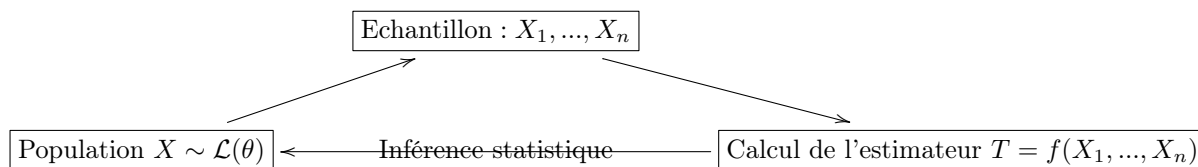


# Chapitre 1

## Estimateurs et intervalles de confiance

### 1.1 Principe général



1. Dans une population donnée, on s'intéresse à un caractère représenté par une v.a.  $X$  qui suit une loi de probabilité  $\mathcal{L}(\theta)$  qui dépend d'un paramètre  $\theta$  que l'on veut déterminer (par exemple  $\theta = \mu =$  moyenne de  $X$ ).
2. Pour cela, on prend un échantillon de la population, donc une réalisation de  $X_1, \dots, X_n$  qui sont des v.a. indépendantes et identiques à  $X$ .
3. On utilise un estimateur  $T = f(X_1, \dots, X_n)$  pour obtenir une valeur approchée de  $\theta$  (par exemple, pour  $\theta = \mu$ , on utilise  $T = \overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ ).
4. On détermine ainsi, par inférence statistique, la caractéristique de la population totale.

On distingue ensuite, et selon les besoins, entre deux procédés :

- **L'estimation ponctuelle**, qui consiste à prendre la valeur donnée par l'estimateur comme valeur approchée de  $\theta$  et à étudier l'erreur commise.
- **L'estimation par intervalle de confiance (IDC)**, qui consiste à trouver un intervalle dans lequel se trouve la valeur exacte de  $\theta$  avec un niveau de confiance donné.

### 1.2 Estimation ponctuelle

#### 1.2.1 Qualités d'un estimateur

Une des premières qualités demandées à un estimateur est d'être *sans biais*. En effet, l'erreur commise, en prenant la valeur estimée à la place de la valeur exacte est :  $T - \theta = [T - E(T)] + [E(T) - \theta]$ . Le premier terme correspond à la fluctuation naturelle d'une v.a. autour de sa moyenne, le deuxième terme est systématique et peut et doit être réduit.

##### 1. Biais :

On appelle *biais* d'un estimateur  $T$  du paramètre  $\theta$ , l'expression :

$$b(T) = E(T) - \theta.$$

Un estimateur est dit *sans biais* si  $b(T) = 0$ .

## 2. Convergence :

Le *risque quadratique* d'un estimateur ou *erreur quadratique moyenne* est donné par  $R_\theta(T) = E((T - \theta)^2) = V(T) + b(T)^2$ .

On dit d'un estimateur qu'il est *convergent* lorsque  $\lim_{n \rightarrow \infty} R_\theta(T) = 0$ .

Ceci est notamment le cas si l'estimateur est sans biais et que sa variance tend vers 0.

## 3. Efficacité :

Si  $T_1$  et  $T_2$  sont deux estimateurs sans biais du même paramètre  $\theta$ , on dit que  $T_1$  est *plus efficace* que  $T_2$  si  $V(T_1) < V(T_2)$ .

L'estimateur est dit *efficace* si sa variance est la plus petite parmi les estimateurs sans biais.

## 1.2.2 Estimateurs usuels

### 1. Proportion :

Lorsque  $X$  suit une loi de Bernouilli,  $X \sim \mathcal{B}(p)$ , et qu'on veut estimer le paramètre  $\theta = p$ , on utilise la fréquence empirique

$$F_n = \overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

— Cet estimateur est **sans biais**, car  $E(F_n) = E(X_1) = p = \theta$ .

— Il est **convergent** car le risque quadratique

$$R_\theta(F_n) = V(F_n) = \frac{V(X_1)}{n} = \frac{p(1-p)}{n} \text{ tend bien vers } 0.$$

— Sa loi asymptotique est donnée par le théorème de la limite centrale (TCL) :  $F_n \xrightarrow{\mathcal{L}} \mathcal{N}(p, \frac{p(1-p)}{n})$ .

### 2. Moyenne :

Si  $X$  suit une loi de moyenne  $\mu$  et de variance  $\sigma^2$ , on utilise la moyenne empirique  $\overline{X}_n$  pour estimer  $\theta = \mu$  avec

$$\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

où  $X_1, \dots, X_n$  est un échantillon correspondant.

— C'est un estimateur **sans biais**.  $E(\overline{X}_n) = \mu = \theta$ .

— Il est **convergent** car le risque quadratique  $R_\theta(\overline{X}_n) = V(\overline{X}_n) = \frac{\sigma^2}{n}$  tend bien vers 0.

— Le TCL donne :  $\overline{X}_n \xrightarrow{\mathcal{L}} \mathcal{N}(\mu, \frac{\sigma^2}{n})$ .

### 3. Variance :

Pour une v.a.  $X$  qui suit une loi de moyenne  $\mu$  et de variance  $\sigma^2$ , lorsqu'on veut estimer la variance ( $\theta = \sigma^2$ ), on utilise la variance empirique :

$$S_1^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \overline{X}_n)^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - (\overline{X}_n)^2$$

—  $E(S_1^2) = \frac{n-1}{n} \sigma^2 \neq \theta$  c'est un estimateur **biaisé**.

C'est pour cette raison qu'on introduit la variance empirique corrigée :

$$S^2 = \frac{n}{n-1} S_1^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \overline{X}_n)^2$$

qui est un estimateur sans biais de la variance.

— La variance empirique converge bien vers  $\sigma^2$  en probabilité et presque sûrement, mais la convergence en moyenne quadratique n'a lieu que si  $X$  et donc  $X_1, \dots, X_n$  suivent une loi normale.

## 1.3 Intervalles de confiance

### 1.3.1 Principe général

Le but est d'obtenir un intervalle  $[a, b]$  dans lequel se trouve le paramètre avec un niveau de risque  $\alpha$  donné.

$$P(a \leq \theta \leq b) = 1 - \alpha$$

$\alpha$  étant le niveau de risque pris, en général  $\alpha = 10\%$ ,  $5\%$  ou  $1\%$ .

$a$  et  $b$  sont obtenus en partageant le risque en  $\alpha = \alpha_1 + \alpha_2$  avec  $\alpha_1 = P(X < a)$  et  $\alpha_2 = P(X > b)$  dans le cas d'un IDC bilatéral.

Cet intervalle bilatéral peut être symétrique ( $\alpha_1 = \alpha_2$ ) ou dissymétrique.

### 1.3.2 IDC pour une proportion

On a vu que lorsque  $X \sim \mathcal{B}(p)$ , on utilisait la fréquence empirique  $F_n$  comme estimateur de  $\theta = p$ . Pour  $n$  assez grand, le TCL nous permettait d'affirmer que  $F_n \sim \mathcal{N}(\mu, \sigma_1^2)$  avec  $\mu = p$  et  $\sigma_1^2 = \frac{p(1-p)}{n} \simeq \frac{f_n(1-f_n)}{n}$

On en déduit que  $Z = \frac{F_n - p}{\sigma_1}$  suit la loi normale centrée réduite  $Z \sim \mathcal{N}(0, 1)$ .

Pour obtenir un IDC symétrique au niveau de risque  $\alpha$ , on utilise la valeur  $t$  lue dans la table de la fonction de répartition de la loi normale.

$t$  vérifie  $F_Z(t) = 1 - \frac{\alpha}{2}$  et aussi  $P(-t \leq Z \leq t) = 1 - \alpha$ .

Si  $f_n$  est la valeur estimée donnée par l'échantillon, on obtient :

$$1 - \alpha = P(-t \leq Z \leq t) = P\left(-t \leq \frac{f_n - p}{\sigma_1} \leq t\right) = P(f_n - t\sigma_1 \leq p \leq f_n + t\sigma_1)$$

L'intervalle de confiance recherché est donc :

$$I = [f_n - t\sigma_1 ; f_n + t\sigma_1]$$

### 1.3.3 IDC pour une moyenne

L'estimateur utilisé est évidemment  $\overline{X}_n$ . Sa loi asymptotique dépend du fait que la variance soit connue ou non.

Il y a, en gros, deux cas selon que l'on connaît ou non la variance de  $X$ .

#### 1. 1<sup>er</sup> cas : $\sigma$ connue et $X$ suit la loi normale ou $n$ est assez grand ( $n > 30$ )

Dans ce cas on peut considérer que :  $\overline{X}_n \sim \mathcal{N}(\mu, \frac{\sigma^2}{n})$  ou encore que

$$\frac{\overline{X}_n - \mu}{\sigma} \sqrt{n} \sim \mathcal{N}(0, 1)$$

Si  $\bar{x}_n$  est la valeur estimée donnée par l'échantillon, et pour un niveau de risque  $\alpha$ , on obtient comme intervalle de confiance :

$$I = \left[ \bar{x}_n - t \frac{\sigma}{\sqrt{n}} ; \bar{x}_n + t \frac{\sigma}{\sqrt{n}} \right]$$

avec  $t$  qui vérifie  $F_Z(t) = 1 - \frac{\alpha}{2}$  et aussi  $P(-t \leq Z \leq t) = 1 - \alpha$ .

#### 2. 2<sup>me</sup> cas : $X$ suit une loi normale, $\sigma$ inconnue :

Dans ce cas, on remplace  $\sigma^2$  par son estimateur corrigé  $S^2$  et on utilise :

$$\frac{\overline{X}_n - \mu}{S} \sqrt{n} \sim T_{n-1}$$

loi de Student à  $(n - 1)$  degrés de liberté qui est symétrique tout comme la loi normale.  
 Pour un niveau de risque  $\alpha$  donné, à l'aide des valeurs estimées  $\bar{x}_n, s$  et de la valeur  $t$  équivalent de  $t$  pour la loi de Student, on obtient l'IDC suivant :

$$I = \left[ \bar{x}_n - t \frac{s}{\sqrt{n}} ; \bar{x}_n + t \frac{s}{\sqrt{n}} \right]$$

### 3. Autres cas

Dans tous les autres cas, on ne peut rien dire.

#### 1.3.4 IDC pour une variance

On ne peut rien dire des lois suivies par les deux estimateurs de la variance, sauf dans le cas gaussien. Lorsque  $X$  et donc  $X_1, X_2, \dots, X_n$  suivent une loi  $\mathcal{N}(\mu, \sigma^2)$ , on distingue deux cas :

##### 1. 1<sup>er</sup> cas : moyenne $\mu$ connue :

L'estimateur utilisé pour la variance est alors :  $V_1 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$  et :

$$n \frac{V_1}{\sigma^2} \sim \chi_n^2$$

loi du khi-2 à  $n$  degrés de liberté.

Pour un risque de niveau  $\alpha = \alpha_1 + \alpha_2$ , l'IDC est obtenu à l'aide des valeurs :  $v_1$  valeur estimée de  $V_1$  et des valeurs  $k_1$  et  $k_2$  relevées sur la table du  $\chi_n^2$  grâce aux relations :  $\alpha_1 = P(\chi_n^2 < k_1)$  et  $\alpha_2 = P(\chi_n^2 > k_2)$

$$1 - \alpha = P(k_1 \leq \chi_n^2 \leq k_2) = P\left(k_1 \leq n \frac{v_1}{\sigma^2} \leq k_2\right) = P\left(n \frac{v_1}{k_2} \leq \sigma^2 \leq n \frac{v_1}{k_1}\right)$$

L'intervalle de confiance pour  $\sigma^2$  est donc :

$$\sigma^2 \in \left[ n \frac{v_1}{k_2} ; n \frac{v_1}{k_1} \right]$$

Un passage à la racine carrée permet alors d'obtenir l'IDC pour  $\sigma$ .

##### 2. 2<sup>eme</sup> cas : moyenne $\mu$ inconnue :

L'estimateur utilisé pour la variance est alors :  $S_1^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  ou  $S^2 = \frac{n}{n-1} S_1^2$ , on alors :

$$n \frac{S_1^2}{\sigma^2} = (n-1) \frac{S^2}{\sigma^2} \sim \chi_{n-1}^2$$

loi du khi-2 à  $(n - 1)$  degrés de liberté.

Pour un risque de niveau  $\alpha = \alpha_1 + \alpha_2$ , l'IDC est obtenu à l'aide des valeurs :  $s$  valeur estimée de  $S$  et des valeurs  $k_1$  et  $k_2$  relevées sur la table du  $\chi_{n-1}^2$  grâce aux relations :  $\alpha_1 = P(\chi_{n-1}^2 < k_1)$  et  $\alpha_2 = P(\chi_{n-1}^2 > k_2)$

On obtient alors comme intervalle de confiance pour  $\sigma^2$  :

$$\sigma^2 \in \left[ (n-1) \frac{s^2}{k_2} ; (n-1) \frac{s^2}{k_1} \right]$$

Un passage à la racine carrée permet alors d'obtenir l'IDC pour  $\sigma$ .

**Le cas d'un intervalle unilatéral se traite de la même manière.**