



# 15

## STATISTIQUES

### ESTIMATIONS PONCTUELLE ET PAR INTERVALLE DE CONFIANCE

---

#### INTRODUCTION...

Dans le chapitre 8, nous avons fourni des outils d'analyse de données numériques : moyenne, écart-type, médiane, quartiles, tableaux, graphiques... On se plaçait alors avec un certain nombre de données, que l'on analysait dans le but d'en extraire les informations essentielles : on parle de **statistique descriptive**.

Dans le chapitre qui suit, nous allons tenter, à partir de données obtenues sur un échantillon, d'en déduire des informations sur la population globale. On parle alors de **statistique inférentielle**. La statistique inférentielle est la statistique des sondages : on interroge un échantillon réduit de la population, et on tente d'estimer alors la tendance générale sur population entière. L'enjeu ici est de préciser la qualité des estimations proposées.

## POUR BIEN DÉMARRER...

### 1. Rappeler l'inégalité de Bienaymé-Tchebychev.

Si  $X$  admet une variance, alors :

$$\forall a \in \mathbb{R}_*^+, \mathbb{P}(|X - \mathbb{E}(X)| \geq a) \leq \frac{\mathbb{V}(X)}{a^2}$$

### 2. Rappeler la loi faible des grands nombres.

Soit  $(X_k)_{k \in \mathbb{N}^*}$  une suite de variables aléatoires sur  $(\Omega, \mathcal{A}, \mathbb{P})$ . Pour tout  $n \in \mathbb{N}^*$ , on note  $\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k$ .

Si  $(X_k)_{k \in \mathbb{N}^*}$  est une suite de variables aléatoires **indépendantes**, admettant toutes la **même espérance  $m$**  et la **même variance  $\sigma^2$**  (c'est le cas si elles ont toutes la même loi) alors :

$$\forall \varepsilon > 0, \lim_{n \rightarrow +\infty} \mathbb{P}(|\bar{X}_n - m| \geq \varepsilon) = 0$$

### 3. Rappeler le théorème central limite.

Soit  $(X_k)_{k \in \mathbb{N}^*}$  une suite de variables aléatoires **indépendantes, de même loi** (espérance  $m$ , variance  $\sigma^2$  non nulle).

On pose, pour tout  $n \in \mathbb{N}^*$  :

$$\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k ; \bar{X}_n^* = \frac{\bar{X}_n - \mathbb{E}(\bar{X}_n)}{\sigma(\bar{X}_n)} = \sqrt{n} \frac{\bar{X}_n - m}{\sigma}$$

Dans ce cas, la suite  $(\bar{X}_n^*)_{n \in \mathbb{N}^*}$  converge en loi vers une variable aléatoire suivant la loi  $\mathcal{N}(0; 1)$ .

Autrement dit :

$$\forall x \in \mathbb{R}, \lim_{n \rightarrow +\infty} \mathbb{P}(\bar{X}_n^* \leq x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt$$

Ou encore :

$$\forall (a, b) \in \mathbb{R}^2, \left( a < b \implies \lim_{n \rightarrow +\infty} \mathbb{P}(a \leq \bar{X}_n^* \leq b) = \int_a^b \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt \right)$$

Lorsque  $n$  est *suffisamment grand*, on pourra considérer que :  $\mathbb{P}(a \leq \bar{X}_n^* \leq b) \simeq \mathbb{P}(a \leq Z \leq b)$  où  $Z \hookrightarrow \mathcal{N}(0; 1)$ .

# I INTRODUCTION

Plaçons nous dans la peau d'enquêteurs souhaitant réaliser un sondage sur l'avis des français concernant la qualité de l'enseignement privé par rapport à l'enseignement public.

On interroge alors des personnes, choisies de façon la plus aléatoire possible et en veillant à ce que cet échantillon soit le plus représentatif possible de la population française. A la question "Pensez-vous que l'on suit une meilleure scolarité dans un établissement privé que dans un établissement public ?", la personne sondée répond "oui" ou "non".

On attribue à "oui" la valeur 1 et à "non" la valeur 0.

On note  $\theta$  la proportion de français votant "oui" et l'on considère donc une variable aléatoire  $X$  sur un certain espace probabilisable, telle que  $X \leftrightarrow \mathcal{B}(\theta)$ . **L'objectif est alors d'estimer au mieux la valeur de  $\theta$ .**

Pour cela, revenons à notre sondage...

On considère une suite de variables aléatoires  $(X_k)_{k \in \mathbb{N}^*}$ , toutes de même loi que  $X$ , telles que : pour tout  $k \in \mathbb{N}^*$ ,  $X_k$  prend la valeur 1 si la  $k$ -ième personne interrogée répond "oui" à la question, et 0 sinon.

Pour tout  $k \in \mathbb{N}^*$ , on notera  $x_k$  la réalisation de la variable aléatoire  $X_k$  dans le sondage effectué.

Le caractère aléatoire ne repose pas dans la réponse de l'individu choisi, mais bien dans le fait que l'individu fasse partie du panel de sondés ou non. Bien évidemment, on comprend le risque de la méthode : choisir un échantillon qui ne soit pas suffisamment représentatif de la population et qui ne nous fournit donc pas une bonne estimation de l'avis général de la population.

L'objectif est de pouvoir conclure un sondage par une phrase du type :

"la valeur observée sur l'échantillon est *proche* de la valeur théorique de  $\theta$  avec *tel niveau de confiance*"

L'enjeu mathématique est alors de préciser le "*proche*" et le "*tel niveau de confiance*" de la phrase précédente.

On considère un phénomène aléatoire, modélisé par une variable aléatoire  $X$  définie sur un espace probabilisable  $(\Omega, \mathcal{A})$ , dont la loi appartient à une famille de lois bien précise, dépendant d'un ou deux paramètres réels que l'on cherche à déterminer.

## EXEMPLES 1

**E1** On modélise le phénomène par une variable aléatoire  $X$  suivant une loi de Bernoulli, comme dans l'introduction ci-dessus, mais on ne connaît pas le paramètre de cette loi, que l'on note  $\theta$  et qui est à chercher dans l'ensemble  $\Theta = [0; 1]$ .

**E2** On modélise la durée de vie d'une ampoule par une variable aléatoire  $X$  suivant une loi exponentielle, dont on ne connaît pas le paramètre, noté  $\theta$ , qui est à chercher dans l'ensemble  $\Theta = \mathbb{R}_*^+$ .

**E3** On modélise le poids des nouveaux-nés par une variable aléatoire  $X$  suivant une loi normale, dont on ne connaît pas le couple des paramètres, noté  $\theta = (\theta_1, \theta_2)$ , qui est à chercher dans l'ensemble  $\Theta = \mathbb{R} \times \mathbb{R}_*^+$ .

On pourra proposer une unique valeur possible de ce paramètre, on parle alors d'**estimation ponctuelle**, ou bien tout un intervalle de valeurs possibles, on parle alors d'**estimation par intervalle de confiance**.

Dans tout le chapitre,  $n$  désignera un entier naturel non nul. On considèrera également une variable aléatoire  $X$  définie sur un espace probabilisable  $(\Omega, \mathcal{A})$  dont la loi dépend d'un paramètre  $\theta$  réel prenant ses valeurs dans un ensemble  $\Theta$ .  
On munit l'espace probabilisable d'une famille de probabilités  $(\mathbb{P}_\theta)_{\theta \in \Theta}$ .

### Petite remarque

Le sondage permet de n'enquêter que sur un groupe restreint d'individus, par manque de temps ou de moyen pour interroger la population entière. Sinon, on organise un référendum.

### Petite remarque

Majuscules pour les variables aléatoires, minuscules pour les réalisations.

### Trivialement...

On imagine bien que la taille de l'échantillon interrogé a un impact assez fort sur la qualité de la prévision faite.

### En fait...

En pratique, on commence par les outils de statistique descriptive pour analyser les données et essayer d'imaginer la famille de lois modélisant la situation : normale, exponentielle, Poisson, log-normale, Pareto...  
Après avoir fixé une famille de lois, on cherche le ou les paramètres de cette dernière par les méthodes de statistique inférentielle.

### Petite remarque

Si la loi dépend de deux paramètres inconnus, on cherchera toujours à estimer l'un puis l'autre... On peut donc considérer dans la suite du cours, que la loi ne dépend que d'un seul paramètre à déterminer.

# II ESTIMATION PONCTUELLE

## II.1 ÉCHANTILLON ET ESTIMATEUR

### DÉFINITION 1

### ÉCHANTILLON

Un  **$n$ -échantillon de  $X$**  est un  $n$ -uplet  $(X_1, \dots, X_n)$  de variables aléatoires définies sur  $(\Omega, \mathcal{A})$  indépendantes et de même loi que  $X$  pour toutes les  $\mathbb{P}_\theta$ .

### DÉFINITIONS 2

### ESTIMATEUR, ESTIMATION

Soit  $(X_1, \dots, X_n)$  un  $n$ -échantillon de  $X$ .

**D1** Un **estimateur de  $\theta$**  est une variable aléatoire sur  $(\Omega, \mathcal{A})$  fonction des variables aléatoires  $X_1, \dots, X_n$  et dont l'expression ne fait pas mention de  $\theta$ .

Autrement dit, une variable aléatoire  $T_n$  sur  $(\Omega, \mathcal{A})$  est un estimateur de  $\theta$  lorsqu'il existe une fonction  $\varphi$  définie sur  $X_1(\Omega) \times \dots \times X_n(\Omega)$  et à valeurs dans  $\mathbb{R}$  telle que  $T_n = \varphi(X_1, \dots, X_n)$ .

### Important !

La loi d'un estimateur de  $\theta$  dépend naturellement de  $\theta$  (puisque l'estimateur est fonction des  $X_k$  dont la loi dépend de  $\theta$ ) ; **mais son expression ne doit pas faire intervenir  $\theta$ .**

**Notation**

Si existence, on notera  $\mathbb{E}_\theta(T_n)$  et  $\mathbb{V}_\theta(T_n)$  l'espérance et la variance de  $T_n$  pour  $\mathbb{P}_\theta$ .

**D2** Soit  $\varphi(X_1, \dots, X_n)$  est un estimateur de  $\theta$ . Soit  $(x_1, \dots, x_n) \in X_1(\Omega) \times \dots \times X_n(\Omega)$ . On dit que  $\varphi(x_1, \dots, x_n)$  est une **estimation ponctuelle de  $\theta$** .  
Autrement dit, une estimation ponctuelle de  $\theta$  est une réalisation d'un estimateur de  $\theta$ .

**EXEMPLE 2**

Soit  $X$  une variable aléatoire suivant la loi  $\mathcal{B}(\theta)$ , avec  $\theta \in ]0; 1[$ . On considère  $X_1, \dots, X_n$  des variables aléatoires indépendantes suivant toutes la même loi que  $X$  et on note :  $T_n = \frac{1}{n} \sum_{k=1}^n X_k$ .

Justifions que  $T_n$  est un estimateur de  $\theta$ .

On sait que :

- ✓  $X_1, \dots, X_n$  sont indépendantes et suivent toutes la même loi  $\mathcal{B}(\theta)$ ,
- ✓  $T_n$  est fonction de  $X_1, \dots, X_n$  dont l'expression ne fait pas apparaître  $\theta$ .

Par conséquent :  $T_n$  est un estimateur de  $\theta$ .

De même :  $\prod_{k=1}^n X_k, \max(X_1, \dots, X_n), \min(X_1, \dots, X_n)$  sont des estimateurs de  $\theta$ .

En revanche,  $\sum_{k=1}^n X_k - \theta$  n'est pas un estimateur de  $\theta$ , car son expression fait intervenir  $\theta$ .

**Petite remarque**

Ce n'est pas ce qui manque des estimateurs, il y en a une infinité : la définition est très générale. Mais lesquels sont des *bons* estimateurs ?

**II.2 OUTILS DE MESURE DE LA QUALITÉ D'UN ESTIMATEUR (HP)**

Dans ce paragraphe, nous introduisons des outils, tous devenus hors programme, qui permettent de mesurer la qualité d'un estimateur.

L'idée générale est les réalisations de l'estimateur  $T_n$  soient proches de  $\theta$ . Pour mesurer l'écart entre  $T_n$  et  $\theta$ , on peut utiliser les quantités suivantes :  $T_n - \theta, |T_n - \theta|, (T_n - \theta)^2, \dots$

Mais bien évidemment, ce qui nous intéresse le plus est l'écart moyen... D'où les définitions suivantes :

**DÉFINITIONS 3**

**BIAIS, RISQUE QUADRATIQUE (HP)**

Soit  $T_n$  un estimateur de  $\theta$ .

**D1** Si  $T_n$  admet une espérance, le **biais de l'estimateur  $T_n$** , noté  $b_\theta(T_n)$ , est le réel défini par :

$$b_\theta(T_n) = \mathbb{E}_\theta(T_n - \theta)$$

**D2** Si  $T_n$  admet une variance, le **risque quadratique de l'estimateur  $T_n$** , noté  $r_\theta(T_n)$ , est le réel défini par :

$$r_\theta(T_n) = \mathbb{E}_\theta((T_n - \theta)^2)$$

**Vocabulaire**

- Si  $b_\theta(T_n) = 0$ , on dit que l'estimateur  $T_n$  est **sans biais**.
- Si  $\lim_{n \rightarrow +\infty} b_\theta(T_n) = 0$ , on dit que l'estimateur  $T_n$  est **asymptotiquement sans biais**.

**PROPRIÉTÉS 1**

**(HP)**

Avec les notations précédentes :

**P1** Si  $T_n$  admet une espérance, alors :  $b_\theta(T_n) = \mathbb{E}_\theta(T_n) - \theta$ .

**P2** Si  $T_n$  admet une variance, alors :  $r_\theta(T_n) = \mathbb{V}_\theta(T_n) + (b_\theta(T_n))^2$ .

**Petite remarque**

Si  $T_n$  est sans biais, alors  $r_\theta(T_n) = \mathbb{V}_\theta(T_n)$ .

**\* DÉMONSTRATION :**

**P1.** Supposons que  $T_n$  possède une espérance. Ainsi :

$$\begin{aligned} b_\theta(T_n) &= \mathbb{E}_\theta(T_n - \theta) \\ &= \mathbb{E}_\theta(T_n) - \theta \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{linéarité de l'espérance}$$

**P2.** Supposons que  $T_n$  possède une variance. Ainsi :

$$\begin{aligned} r_\theta(T_n) &= \mathbb{E}_\theta((T_n - \theta)^2) \\ &= \mathbb{E}_\theta(T_n^2 - 2\theta T_n + \theta^2) \\ &= \mathbb{E}_\theta(T_n^2) - 2\theta \mathbb{E}_\theta(T_n) + \theta^2 \\ &= \mathbb{V}_\theta(T_n) + (\mathbb{E}_\theta(T_n))^2 - 2\theta \mathbb{E}_\theta(T_n) + \theta^2 \\ &= \mathbb{V}_\theta(T_n) + (\mathbb{E}_\theta(T_n) - \theta)^2 \\ &= \mathbb{V}_\theta(T_n) + (b_\theta(T_n))^2 \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \begin{array}{l} \text{linéarité de l'espérance, toutes existent} \\ \text{formule de Koenig-Huuggens} \end{array}$$

$$\left. \begin{array}{l} \\ \end{array} \right\} \text{point précédent}$$

\*

Dans l'idéal, on cherche à obtenir un estimateur dont le biais et le risque quadratique sont les plus proches de 0 possible, l'idéal étant un estimateur sans biais et de variance minimale... Mais celui-ci n'existe pas toujours, ou n'est pas simple à trouver.

De façon générale, pour comparer deux estimateurs :

- s'ils ont même biais, on préférera celui de risque quadratique minimal,
- sinon, on pourra parfois préférer un estimateur biaisé à un estimateur sans biais si son risque quadratique est plus faible que la variance de l'estimateur sans biais.

DÉFINITION 4	ESTIMATEUR CONVERGENT (HP)
<p>Soit <math>(T_n)_{n \in \mathbb{N}^*}</math> une suite d'estimateurs de <math>\theta</math>.                  On dit que <math>(T_n)_{n \in \mathbb{N}^*}</math> est convergente (ou plus simplement que l'estimateur <math>T_n</math> est convergent) lorsque :</p> $\forall \theta \in \Theta, \forall \varepsilon > 0, \lim_{n \rightarrow +\infty} \mathbb{P}_\theta( T_n - \theta  \geq \varepsilon) = 0$	

**Confusion d'objets !**  
 On verra parfois la confusion entre estimateur et suite d'estimateurs.

**Autrement dit :**  
 $T_n$  est convergent ssi  $T_n \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \theta$

**EXEMPLE 3**

Soit  $T_n$  un estimateur de  $\theta$  admettant une variance. Montrons que si, pour tout  $\theta \in \Theta$ ,  $\lim_{n \rightarrow +\infty} r_\theta(T_n) = 0$ , alors  $T_n$  est convergent.

Soit  $\theta \in \Theta$ . Supposons que  $\lim_{n \rightarrow +\infty} r_\theta(T_n) = 0$ . Montrons :  $\forall \varepsilon > 0, \lim_{n \rightarrow +\infty} \mathbb{P}_\theta(|T_n - \theta| \geq \varepsilon) = 0$ .

Soit  $\varepsilon > 0$ . Soit  $n \in \mathbb{N}^*$ .

La variable aléatoire  $(T_n - \theta)^2$  :

- ✓ est à valeurs positives,
- ✓ admet une espérance (car  $T_n$  admet une variance).

Ainsi, d'après l'inégalité de Markov :

$$\forall a > 0, \mathbb{P}_\theta((T_n - \theta)^2 \geq a) \leq \frac{\mathbb{E}_\theta(T_n - \theta)^2}{a}$$

Avec  $a = \varepsilon^2 > 0$  :

$$\mathbb{P}_\theta((T_n - \theta)^2 \geq \varepsilon^2) \leq \frac{r_\theta(T_n)}{\varepsilon^2}$$

Mais, par stricte croissance de  $\sqrt{\cdot}$  sur  $\mathbb{R}^+$  et comme  $|\varepsilon| = \varepsilon$  (car  $\varepsilon > 0$ ) :

$$[(T_n - \theta)^2 \geq \varepsilon^2] = [|T_n - \theta| \geq \varepsilon]$$

On a ainsi établi :

$$\forall n \in \mathbb{N}^*, 0 \leq \mathbb{P}_\theta(|T_n - \theta| \geq \varepsilon) \leq \frac{r_\theta(T_n)}{\varepsilon^2}$$

Or  $\lim_{n \rightarrow +\infty} r_\theta(T_n) = 0$ . Donc, par théorème d'encadrement :

$$\lim_{n \rightarrow +\infty} \mathbb{P}_\theta(|T_n - \theta| \geq \varepsilon) = 0$$

**Conclusion :** si, pour tout  $\theta \in \Theta$ ,  $\lim_{n \rightarrow +\infty} r_\theta(T_n) = 0$ , alors  $T_n$  est convergent.

**À retenir...**  
 Cas classique pour établir la convergence d'un estimateur.

Un estimateur sans biais garantit l'absence d'erreur en moyenne, mais peut produire de très mauvaises estimations ponctuelles. Au contraire, un estimateur convergent, biaisé ou non, sera d'autant plus fiable que la taille de l'échantillon est grande. En général, on préfère ces derniers.

### II.3 MOYENNE ET VARIANCE EMPIRIQUES

Dans cette partie, nous allons donner et étudier deux estimateurs usuels et pertinents, l'un pour l'espérance et l'autre pour la variance. Ces estimateurs peuvent ainsi être utilisés pour déterminer les paramètres des lois  $\mathcal{B}(p)$ ,  $\mathcal{P}(\lambda)$ ,  $\mathcal{U}([1; n])$ ,  $\mathcal{E}(\lambda)$ ,  $\mathcal{U}([0; b])$ ,  $\mathcal{N}(\mu; \sigma^2)$ ... puisque chacune des lois fait intervenir espérance et/ou variance en paramètres !

DÉFINITION 5	MOYENNE EMPIRIQUE
<p>Soient <math>X</math> une variable aléatoire admettant une espérance et <math>(X_1, \dots, X_n)</math> un <math>n</math>-échantillon de <math>X</math>.</p> <p>La <b>moyenne empirique</b> de <math>X_1, \dots, X_n</math> est l'estimateur <math>\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k</math>.</p>	

**Petite remarque**  
 C'est l'estimateur le plus naturel et le plus usuel de l'espérance d'une variable aléatoire !

Soient  $X$  une variable aléatoire admettant une espérance et une variance ainsi que  $(X_1, \dots, X_n)$  un  $n$ -échantillon de  $X$ .

La moyenne empirique de  $X_1, \dots, X_n$  est un estimateur sans biais et convergent de  $\mathbb{E}(X)$ .

\* DÉMONSTRATION :

- Remarquons déjà que :
  - ✓  $X_1, \dots, X_n$  sont indépendantes et suivent la même loi,
  - ✓  $\bar{X}_n$  est fonction de  $X_1, \dots, X_n$  dont l'expression ne fait pas apparaître  $\mathbb{E}(X)$ .

**Conclusion :**  $\bar{X}_n$  est un estimateur de  $\mathbb{E}(X)$ .

- Soit  $n \in \mathbb{N}^*$ . La variable aléatoire  $\bar{X}_n$  est une combinaison linéaire de variables aléatoires admettant une espérance, elle admet donc une espérance. Puis :

$$\begin{aligned} \mathbb{E}(\bar{X}_n) &= \mathbb{E}\left(\frac{1}{n} \sum_{k=1}^n X_k\right) && \left. \begin{array}{l} \text{linéarité de l'espérance, toutes existent} \\ \forall k \in \llbracket 1; n \rrbracket, \mathbb{E}(X_k) = \mathbb{E}(X) \end{array} \right\} \\ &= \frac{1}{n} \sum_{k=1}^n \mathbb{E}(X_k) \\ &= \mathbb{E}(X) \end{aligned}$$

**Conclusion :**  $\bar{X}_n$  est un estimateur sans biais de  $\mathbb{E}(X)$ .

- Enfin, puisque les variables aléatoires  $X_1, \dots, X_n$  sont indépendantes, de même espérance et de même variance (car ont la même loi), d'après la loi faible des grands nombres :

$$\forall \varepsilon > 0, \lim_{n \rightarrow +\infty} \mathbb{P}(|\bar{X}_n - \mathbb{E}(X)| \geq \varepsilon) = 0$$

**Conclusion :**  $\bar{X}_n$  est un estimateur convergent de  $\mathbb{E}(X)$ .

\*

EXEMPLES 4

**E1** Dans le cas du sondage initial,  $X$  suivait une loi de Bernoulli de paramètre  $p$  inconnu. Puisque  $\mathbb{E}(X) = p$ , on peut obtenir une estimation du paramètre  $p$  de la loi de Bernoulli en considérant la réalisation de la moyenne empirique sur un échantillon suffisamment grand.

On retrouve notre intuition initiale !

Pour estimer  $\mathbb{V}(X) = p(1 - p)$ , on pourrait par exemple proposer  $\bar{X}_n(1 - \bar{X}_n)$ ...

**E2** On suppose maintenant que  $X$  suit un loi exponentielle de paramètre  $\lambda$  inconnu.

On peut utiliser la moyenne empirique d'un échantillon pour obtenir une estimation de  $\mathbb{E}(X) = \frac{1}{\lambda}$ , ce qui nous fournit ensuite une estimation de  $\lambda$ ...

**E3** On suppose maintenant que  $X$  suit un loi uniforme sur  $[0; b]$ , avec  $b$  inconnu.

On peut utiliser la moyenne empirique d'un échantillon pour obtenir une estimation de  $\mathbb{E}(X) = \frac{b}{2}$ , ce qui nous fournit ensuite une estimation de  $b$ ...

Petite remarque

Tout ce qui précède est à adapter dans le cas où l'on cherche à estimer  $g(\theta)$ , et pas  $\theta$  (où  $g$  est une fonction définie sur  $\Theta$ ). De façon générale, si  $\mathbb{E}(X) = f(\theta)$ , où  $f$  est bijective sur  $\Theta$ , alors une estimation de  $\theta$  sera fournie par une réalisation de  $f^{-1}(\bar{X}_n)$ .

DÉFINITION 6

VARIANCE EMPIRIQUE

Soient  $X$  une variable aléatoire admettant une espérance et une variance, ainsi que  $(X_1, \dots, X_n)$  un  $n$ -échantillon de  $X$ .

La **variance empirique** de  $X_1, \dots, X_n$  est l'estimateur  $V_n = \frac{1}{n} \sum_{k=1}^n (X_k - \bar{X}_n)^2$ .

EXEMPLE 5

Justifions que  $V_n$  possède une espérance et déterminons-la. Donnons, à partir de  $V_n$ , un estimateur  $V'_n$  qui soit un estimateur sans biais de  $\mathbb{V}(X)$ .

Soit  $n \in \mathbb{N}^*$ .

- On a :

$$\begin{aligned} V_n &= \frac{1}{n} \sum_{k=1}^n (X_k - \bar{X}_n)^2 \\ &= \frac{1}{n} \sum_{k=1}^n (X_k^2 - 2\bar{X}_n X_k + \bar{X}_n^2) \\ &= \frac{1}{n} \sum_{k=1}^n X_k^2 - 2\bar{X}_n \frac{1}{n} \sum_{k=1}^n X_k + \bar{X}_n^2 \end{aligned}$$

)  $\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k$

$$\begin{aligned}
&= \frac{1}{n} \sum_{k=1}^n X_k^2 - \bar{X}_n^2 \\
&= \frac{1}{n} \sum_{k=1}^n X_k^2 - \left( \frac{1}{n} \sum_{k=1}^n X_k \right)^2 \\
&= \frac{1}{n} \sum_{k=1}^n X_k^2 - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n X_i X_j \\
&= \frac{1}{n} \sum_{k=1}^n X_k^2 - \frac{1}{n^2} \sum_{i=1}^n X_i^2 + \frac{1}{n^2} \sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n X_i X_j \\
&= \frac{n-1}{n^2} \sum_{k=1}^n X_k^2 + \frac{1}{n^2} \sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n X_i X_j
\end{aligned}$$

Or :

- ✓ pour tout  $k \in \llbracket 1; n \rrbracket$ ,  $X_k^2$  admet une espérance car  $X$  admet une variance ;
- ✓ pour tout  $(i, j) \in \llbracket 1; n \rrbracket$  tel que  $i \neq j$ ,  $X_i$  et  $X_j$  admettent une espérance et sont indépendantes, donc  $X_i X_j$  admet une espérance et  $\mathbb{E}(X_i X_j) = \mathbb{E}(X_i) \mathbb{E}(X_j) = \mathbb{E}(X)^2$ .

Par conséquent,  $V_n$  est une combinaison linéaire de variables aléatoires admettant une espérance ;  $V_n$  admet donc une espérance et, par linéarité de l'espérance :

$$\begin{aligned}
\mathbb{E}(V_n) &= \frac{n-1}{n^2} \sum_{k=1}^n \mathbb{E}(X_k^2) + \frac{1}{n^2} \sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n \mathbb{E}(X_i X_j) \\
&= \frac{n-1}{n^2} \sum_{k=1}^n \mathbb{E}(X^2) + \frac{1}{n^2} \sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n \mathbb{E}(X)^2 \\
&= \frac{n-1}{n} \mathbb{E}(X^2) + \frac{n-1}{n} \mathbb{E}(X)^2 && \left. \begin{array}{l} \curvearrowright \\ \text{formule de Koenig-Huygens} \end{array} \right. \\
&= \frac{n-1}{n} \mathbb{V}(X) \\
&\neq \mathbb{V}(X)
\end{aligned}$$

**Conclusion :** la variance empirique est un estimateur (c'est bien le cas...) biaisé de  $\mathbb{V}(X)$ .

- Posons, pour  $n \geq 2$  :

$$V'_n = \frac{n}{n-1} V_n = \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X}_n)^2$$

- ✓  $V'_n$  est un estimateur de  $\mathbb{V}(X)$  (car  $X_1, \dots, X_n$  sont indépendantes et de même loi ; et que l'expression de  $V'_n$  ne fait pas apparaître  $\mathbb{V}(X)$ ),
- ✓ et d'après le calcul précédent :

$$\mathbb{E}(V'_n) = \mathbb{V}(X)$$

**Conclusion :**  $\frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X}_n)^2$  est un estimateur sans biais de  $\mathbb{V}(X)$ .

✎ **Pour info...**

De nombreux logiciels informatiques (et des calculatrices) utilisent l'estimateur  $V'_n$  (**variance empirique corrigée**) pour la variance plutôt que  $V_n$ .

### PROPRIÉTÉ 3

### QUALITÉ DE LA VARIANCE EMPIRIQUE CORRIGÉE (HP)

Soient  $X$  une variable aléatoire admettant un moment d'ordre 4 ainsi que  $(X_1, \dots, X_n)$  un  $n$ -échantillon de  $X$ . La variance empirique corrigée de  $X_1, \dots, X_n$  est un estimateur sans biais et convergent de  $\mathbb{V}(X)$ .

\* DÉMONSTRATION :

- Par construction  $\mathbb{E}(V'_n) = \mathbb{V}(X)$ .
- Et on a :

$$r_\theta(V'_n) = \mathbb{V}(V'_n) = \dots = \frac{1}{n} \mathbb{E}((X - \mathbb{E}(X))^4) - \frac{n-3}{n(n-1)} \mathbb{V}(X)^2$$

Ainsi :  $\lim_{n \rightarrow +\infty} r_\theta(V'_n) = 0$  et on conclut sur la convergence en utilisant l'inégalité de Markov et le théorème d'encadrement...

✎ **Petite remarque**

Cet estimateur permettrait d'obtenir une estimation du second paramètre d'une loi normale par exemple.

## II.4 ESTIMATEUR DU MAXIMUM DE VRAISEMBLANCE

### Un peu d'intuition... ?

- Si je vous donne une liste de 10000 réalisations indépendantes d'une certaine loi de Bernoulli, comment feriez-vous pour estimer la valeur du paramètre de cette loi ?
- Si je vous donne une liste de 10000 réalisations indépendantes d'une certaine loi de Poisson, comment feriez-vous pour estimer la valeur du paramètre de cette loi ?
- Si je vous donne une liste de 10000 réalisations indépendantes d'une certaine loi exponentielle, comment feriez-vous pour estimer la valeur du paramètre de cette loi ?
- Si je vous donne une liste de 10000 réalisations indépendantes d'une certaine loi uniforme, comment feriez-vous pour estimer la valeur des paramètres de cette loi ?

On considère toujours une variable aléatoire suivant une certaine loi, dont un paramètre est inconnu, et un  $n$ -échantillon  $(X_1, \dots, X_n)$  de  $X$  ayant fourni un  $n$ -uplet de réalisations  $(x_1, \dots, x_n)$ .

Une idée assez naturelle consiste à considérer que la valeur de  $\theta$  qui a permis de générer les observations est celle qui avait la plus grande probabilité de les générer.

### EXEMPLE 6

Commençons par un cas simple où  $n = 1$ . On considère une variable aléatoire  $X$  suivant la loi  $\mathcal{B}(15, p)$ , avec  $p$  inconnu. Une réalisation de cette variable aléatoire fournit la valeur 5.

Question : quelle valeur de  $p$  maximise  $\mathbb{P}([X = 5])$  ?

Puisque  $X \hookrightarrow \mathcal{B}(15, p)$ , on a :

$$\mathbb{P}([X = 5]) = \binom{15}{5} p^5 (1-p)^{10}$$

Posons  $f : x \mapsto \binom{15}{5} x^5 (1-x)^{10}$ .

La fonction  $f$  est polynomiale, donc dérivable sur  $[0; 1]$  et, pour tout  $x \in [0; 1]$  :

$$\begin{aligned} f'(x) &= \binom{15}{5} (5x^4(1-x)^{10} - 10x^5(1-x)^9) \\ &= \binom{15}{5} 5x^4(1-x)^9(1-x-2x) \\ &= \binom{15}{5} 5x^4(1-x)^9(1-3x) \end{aligned}$$

Mais, pour tout  $x \in [0; 1]$ ,  $x^4 \geq 0$  et  $(1-x)^9 \geq 0$ ; et :

$$1 - 3x > 0 \iff x < \frac{1}{3}$$

D'où :

$x$	0	$\frac{1}{3}$	1
$f'(x)$	+	0	-
$f$	0	↗ ... ↘	0

**Conclusion** : la probabilité  $\mathbb{P}([X = 5])$  est maximale lorsque  $p = \frac{1}{3}$ .  
Puisque l'évènement  $[X = 5]$  s'est réalisé, on peut penser que  $p = \frac{1}{3}$ .

### DÉFINITIONS 7

### FONCTION DE VRAISEMBLANCE

Soit  $X$  une variable aléatoire suivant une loi dépendant d'un certain paramètre  $\theta$  inconnu et  $(X_1, \dots, X_n)$  un  $n$ -échantillon de  $X$ . Soit  $(x_1, \dots, x_n) \in X_1(\Omega) \times \dots \times X_n(\Omega)$  un  $n$ -uplet d'observations.

**D1** Si  $X$  est discrète, on appelle **fonction de vraisemblance** pour  $(x_1, \dots, x_n)$  la fonction :

$$L_n : \theta \mapsto \prod_{k=1}^n \mathbb{P}_\theta([X_k = x_k])$$

**D2** Si  $X$  est à densité, de densité  $f_\theta$ , on appelle **fonction de vraisemblance** pour  $(x_1, \dots, x_n)$  la fonction :

$$L_n : \theta \mapsto \prod_{k=1}^n f_\theta(x_k)$$

### Important !

Puisqu'il est plus commode de dériver une somme qu'un produit (surtout pour les fonctions de deux variables), on travaille souvent sur la **log-vraisemblance**, égale à  $\ln \circ L_n$ .  
Les fonctions vraisemblance et log-vraisemblance ont mêmes variations et mêmes extrema (par stricte croissance de  $\ln$ )...



Soit  $X$  une variable aléatoire suivant une loi dépendant d'un certain paramètre  $\theta$  inconnu et  $(X_1, \dots, X_n)$  un  $n$ -échantillon de  $X$ . Soit  $(x_1, \dots, x_n) \in X_1(\Omega) \times \dots \times X_n(\Omega)$  un  $n$ -uplet d'observations. On note  $L_n$  la fonction de vraisemblance associée et on suppose qu'elle admet un maximum sur  $\Theta$ , atteint en un unique réel.

**D1** L'estimation du maximum de vraisemblance de  $\theta$  est le réel  $\theta_n^*$  en lequel  $L_n$  admet son maximum.

**D2** Si  $g$  est la fonction telle que  $\theta_n^* = g(x_1, \dots, x_n)$ , alors l'estimateur du maximum de vraisemblance de  $\theta$  est la variable aléatoire  $T_n = g(X_1, \dots, X_n)$ .

**EXEMPLE 7**

Déterminons l'estimateur du maximum de vraisemblance du paramètre  $p$  d'une loi de Bernoulli. Pour cela, considérons une variable aléatoire  $X$  suivant la loi  $\mathcal{B}(p)$ , avec  $p \in ]0; 1[$  inconnu. Considérons également, pour tout  $n \in \mathbb{N}^*$ ,  $(X_1, \dots, X_n)$  un  $n$ -échantillon de  $X$  et  $(x_1, \dots, x_n)$  une réalisation de cet échantillon. Soit  $n \in \mathbb{N}^*$ . On a :

$$\forall k \in \llbracket 1; n \rrbracket, \forall x_k \in \{0; 1\}, \mathbb{P}([X_k = x_k]) = \begin{cases} p & \text{si } x_k = 1 \\ 1 - p & \text{si } x_k = 0 \end{cases} = p^{x_k} (1 - p)^{1 - x_k}$$

Ainsi, en notant  $L_n$  la fonction de vraisemblance pour  $(x_1, \dots, x_n)$  et  $\varphi_n(p) = \ln(L_n(p))$ , on a, pour tout  $p \in ]0; 1[$  :

$$\begin{aligned} L_n(p) &= \prod_{k=1}^n \mathbb{P}([X_k = x_k]) \\ &= \prod_{k=1}^n p^{x_k} (1 - p)^{1 - x_k} \end{aligned}$$

D'où, pour tout  $p \in ]0; 1[$ ,  $L_n(p) > 0$  et :

$$\begin{aligned} \varphi_n(p) &= \ln(L_n(p)) \\ &= \ln\left(\prod_{k=1}^n p^{x_k} (1 - p)^{1 - x_k}\right) && \hookrightarrow \forall k \in \llbracket 1; n \rrbracket, p^{x_k} (1 - p)^{1 - x_k} > 0 \\ &= \sum_{k=1}^n \ln(p^{x_k} (1 - p)^{1 - x_k}) && \hookrightarrow \forall k \in \llbracket 1; n \rrbracket, p > 0, 1 - p > 0 \\ &= \sum_{k=1}^n (x_k \ln(p) + (1 - x_k) \ln(1 - p)) \\ &= \left(\sum_{k=1}^n x_k\right) \ln(p) + \left(n - \sum_{k=1}^n x_k\right) \ln(1 - p) && \hookrightarrow \text{en notant } s_n = \sum_{k=1}^n x_k \\ &= s_n \ln(p) + (n - s_n) \ln(1 - p) \end{aligned}$$

La fonction  $\varphi_n$  est donc dérivable sur  $]0; 1[$  et, pour tout  $p \in ]0; 1[$  :

$$\begin{aligned} \varphi_n'(p) &= s_n \frac{1}{p} - (n - s_n) \frac{1}{1 - p} \\ &= \frac{s_n(1 - p) - (n - s_n)p}{p(1 - p)} \\ &= \frac{s_n - np}{p(1 - p)} \end{aligned}$$

D'où :

$p$	0	$\frac{s_n}{n}$	1
$\varphi_n'(p)$	+	0	-
$\varphi_n$	↗ ... ↘		

La fonction  $\varphi_n$  admet donc un maximum en  $\frac{s_n}{n} = \frac{1}{n} \sum_{k=1}^n x_k$ .

Or, par stricte croissance de  $\ln$  sur  $\mathbb{R}_*^+$ , les fonctions  $\varphi_n$  et  $L_n$  ont les mêmes variations sur  $]0; 1[$ .

Par conséquent : la fonction  $L_n$  admet un maximum en  $\frac{1}{n} \sum_{k=1}^n x_k$ .

**Conclusion** : l'estimateur du maximum de vraisemblance de  $p$  est  $\frac{1}{n} \sum_{k=1}^n X_k = \bar{X}_n$ .

**Sympa !**

Le cadre du cours était celui d'un unique paramètre réel  $\theta$ . La fonction  $L_n$  est ainsi une fonction d'une variable. En revanche, si l'on cherche les deux paramètres d'une loi,  $L_n$  sera une fonction de deux variables... Vous voyez les liens ?!

**Pour info...**

L'estimateur du maximum de vraisemblance est toujours convergent, mais pas nécessairement sans biais.

**Pénible !**

On utilise souvent la même lettre pour le paramètre à déterminer et pour la variable de la fonction de vraisemblance...

**Petite remarque**

On retrouve la moyenne empirique... mais ce ne sera pas toujours le cas !

# III ESTIMATION PAR INTERVALLE DE CONFIANCE

## III.1 INTERVALLE DE CONFIANCE (EXACT)

### DÉFINITION 9

### INTERVALLE DE CONFIANCE (EXACT)

Soit  $X$  une variable aléatoire suivant une loi dépendant d'un certain paramètre  $\theta$  inconnu. Soient  $\alpha \in ]0; 1[$  ainsi que  $(U_n)_{n \in \mathbb{N}^*}$  et  $(V_n)_{n \in \mathbb{N}^*}$  deux suites de variables aléatoires.

On dit que  $[U_n, V_n]$  est un intervalle de confiance de  $\theta$  au niveau de confiance  $1 - \alpha$  (ou au risque  $\alpha$ ), lorsque :

- ✓  $(U_n)_{n \in \mathbb{N}^*}$  et  $(V_n)_{n \in \mathbb{N}^*}$  sont deux suites d'estimateurs de  $\theta$  telles que :  $\forall n \in \mathbb{N}^*, \mathbb{P}_\theta([U_n \leq V_n]) = 1$  ;
- ✓  $\mathbb{P}_\theta([\theta \in [U_n, V_n]]) \geq 1 - \alpha$

### ✗ Attention !

$\theta$  n'est pas une variable aléatoire ! En revanche,  $[U_n, V_n]$  est un intervalle aléatoire !  
On garde en tête :  
 $[\theta \in [U_n, V_n]] = [U_n \leq \theta] \cap [V_n \geq \theta]$   
et on voit mieux, sous cette forme, ce qui est aléatoire et ce qui ne l'est pas !

Si  $u_n$  et  $v_n$  sont des réalisations de  $U_n$  et  $V_n$  respectivement, il est faux de dire " $\theta$  est compris entre  $u_n$  et  $v_n$  avec une probabilité au moins égale à  $1 - \alpha$ ".  
En effet, rien d'aléatoire dans  $\theta, u_n$  et  $v_n$ .  
On dira : "on a une confiance de  $1 - \alpha$  dans le fait que  $\theta$  soit compris entre  $u_n$  et  $v_n$ ".

Voyons déjà un premier exemple :

### EXEMPLE 8

On considère une variable aléatoire  $X$  suivant une loi  $\mathcal{N}(\mu; \sigma^2)$ , où  $m$  est inconnu et  $\sigma^2$  connu non nul. On dispose d'un  $n$ -échantillon  $(X_1, \dots, X_n)$  de  $X$ .

### ES Rappel...

Si  $X_1, \dots, X_n$  sont indépendantes telles que pour tout  $k \in \llbracket 1; n \rrbracket, X_k \hookrightarrow \mathcal{N}(m_k, \sigma_k^2)$ , alors  $\sum_{k=1}^n a_k X_k$  suit la loi  $\mathcal{N}\left(\sum_{k=1}^n a_k m_k; \sum_{k=1}^n a_k^2 \sigma_k^2\right)$ .

- Donnons la loi de  $\bar{X}_n$  et celle de  $\bar{X}_n^* = \sqrt{n} \frac{\bar{X}_n - m}{\sigma}$ .  
Puisque  $X_1, \dots, X_n$  sont indépendantes et de même loi que  $X$ , on a :

$$\bar{X}_n \hookrightarrow \mathcal{N}\left(m, \frac{\sigma^2}{n}\right)$$

Par conséquent :

$$\frac{\bar{X}_n - m}{\sqrt{\frac{\sigma^2}{n}}} \hookrightarrow \mathcal{N}(0; 1)$$

Autrement dit :

$$\bar{X}_n^* \hookrightarrow \mathcal{N}(0; 1)$$

- Soit  $\alpha \in ]0; 1[$ . Justifions l'existence d'un unique réel  $t_\alpha$  tel que  $\Phi(t_\alpha) = 1 - \frac{\alpha}{2}$ , où  $\Phi$  désigne la fonction de répartition d'une variable aléatoire suivant la loi  $\mathcal{N}(0; 1)$ .  
Notons  $\varphi$  la densité continue sur  $\mathbb{R}$  associée à la loi  $\mathcal{N}(0; 1)$ . Puisque  $\varphi$  est continue sur  $\mathbb{R}$ , la fonction  $\Phi$  est  $\mathcal{C}^1$  sur  $\mathbb{R}$  et pour tout  $x \in \mathbb{R}$  :

$$\Phi'(x) = \varphi(x) > 0$$

Par conséquent :

- ✓  $\Phi$  est continue sur  $\mathbb{R}$ ,
- ✓  $\Phi$  est strictement croissante sur  $\mathbb{R}$ .

Donc par théorème de bijection,  $\Phi$  est bijective de  $\mathbb{R}$  dans  $\Phi(\mathbb{R})$ .

Or  $\Phi(\mathbb{R}) = ]0; 1[$ .

Ensuite, puisque  $\alpha \in ]0; 1[$ , on a  $1 - \frac{\alpha}{2} \in ]0; 1[$ .

Conclusion : il existe un unique réel  $t_\alpha$  tel que  $\Phi(t_\alpha) = 1 - \frac{\alpha}{2}$ .

### Pourquoi ?

Car  $\Phi$  est continue, strictement croissante, de limite 0 en  $-\infty$  et 1 en  $+\infty$ .

- Déduisons-en que  $\left[\bar{X}_n - t_\alpha \frac{\sigma}{\sqrt{n}}, \bar{X}_n + t_\alpha \frac{\sigma}{\sqrt{n}}\right]$  est un intervalle de confiance de  $m$  au niveau de confiance  $1 - \alpha$ .

- ✓ En notant  $U_n = \bar{X}_n - t_\alpha \frac{\sigma}{\sqrt{n}}$  et  $V_n = \bar{X}_n + t_\alpha \frac{\sigma}{\sqrt{n}}$ , on sait déjà que  $U_n$  et  $V_n$  sont des estimateurs de  $m$  (car fonctions d'un  $n$ -échantillon de  $X$  dont les expressions ne font pas apparaître  $m$ ) et que  $\mathbb{P}([U_n \leq V_n]) = 1$  (car  $t_\alpha \frac{\sigma}{\sqrt{n}} > 0$ ).

- ✓ Montrons ensuite que  $\mathbb{P}\left(\left[m \in \left[\bar{X}_n - t_\alpha \frac{\sigma}{\sqrt{n}}, \bar{X}_n + t_\alpha \frac{\sigma}{\sqrt{n}}\right]\right]\right) \geq 1 - \alpha$ .

On a :

$$\begin{aligned} \mathbb{P}\left(\left[m \in \left[\bar{X}_n - t_\alpha \frac{\sigma}{\sqrt{n}}, \bar{X}_n + t_\alpha \frac{\sigma}{\sqrt{n}}\right]\right]\right) &= \mathbb{P}\left(\left[\bar{X}_n - t_\alpha \frac{\sigma}{\sqrt{n}} \leq m \leq \bar{X}_n + t_\alpha \frac{\sigma}{\sqrt{n}}\right]\right) \\ &= \mathbb{P}\left(\left[-t_\alpha \frac{\sigma}{\sqrt{n}} \leq m - \bar{X}_n \leq t_\alpha \frac{\sigma}{\sqrt{n}}\right]\right) \end{aligned} \quad \left. \vphantom{\mathbb{P}} \right) \frac{\sigma}{\sqrt{n}} > 0$$

$$\begin{aligned}
&= \mathbb{P} \left( \left[ -t_\alpha \leq \frac{m - \bar{X}_n}{\frac{\sigma}{\sqrt{n}}} \leq t_\alpha \right] \right) \\
&= \mathbb{P} \left( \left[ t_\alpha \geq \sqrt{n} \frac{\bar{X}_n - m}{\sigma} \geq -t_\alpha \right] \right) \\
&= \mathbb{P} \left( [-t_\alpha \leq \bar{X}_n^* \leq t_\alpha] \right) \\
&= \Phi(t_\alpha) - \Phi(-t_\alpha) \\
&= 2\Phi(t_\alpha) - 1 \\
&= 1 - \alpha
\end{aligned}$$

$\bar{X}_n^* \hookrightarrow \mathcal{N}(0, 1)$ , donc  $\bar{X}_n^*$  est à densité  
 $\forall x \in \mathbb{R}, \Phi(-x) = 1 - \Phi(x)$   
 $\Phi(t_\alpha) = 1 - \frac{\alpha}{2}$

**Conclusion :**  $\left[ \bar{X}_n - t_\alpha \frac{\sigma}{\sqrt{n}}, \bar{X}_n + t_\alpha \frac{\sigma}{\sqrt{n}} \right]$  est un intervalle de confiance de  $m$  au niveau de confiance  $1 - \alpha$ .

Plus l'intervalle de confiance est réduit, plus le risque augmente; alors que plus l'intervalle de confiance a une grande amplitude, plus le risque est faible. En pratique, il faudra trouver un compromis entre les deux.

**ES Pour info...**  
 En pratique, un risque de 0,05 est acceptable... on cherchera donc souvent un intervalle de confiance de niveau de confiance 0,95 (95%).

**♣ MÉTHODE 1 ♣ Recherche d'un intervalle de confiance exact pour  $\mathbb{E}(X)$  avec l'inégalité de Bienaymé-Tchebychev.**  $X$  admet une espérance  $m$  inconnue et une variance  $\sigma^2$ . On considère  $T_n$  un estimateur de  $m$  tel que  $\mathbb{E}(T_n) = m$  (estimateur sans biais). Résumons les étapes de l'obtention de l'intervalle de confiance :

1. Inégalité de Bienaymé-Tchebychev :

$$\forall \varepsilon > 0, \mathbb{P}(|T_n - m| \geq \varepsilon) \leq \frac{\mathbb{V}(T_n)}{\varepsilon^2}$$

2. Si  $\mathbb{V}(T_n)$  dépend de  $m$ , on majore  $\mathbb{V}(T_n)$  par  $v_n$ , indépendant de  $m$  (sinon, on prend  $v_n = \mathbb{V}(T_n)$ ) et ainsi :

$$\forall \varepsilon > 0, \mathbb{P}(|T_n - m| \geq \varepsilon) \leq \frac{v_n}{\varepsilon^2}$$

3. Par passage à l'évènement contraire, on obtient :

$$\forall \varepsilon > 0, \mathbb{P}(|T_n - m| < \varepsilon) \geq 1 - \frac{v_n}{\varepsilon^2}$$

Autrement dit :

$$\forall \varepsilon > 0, \mathbb{P}([T_n - \varepsilon < m < T_n + \varepsilon]) \geq 1 - \frac{v_n}{\varepsilon^2}$$

4. Et puisque  $[T_n - \varepsilon < m < T_n + \varepsilon] \subset [T_n - \varepsilon \leq m \leq T_n + \varepsilon]$ , on obtient :

$$\forall \varepsilon > 0, \mathbb{P}([m \in [T_n - \varepsilon; T_n + \varepsilon]]) \geq 1 - \frac{v_n}{\varepsilon^2}$$

Reste à choisir  $\varepsilon$  tel que, pour le  $\alpha$  donné, on ait  $\frac{v_n}{\varepsilon^2} = \alpha$ ... et ainsi :

$$\mathbb{P}([m \in [T_n - \varepsilon; T_n + \varepsilon]]) \geq 1 - \alpha$$

**★ Classique ! ★**  
 C'est très classique (et le seul cas qui semble être au programme pour cette partie du cours). On prendra souvent (toujours ?) la moyenne empirique comme estimateur de la moyenne.

**✗ Attention !**  
 Le niveau de confiance ne peut pas dépendre du paramètre  $m$  que l'on cherche à estimer.

**Petite remarque**  
 L'intervalle de confiance proposé est centré en  $T_n$ , d'amplitude  $2\varepsilon$ .

**EXEMPLE 9**

On considère une variable aléatoire suivant une loi de Bernoulli de paramètre  $p$  inconnu ainsi qu'un  $n$ -échantillon  $(X_1, \dots, X_n)$  de  $X$ . On note  $\bar{X}_n$  la moyenne empirique de cet échantillon.

- Montrons :  $\forall \varepsilon > 0, \mathbb{P}([\bar{X}_n - \varepsilon \leq p \leq \bar{X}_n + \varepsilon]) \geq 1 - \frac{1}{4n\varepsilon^2}$ .  
 Soit  $\varepsilon > 0$ .

\* On sait que  $\bar{X}_n$  admet une espérance et une variance et que :

$$\mathbb{E}(\bar{X}_n) = p ; \mathbb{V}(\bar{X}_n) = \frac{p(1-p)}{n}$$

\* D'après l'inégalité de Bienaymé-Tchebychev, licite car  $\bar{X}_n$  admet une variance et que  $\varepsilon > 0$  :

$$\mathbb{P}(|\bar{X}_n - \mathbb{E}(\bar{X}_n)| \geq \varepsilon) \leq \frac{\mathbb{V}(\bar{X}_n)}{\varepsilon^2}$$

Autrement dit :

$$\mathbb{P}(|\bar{X}_n - p| \geq \varepsilon) \leq \frac{p(1-p)}{n\varepsilon^2}$$

Mais, on sait que la fonction  $x \mapsto x(1-x)$  est majorée par  $\frac{1}{4}$ ; d'où :

$$\mathbb{P}(|\bar{X}_n - m| \geq \varepsilon) \leq \frac{1}{4n\varepsilon^2}$$

Or :

$$\begin{aligned} \mathbb{P}(|\bar{X}_n - \rho| \geq \varepsilon) &= 1 - \mathbb{P}(|\bar{X}_n - \rho| < \varepsilon) \\ &= 1 - \mathbb{P}(\bar{X}_n - \varepsilon < \rho < \bar{X}_n + \varepsilon) \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \varepsilon > 0$$

D'où :

$$\mathbb{P}(\bar{X}_n - \varepsilon < \rho < \bar{X}_n + \varepsilon) \geq 1 - \frac{1}{4n\varepsilon^2}$$

Enfin, on a  $[\bar{X}_n - \varepsilon < \rho < \bar{X}_n + \varepsilon] \subset [\bar{X}_n - \varepsilon \leq \rho \leq \bar{X}_n + \varepsilon]$  et par croissance de  $\mathbb{P}$ , on obtient par transitivité :

$$\mathbb{P}(\bar{X}_n - \varepsilon \leq \rho \leq \bar{X}_n + \varepsilon) \geq 1 - \frac{1}{4n\varepsilon^2}$$

#### Petites remarques

- En se fixant un risque  $\alpha = \frac{1}{4n\varepsilon^2}$ , plus  $\varepsilon$  diminue et plus  $n$  doit être grand pour conserver le même risque.
- En prenant  $\alpha = \frac{1}{4n\varepsilon^2}$ , on voit aussi que, à  $n$  fixé, plus  $\varepsilon$  est petit et plus  $\alpha$  est grand... et plus  $\varepsilon$  est grand, plus  $\alpha$  est petit.

- Déduisons-en que  $\left[ \bar{X}_n - \sqrt{\frac{5}{n}}, \bar{X}_n + \sqrt{\frac{5}{n}} \right]$  est un intervalle de confiance de  $p$  au niveau de confiance 95%.

✓ En notant  $U_n = \bar{X}_n - \sqrt{\frac{5}{n}}$  et  $V_n = \bar{X}_n + \sqrt{\frac{5}{n}}$ , on sait déjà que  $U_n$  et  $V_n$  sont des estimateurs de  $p$  (car fonctions d'un  $n$ -échantillon de  $X$  dont les expressions ne font pas apparaître  $p$ ) et que  $\mathbb{P}(U_n \leq V_n) = 1$ .

✓ Montrons ensuite que  $\mathbb{P}\left(\left[p \in \left[\bar{X}_n - \sqrt{\frac{5}{n}}, \bar{X}_n + \sqrt{\frac{5}{n}}\right]\right]\right) \geq 0,95$ .

On a :

$$\begin{aligned} \mathbb{P}\left(\left[p \in \left[\bar{X}_n - \sqrt{\frac{5}{n}}, \bar{X}_n + \sqrt{\frac{5}{n}}\right]\right]\right) &= \mathbb{P}\left(\left[\bar{X}_n - \sqrt{\frac{5}{n}} \leq p \leq \bar{X}_n + \sqrt{\frac{5}{n}}\right]\right) \\ &\geq 1 - \frac{1}{20} = 0,95 \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{point précédent avec } \varepsilon = \sqrt{\frac{5}{n}}$$

**Conclusion :**  $\left[ \bar{X}_n - \sqrt{\frac{5}{n}}, \bar{X}_n + \sqrt{\frac{5}{n}} \right]$  est un intervalle de confiance de  $p$  au niveau de confiance 95%.

- Au second tour d'une élection présidentielle, les citoyens ont le choix entre deux candidats  $A$  et  $B$ . Un institut de sondage réalise un sondage auprès de 2000 personnes, dont 1150 affirment vouloir voter pour le candidat  $A$ . Peut-on affirmer, avec un risque d'erreur de 5% que la candidat  $A$  sera élu ?

Considérons  $X$  une variable aléatoire suivant la loi de Bernoulli de paramètre  $p$ , où  $p$  est la probabilité que le candidat  $A$  soit élu.

Pour tout  $k \in \mathbb{N}^*$ ,  $X_k$  est la variable aléatoire suivant la loi de Bernoulli de succès "la  $k$ -ième personne interrogée vote le candidat  $A$ ".

En supposant le panel du sondage bien choisi, on peut considérer que  $(X_1, \dots, X_{2000})$  est un échantillon de  $X$ .

Notons  $\bar{x}_{2000}$  la réalisation de  $\bar{X}_{2000}$  obtenue sur cet échantillon :  $\bar{x}_{2000} = \frac{1150}{2000}$ .

Dans ce cas, d'après ce qui précède, on peut estimer que  $p$  se situe dans l'intervalle  $\left[ \bar{x}_{2000} - \sqrt{\frac{5}{2000}}, \bar{x}_{2000} + \sqrt{\frac{5}{2000}} \right]$

avec un niveau de confiance 95%.

Mais :

$$\begin{aligned} \left[ \bar{x}_{2000} - \sqrt{\frac{5}{2000}}, \bar{x}_{2000} + \sqrt{\frac{5}{2000}} \right] &= \left[ \frac{1150}{2000} - \sqrt{\frac{1}{400}}, \frac{1150}{2000} + \sqrt{\frac{1}{400}} \right] \\ &= [0,575 - 0,05; 0,575 + 0,05] \\ &= [0,525; 0,625] \end{aligned}$$

On peut donc affirmer, avec un niveau de confiance de 95%, que la probabilité que le candidat  $A$  soit élu se situe dans l'intervalle  $[0,525; 0,625]$ .

**Conclusion :** on peut affirmer, avec un risque d'erreur de 5%, que la candidat  $A$  sera élu.

DÉFINITION 10

INTERVALLE DE CONFIANCE ASYMPTOTIQUE

Soit  $X$  une variable aléatoire suivant une loi dépendant d'un certain paramètre  $\theta$  inconnu. Soient  $\alpha \in ]0; 1[$  ainsi que  $(U_n)_{n \in \mathbb{N}^*}$  et  $(V_n)_{n \in \mathbb{N}^*}$  deux suites de variables aléatoires.

On dit que  $[U_n, V_n]$  est un intervalle de confiance asymptotique de  $\theta$  au niveau de confiance  $1 - \alpha$  (ou au risque  $\alpha$ ), lorsque :

- ✓  $(U_n)_{n \in \mathbb{N}^*}$  et  $(V_n)_{n \in \mathbb{N}^*}$  sont deux suites d'estimateurs de  $\theta$  telles que :  $\forall n \in \mathbb{N}^*, \mathbb{P}_\theta([U_n \leq V_n]) = 1$  ;
- ✓  $\lim_{n \rightarrow +\infty} \mathbb{P}_\theta([\theta \in [U_n, V_n]]) \geq 1 - \alpha$

À retenir...

De façon plus générale, comme nous le verrons en exercice dans certains cas, une convergence en loi peut nous fournir un intervalle de confiance asymptotique.

♣ MÉTHODE 2 ♣ Recherche d'un intervalle de confiance asymptotique pour  $\mathbb{E}(X)$  avec le TCL.

$X$  admet une espérance  $m$  inconnue et une variance  $\sigma^2$  non nulle. On travaille avec l'estimateur  $\bar{X}_n$  de  $m$  (estimateur sans biais). Résumons les étapes de l'obtention de l'intervalle de confiance asymptotique :

- En posant  $\bar{X}_n^* = \sqrt{n} \frac{\bar{X}_n - m}{\sigma}$ , on a, par le TCL :

$$\bar{X}_n^* \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} Z$$

où  $Z \hookrightarrow \mathcal{N}(0; 1)$ .

- Ainsi :

$$\forall x \geq 0, \lim_{n \rightarrow +\infty} \mathbb{P}([-x \leq \bar{X}_n^* \leq x]) = \mathbb{P}([-x \leq Z \leq x]) = 2\Phi(x) - 1$$

- Or :

$$\forall x \geq 0, \mathbb{P}([-x \leq \bar{X}_n^* \leq x]) = \dots = \mathbb{P}\left(\left[\bar{X}_n - \frac{\sigma x}{\sqrt{n}} \leq m \leq \bar{X}_n + \frac{\sigma x}{\sqrt{n}}\right]\right)$$

D'où :

$$\lim_{n \rightarrow +\infty} \mathbb{P}\left(\left[\bar{X}_n - \frac{\sigma x}{\sqrt{n}} \leq m \leq \bar{X}_n + \frac{\sigma x}{\sqrt{n}}\right]\right) = 2\Phi(x) - 1$$

- Si  $\sigma$  dépend de  $m$ , on majore  $\sigma$  par  $s$ , indépendant de  $m$  (sinon, on prend  $s = \sigma$ ) et ainsi :

$$\left[\bar{X}_n - \frac{\sigma x}{\sqrt{n}} \leq m \leq \bar{X}_n + \frac{\sigma x}{\sqrt{n}}\right] \subset \left[\bar{X}_n - \frac{s x}{\sqrt{n}} \leq m \leq \bar{X}_n + \frac{s x}{\sqrt{n}}\right]$$

D'où, par croissance de  $\mathbb{P}$  :

$$\lim_{n \rightarrow +\infty} \mathbb{P}\left(\left[\bar{X}_n - \frac{s x}{\sqrt{n}} \leq m \leq \bar{X}_n + \frac{s x}{\sqrt{n}}\right]\right) \geq 2\Phi(x) - 1$$

- Reste à choisir  $x$  tel que, pour le  $\alpha$  donné, on ait  $2\Phi(x) - 1 = 1 - \alpha$ .

Mais :

$$2\Phi(x) - 1 = 1 - \alpha \iff \Phi(x) = 1 - \frac{\alpha}{2} \iff x = \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)$$

$\left. \begin{array}{l} \Phi \text{ est bijective de } \mathbb{R} \text{ dans } ]0; 1[ \text{ et} \\ 1 - \frac{\alpha}{2} \in ]0; 1[ \end{array} \right\}$

En notant  $t_\alpha = \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)$ , on obtient :

$$\lim_{n \rightarrow +\infty} \mathbb{P}\left(\left[m \in \left[\bar{X}_n - t_\alpha \frac{s}{\sqrt{n}}, \bar{X}_n + t_\alpha \frac{s}{\sqrt{n}}\right]\right]\right) \geq 2\Phi(t_\alpha) - 1 = 1 - \alpha$$

★ Classique ! ★

C'est très classique (et le seul cas qui semble être au programme pour cette partie du cours).

✗ Attention !

L'intervalle de confiance est donné par deux estimateurs, qui ne peuvent pas dépendre du paramètre  $m$  recherché ! Autrement dit, l'expression des bornes de l'intervalle de confiance ne doit pas faire apparaître  $m$  !

ES Pour info...

Avec  $\alpha = 0,05$ , on a :  $\Phi^{-1}\left(1 - \frac{\alpha}{2}\right) \simeq 1,96$ .

Petite remarque

L'intervalle de confiance asymptotique proposé est centré en  $\bar{X}_n$ , d'amplitude  $2t_\alpha \frac{s}{\sqrt{n}}$ .

ES Pour info...

De façon générale,  $\sigma$  est inconnu et on peut alors utiliser un estimateur de l'écart-type dans l'intervalle de confiance asymptotique...

EXEMPLE 10

On considère une variable aléatoire suivant une loi de Bernoulli de paramètre  $p$  inconnu ainsi qu'un  $n$ -échantillon  $(X_1, \dots, X_n)$  de  $X$ . On note  $\bar{X}_n$  la moyenne empirique de cet échantillon.

- Justifions que :

$$\forall x \geq 0, \lim_{n \rightarrow +\infty} \mathbb{P}\left(\left[-x \leq \sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \leq x\right]\right) = 2\Phi(x) - 1$$

où  $\Phi$  est la fonction de répartition d'une variable aléatoire suivant la loi  $\mathcal{N}(0; 1)$ .

On sait que :

- ✓  $(X_k)_{k \in \mathbb{N}^*}$  est une suite de variables aléatoires indépendantes, de même loi, admettant la même espérance  $p$  et même variance  $p(1-p)$  non nulle ;

✓ pour tout  $n \in \mathbb{N}^*$ ,  $\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k$  et

$$\begin{aligned} \bar{X}_n^* &= \frac{\bar{X}_n - \mathbb{E}(\bar{X}_n)}{\sigma(\bar{X}_n)} \\ &= \sqrt{n} \frac{\bar{X}_n - p}{\sigma} \\ &= \sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \end{aligned}$$

Ainsi, d'après le théorème central limite :

$$\bar{X}_n^* \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} Z$$

où  $Z$  est une variable aléatoire suivant la loi  $\mathcal{N}(0; 1)$ .

Par conséquent, pour tout  $x \in \mathbb{R}^+$

$$\begin{aligned} \lim_{n \rightarrow +\infty} \mathbb{P}([-x \leq \bar{X}_n^* \leq x]) &= \mathbb{P}([-x \leq Z \leq x]) \\ &= \Phi(x) - \Phi(-x) \\ &= 2\Phi(x) - 1 \end{aligned} \quad \left. \begin{array}{l} \left. \begin{array}{l} \left. \begin{array}{l} Z \mapsto \mathcal{N}(0; 1), \text{ donc } Z \text{ est à densité} \\ \Phi(-x) = 1 - \Phi(x) \end{array} \right\} \right. \end{array} \right\}$$

**Conclusion :**  $\forall x \geq 0, \lim_{n \rightarrow +\infty} \mathbb{P} \left( \left[ -x \leq \sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \leq x \right] \right) = 2\Phi(x) - 1.$

- Soit  $\alpha \in ]0; 1[$ . Démontrons l'existence d'un unique réel  $t_\alpha$  tel que  $\Phi(t_\alpha) = 1 - \frac{\alpha}{2}$ . On donne :  $t_{0,05} \simeq 1,96$ . Voir exemple 8.

- Déduisons-en un intervalle de confiance asymptotique de  $p$  au niveau de confiance 95%.

✓ D'après ce qui précède avec  $x = t_{0,05} \geq 0$ , on a :

$$\lim_{n \rightarrow +\infty} \mathbb{P} \left( \left[ -t_{0,05} \leq \sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \leq t_{0,05} \right] \right) = 2\Phi(t_{0,05}) - 1 = 1 - 0,05 = 0,95$$

Or :

$$\begin{aligned} \mathbb{P} \left( \left[ -t_{0,05} \leq \sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \leq t_{0,05} \right] \right) &= \mathbb{P} \left( \left[ \bar{X}_n - t_{0,05} \frac{\sqrt{p(1-p)}}{\sqrt{n}} \leq p \leq \bar{X}_n + t_{0,05} \frac{\sqrt{p(1-p)}}{\sqrt{n}} \right] \right) \\ &\leq \mathbb{P} \left( \left[ \bar{X}_n - t_{0,05} \frac{1}{2\sqrt{n}} \leq p \leq \bar{X}_n + t_{0,05} \frac{1}{2\sqrt{n}} \right] \right) \\ &\leq \mathbb{P} \left( \left[ \bar{X}_n - \frac{1}{\sqrt{n}} \leq p \leq \bar{X}_n + \frac{1}{\sqrt{n}} \right] \right) \end{aligned}$$

↳  $p(1-p) \leq \frac{1}{4}$ , d'où l'inclusion d'évènements et par croissance de  $\mathbb{P} \dots$   
↳  $t_{0,05} \leq 2$  et même raisonnement

On obtient donc :

$$\lim_{n \rightarrow +\infty} \mathbb{P} \left( \left[ \bar{X}_n - \frac{1}{\sqrt{n}} \leq p \leq \bar{X}_n + \frac{1}{\sqrt{n}} \right] \right) \geq 0,95$$

- ✓ Enfin,  $\bar{X}_n - \frac{1}{\sqrt{n}}$  et  $\bar{X}_n + \frac{1}{\sqrt{n}}$  sont des estimateurs de  $p$  (car fonctions d'un  $n$ -échantillon de  $X$  dont les expressions ne font pas apparaître  $p$ ) et  $\mathbb{P} \left( \left[ \bar{X}_n - \frac{1}{\sqrt{n}} \leq \bar{X}_n + \frac{1}{\sqrt{n}} \right] \right) = 1.$

**Conclusion :** l'intervalle  $\left[ \bar{X}_n - \frac{1}{\sqrt{n}}; \bar{X}_n + \frac{1}{\sqrt{n}} \right]$  est un intervalle de confiance asymptotique de  $p$  au niveau de confiance 95%.

- Quelle taille d'échantillon devons-nous avoir pour que l'intervalle de confiance asymptotique donné ait une amplitude inférieure ou égale à 0,1 ?

L'amplitude de  $\left[ \bar{X}_n - \frac{1}{\sqrt{n}}; \bar{X}_n + \frac{1}{\sqrt{n}} \right]$  est égale à  $\frac{2}{\sqrt{n}}$ . Or, pour tout  $n \in \mathbb{N}^*$  :

$$\frac{2}{\sqrt{n}} \leq 0,1 \iff n \geq 400$$

**Conclusion :** il faudrait un échantillon d'au moins 400 individus pour que l'amplitude de cet intervalle soit inférieure ou égale à 0,1.

- Comparons cet intervalle avec celui obtenu dans l'exemple 9.

\* Puisque  $\sqrt{5} > 1$ , on a :

$$\left[ \bar{X}_n - \frac{1}{\sqrt{n}}; \bar{X}_n + \frac{1}{\sqrt{n}} \right] \subset \left[ \bar{X}_n - \sqrt{\frac{5}{n}}; \bar{X}_n + \sqrt{\frac{5}{n}} \right]$$

\* En comparaison, il fallait un échantillon de taille au moins 2000 pour obtenir une amplitude inférieure ou égale à 0,1 avec l'intervalle de confiance exact obtenu dans l'exemple 9.

**Petite remarque**  
Quand  $n = 100$ , on obtient un intervalle de confiance de la forme  $[\bar{X}_{100} - 0,1; \bar{X}_{100} + 0,1]$  qui est relativement grand pour estimer un paramètre  $p \in ]0; 1[$ .