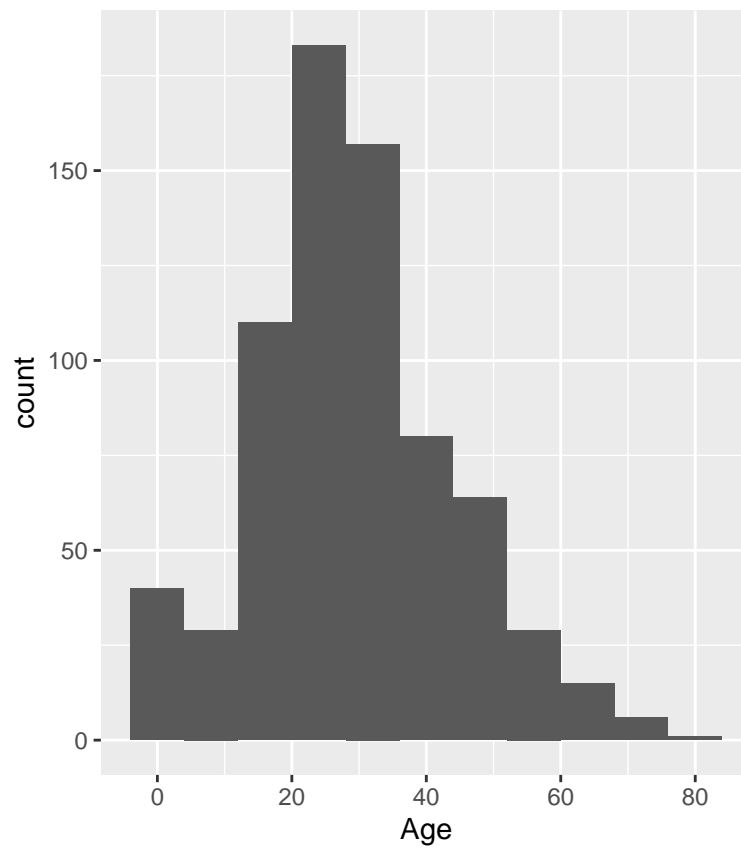


Figure 2.2: Même figure avec une largeur de 8 ans

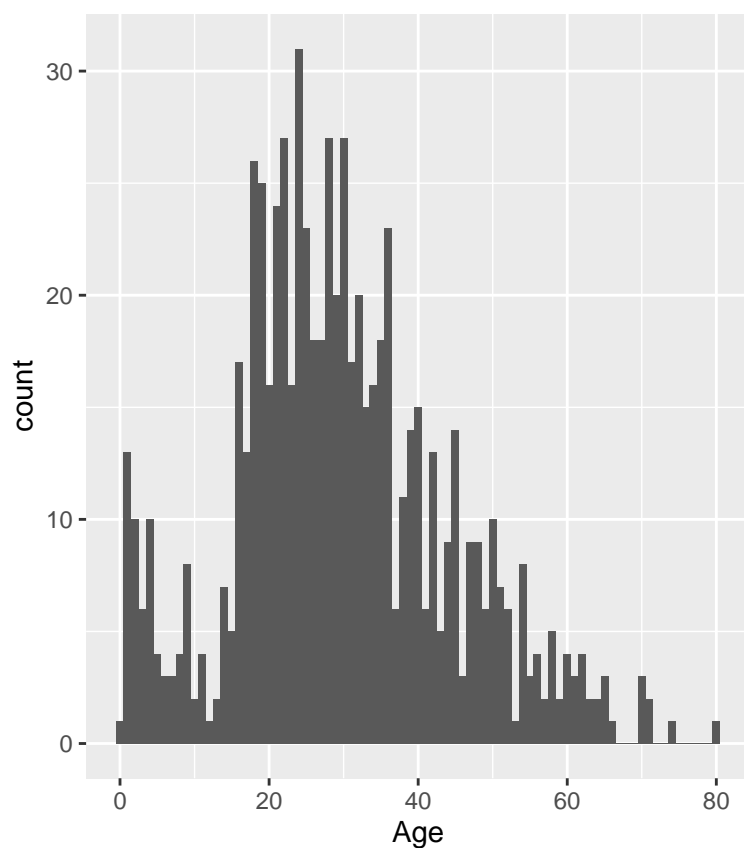
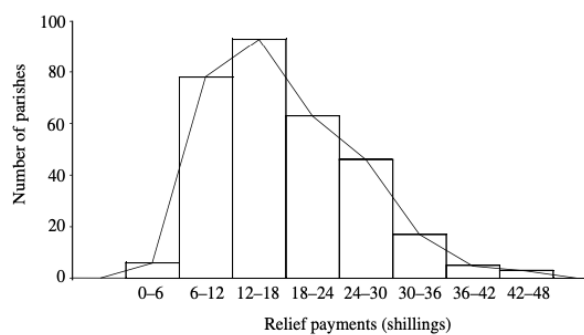


Figure 2.3: Même figure avec une largeur d'un an

Figure 2.4: *Per capita* relief payments in 311 parishes in 1831 (Fenstein & Thomas, p. 41)

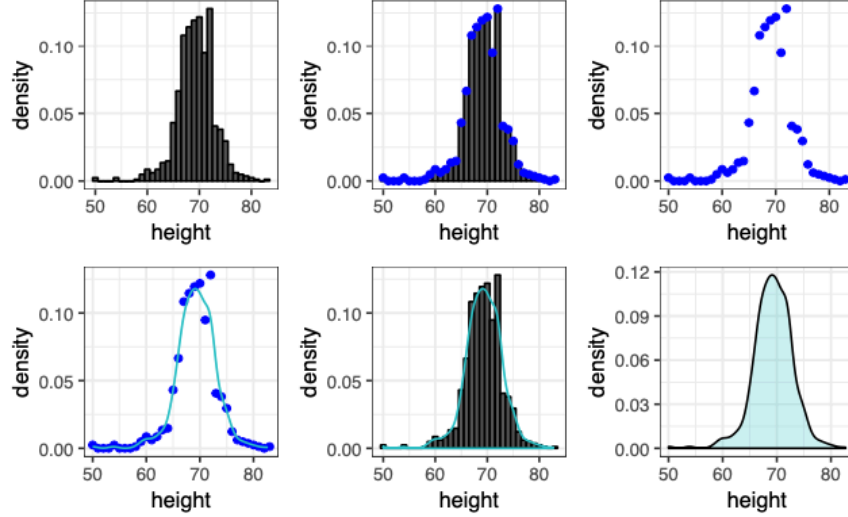


Figure 2.5: D'un histogramme à une densité de probabilité

### 2.2.2.3 Asymétrie (*skewness*)

Une manière de caractériser les distribution est leur symétrie par rapport à la moyenne. Les valeurs peuvent en effet être réparties de manière symétrique de part et d'autre de la moyenne, ou bien de manière asymétrique (*skewed*).

Différentes mesures permettent de quantifier l'asymétrie d'une distribution.

- Elles doivent être indépendantes de l'unité de mesure
- Et elle doivent être nulles lorsque la distribution est symétrique.

Un exemple de coefficient d'asymétrie est le suivant, mais il en existe d'autres :

$$Skewness = \frac{3 * (Mean - Median)}{\sigma}$$

Un exemple de variable dont la distribution est asymétrique est la distribution des revenus dans la population française. L'histogramme suivant représente la distribution du niveau de vie (c'est le revenu des ménages divisé par le nombre d'unités de consommation). Le niveau de vie médian (compris dans la portion verte du graphique) est inférieur au niveau de vie moyen, car les ménages au niveaux de vie très élevés (en bleu) tirent la moyenne vers le haut, mais n'ont pas d'effet sur la médiane.



**Figure 2.4**  
Symmetrical and  
skewed frequency  
curves

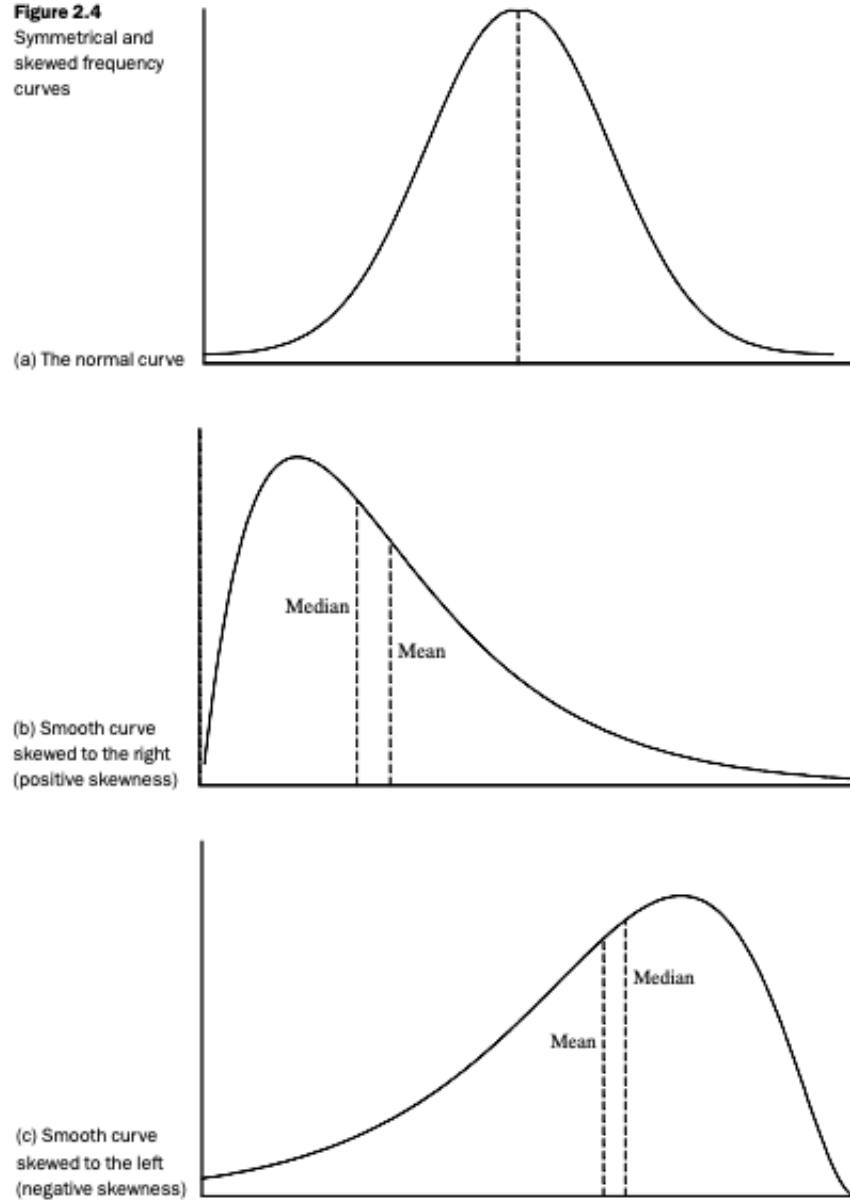


Figure 2.6: Distributions symétriques et asymétriques [Feinstein and Thomas, 2002, p.54]

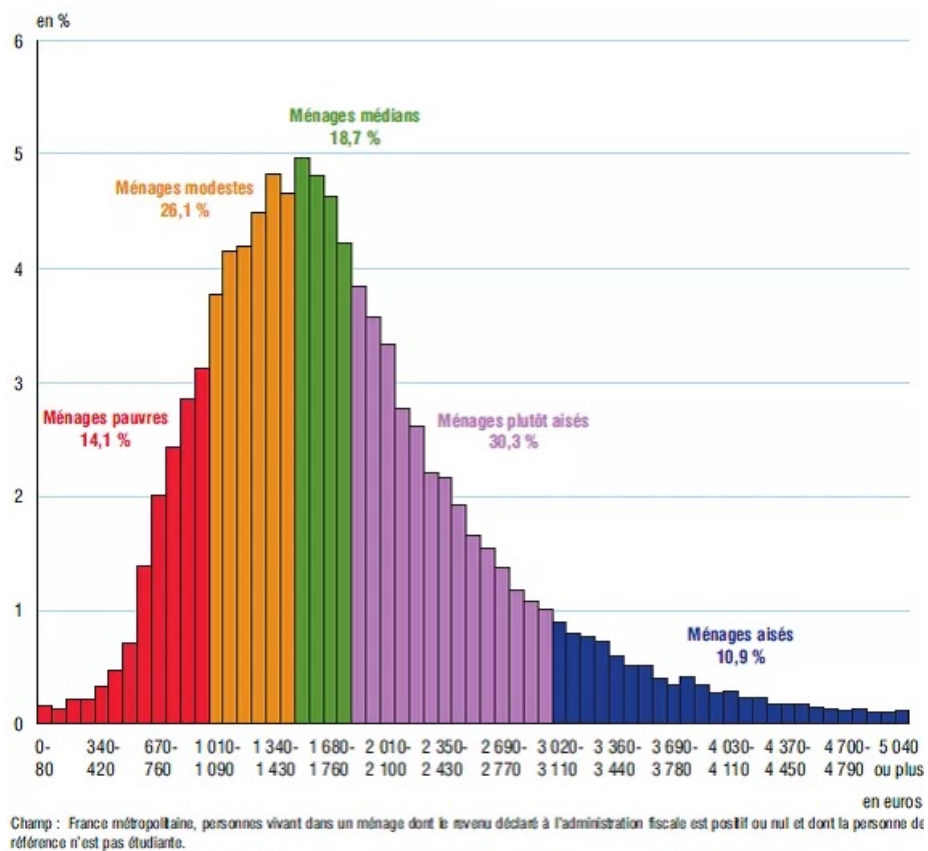


Figure 2.7: Distribution des niveaux de vie mensuels en 2014 en France (Source : Insee, Portrait social 2014)

## 2.3 La loi normale : une distribution importante

La loi normale est une distribution théorique, définie à partir de son expression mathématique. Mais bien que théorique, c'est une distribution très importante, car elle est souvent utilisée comme approximation de distributions réelles. Je vous la présente ici rapidement, on la retrouvera dans des prochaines séances.

Pour définir une loi normale, il faut connaître deux constantes : sa moyenne  $X_m$  et l'écart type  $\sigma$ . L'équation donne la valeur de  $Y$  (la hauteur de la courbe, qui apparaît sur l'axe des ordonnées) pour toute valeur de  $X$  (mesuré sur l'axe des abscisses) :

$$Y(X) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(X - X_m)^2}{2 * \sigma^2}\right)$$

La fonction  $\exp()$  qui apparaît dans la formule est la **fonction exponentielle**. Si vous ne connaissez pas cette fonction, sachez qu'elle est définie par le fait qu'il s'agit de l'unique fonction  $f(x)$  qui est toujours égale à sa dérivée (la fonction dérivée est celle qui mesure la pente de la courbe en chaque point, on la note  $f'(x)$  : elle est positive lorsque  $f$  est croissante, et négative lorsqu'elle est décroissante) et qui est égale à 1 lorsque  $x = 0$ . Comme elle est toujours égale à sa dérivée, plus  $x$  est élevé, plus la fonction exponentielle doit avoir une dérivée élevée, donc plus elle doit croître rapidement.

```
curve(exp(x), from=-5, to=5, , xlab="x", ylab="y")
```

Dans la distribution de la loi normale, la fonction exponentielle contient une expression qui est toujours inférieure ou égale à zéro. Son maximum est donc atteint lorsque  $X$  est égal à sa moyenne  $X_m$ , auquel cas  $Y(X_m) = \frac{1}{\sigma\sqrt{2\pi}}$ . Plus  $X$  va s'éloigner de sa moyenne, plus  $Y(X)$  sera faible, on dit que la distribution tend vers 0 lorsque  $X$  tend vers "moins l'infini" ou "plus l'infini". Le graphe de la loi normale ressemble donc à un dos d'âne, ce qui explique qu'on l'appelle aussi "la courbe en cloche".

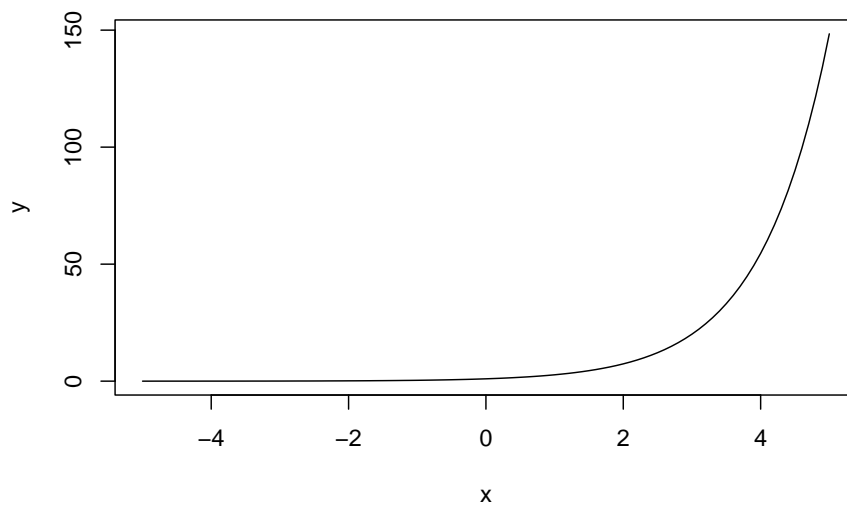
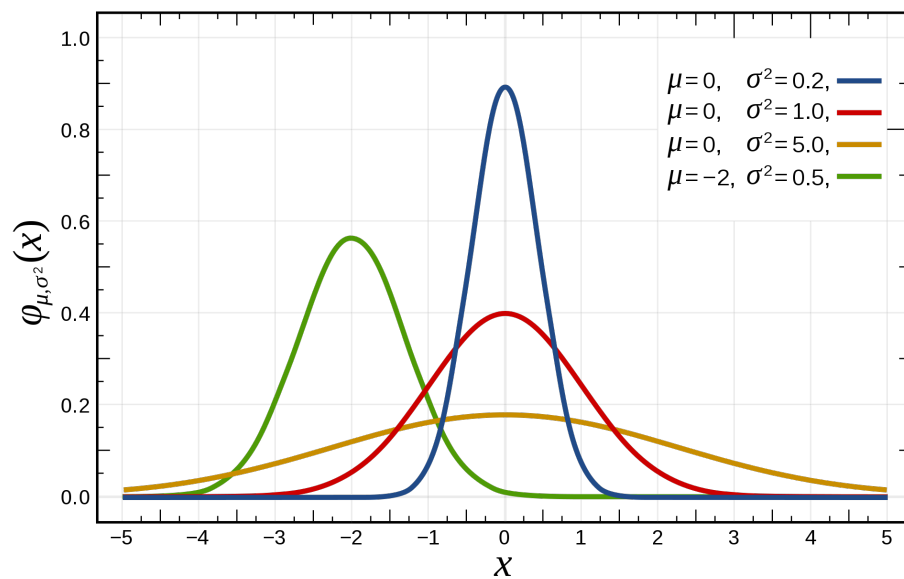


Figure 2.8: Graphe de la fonction exponentielle entre -5 et 5



Ce dernier graphe permet de constater que, si les lois normales ont toutes la même allure, leur forme dépend de la moyenne et de l'écart-type de la

distribution. Comme déjà évoqué, la moyenne indique le maximum de la courbe. L'écart-type détermine lui la "largeur" de la bosse, c'est-à-dire à quel point les données s'étalent autour de la moyenne.

Une propriété importante de la loi normale est que, quelque soit sa moyenne et son écart-type, il y a toujours une même proportion d'observations qui seront distribués à une certaine distance de la moyenne (que l'on peut mesurer en calculant l'aire sous la courbe), mesurée en nombre d'écart-type.

Par exemple :

- **90% des observations** sont situés à moins de **1,645 écarts-type** autour de la moyenne, laissant 5% de chaque côté.
- **95% des observations** sont situés à moins de **1,96 écarts-type** autour de la moyenne, laissant 2,5% de chaque côté.
- **99% des observations** sont situés à moins de **2,58 écarts-type** autour de la moyenne, laissant 0,5% de chaque côté.

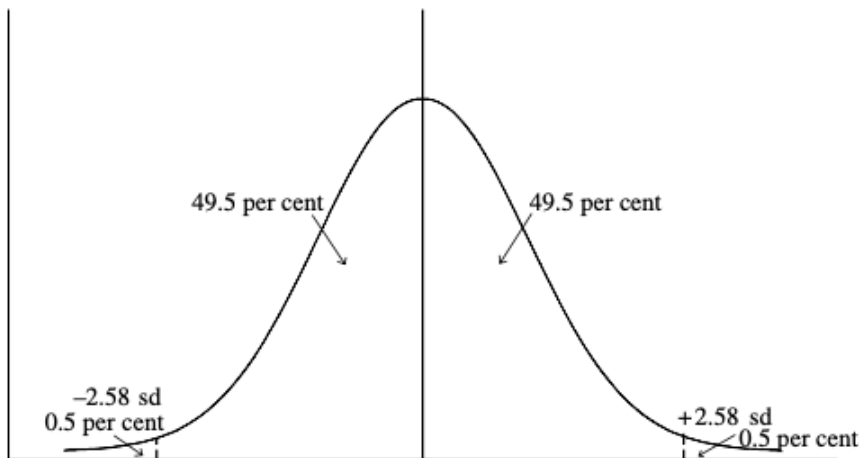


Figure 2.9: Aire sous la courbe



# Références





# Bibliography

Pierre Bourdieu. *La distinction: critique sociale du jugement*. Les Editions de minuit, Paris, France, 1979. ISSN: 0768-049X.

Fanny Bugeja-Bloch and Marie-Paule Couto. *Les méthodes quantitatives*. Que sais-je ? PUF, 2021. OCLC: 1285669386.

Alain Desrosières. Décrire l'État ou explorer la société : les deux sources de la statistique publique. *Geneses*, no 58(1):4–27, 2005. URL <https://www.cairn.info/journal-geneses-2005-1-page-4.htm>. Bibliographie\_available: 0 Cairn-domain: www.cairn.info Cite Par\_available: 1 Publisher: Belin.

Charles H. Feinstein and Mark Thomas. *Making History Count: A Primer in Quantitative Methods for Historians*. Cambridge University Press, Cambridge ; New York, 08 2002.

Claire Lemerrier and Claire Zalc. *Méthodes quantitatives pour l'historien*. Number 507 in Repères. la Découverte, Paris, 2008.

Jean Peneff. *L'hôpital en urgence: étude par observation participante*. Métailié : Diffusion, Seuil, Paris, 1992.

Hadley Wickham. Tidy data. *Journal of Statistical Software*, 59:1–23, 09 2014. doi: 10.18637/jss.v059.i10. URL <https://doi.org/10.18637/jss.v059.i10>.