

*If there is a 50-50 chance that something can go wrong,
then nine times out of 10 it will.*

PAUL HARVEY

Animateur radio américain (1918-2009)

1 Estimation ponctuelle

1.1 Principe : modéliser, estimer, tester/prédire

On considère un phénomène aléatoire et on s'intéresse à une variable aléatoire réelle X qui pourrait le décrire. On suppose que la loi de probabilité de X n'est pas complètement spécifiée et appartient à une famille de lois dépendant d'un paramètre θ décrivant un ensemble Θ . Le paramètre θ est une quantité inconnue, fixée dans toute l'étude, que l'on cherche à déterminer ou pour laquelle on cherche une information partielle.

Le problème de l'*estimation ponctuelle* consiste alors à préciser la vraie valeur du paramètre θ (ou plus généralement d'une fonction $g(\theta)$) à partir d'un échantillon de données x_1, \dots, x_n obtenues en observant n fois le phénomène. Cette fonction du paramètre représentera en général une valeur caractéristique de la loi inconnue comme son espérance, sa variance, son étendue¹, etc.

Exemples.

- *Exemple 1. Nombre de buts en Ligue 1.*

→ Voici le nombre de buts par journée lors de la saison 2021/2022 :

26, 34, 29, 31, 28, 25, 32, 30, 25, 26, 29, 29, 29, 30, 21, 29, 29, 27,

27, 15, 26, 35, 30, 23, 23, 22, 21, 19, 23, 32, 38, 30, 26, 35, 30, 36, 30, 37.

Ces nombres constituent l'*échantillon de données*.

→ On pose un *modèle* en supposant que le nombre de buts lors d'une journée est une variable aléatoire X qui suit une loi de Poisson. Notons ce paramètre λ .

→ Pour donner une valeur à λ , (on dira *estimer*), on peut utiliser le fait que X admet une espérance $E(X) = \lambda$. Ainsi, on peut préciser λ en prenant la valeur moyenne du nombre de buts (à savoir ici $\lambda = 2.81$).

→ Une fois le modèle posé, on peut le *tester* s'il donne des résultats cohérents et ainsi l'utiliser pour faire de la *prédiction*.

- *Exemple 2. Durée de vie d'un composant électrique.*

On suppose que l'emploi de ce composant se fait dans des conditions normales d'utilisation et on néglige les phénomènes d'usure. Autrement dit, on suppose que la variable aléatoire X qui donne la durée de vie du composant (en

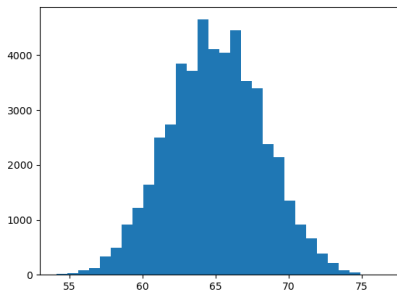
1. Si $X(\Omega) = [a; b]$, l'étendue est définie par la quantité $b - a$.

heure) est une loi sans mémoire. On a vu (voir chapitre LOIS USUELLES À DENSITÉ) que X suit alors une loi exponentielle. De plus, on appelle demi-vie de X le réel t tel que $\mathbf{P}(X \leq t) = \mathbf{P}(X \geq t)$ (on parle aussi de médiane). On vérifie que que $t = \ln 2 / \lambda$.

On teste une centaine de composants et au bout de 1000 heures d'utilisation, la moitié des composants ne fonctionnent plus. On fait alors le choix de λ tel que $1000 = \ln 2 / \lambda$, soit $\lambda = \ln(2) / 1000 \approx 6.9 \cdot 10^{-4}$.

• *Exemple 3.*

Vincent, éleveur à Armentières sur Ourcq (02) produit des œufs de poules élevées en plein air. Il souhaite anticiper sa production de gros œufs pour l'année 2023 en analysant sa production actuelle. Un très gros œuf a un poids supérieur ou égal à 73 grammes. Voici la répartition de la production des 51 243 œufs de l'année.



Console

```
>>> len(L)
51243

>>> np.mean(L) # poids moyen d'un oeuf
65.115928

>>> np.std(L) # écart-type
3.2336541402592824
```

Vincent suppose que le poids moyen d'un œuf correspond à une variable aléatoire X qui a une loi normale de paramètre (μ, σ^2) avec $\mu = 65.1$ et $\sigma = 3.23$.

La suite de ce chapitre propose de donner un cadre théorique rigoureux à ces trois exemples et de juger aussi de la pertinence de cette approche. Nous discuterons essentiellement de l'étape d'estimation. Notamment nous verrons deux types d'estimations :

- L'estimation ponctuelle .
- L'estimation par intervalle de confiance.

1.2 Définitions et exemples

Dans la suite, on se fixe :

- un espace probabilisable (Ω, \mathcal{A}) .
- un espace des paramètres Θ qui est une partie de \mathbb{R}^n .
On suppose de plus que pour chaque paramètre $\theta \in \Theta$, il existe une probabilité \mathbf{P}_θ définie sur (Ω, \mathcal{A}) .
- une application X qui est bien une variable aléatoire sur tous les espaces probabilisés $(\Omega, \mathcal{A}, \mathbf{P}_\theta)$ (où $\theta \in \Theta$).

Notation. Sous réserve d'existence :

- $\mathbf{E}_\theta(X)$ désigne l'espérance de X pour la probabilité \mathbf{P}_θ .
- $\mathbf{V}_\theta(X)$ est la variance de X pour la probabilité \mathbf{P}_θ .

Exemples.

• *Exemple 1.*

Dans ce cas, la variable X suit une loi de Poisson de paramètre θ . On a $\Theta = \mathbb{R}_*^+$ et

$$\forall \theta \in \mathbb{R}_*^+, \quad \forall n \in \mathbb{N}, \quad \mathbf{P}_\theta(X = n) = e^{-\theta} \frac{\theta^n}{n!}.$$

On a de plus, $\mathbf{E}_\theta(X_i) = \theta$ et $\mathbf{V}_\theta(X_i) = \theta$.

- *Exemple 2.* Dans le cas où X suit une loi exponentielle de paramètre $\lambda = \theta$, on a

$$\forall \theta \in \mathbb{R}_+^*, \quad \forall t \in \mathbb{R}_+^*, \quad \mathbf{P}_\theta(X \leq t) = 1 - e^{-\theta t}.$$

- *Exemple 3.* Dans ce dernier cas, X suit une loi normale $\mathcal{N}(\mu, \sigma^2)$ dont on ne connaît ni l'espérance ni la variance, alors $\theta = (\mu, \sigma)$, et $\Theta = \mathbb{R} \times \mathbb{R}_+^*$.

DÉFINITION

échantillon

Soient $X: \Omega \rightarrow \mathbb{R}$ une variable aléatoire définie sur (Ω, \mathcal{A}) et $n \in \mathbb{N}^*$.

On appelle **n -échantillon** de la loi de X toute famille (X_1, \dots, X_n) telle que :

- Les applications X_1, \dots, X_n sont des variables aléatoires sur (Ω, \mathcal{A}) .
- Les variables X_1, \dots, X_n sont \mathbf{P}_θ -indépendantes et de même loi que X pour tout $\theta \in \Theta$.

Vocabulaire.

- On dit aussi que la loi de X est la loi parente (ou encore loi mère) de l'échantillon.
- On note souvent (X_1, \dots, X_n) *i.i.d* pour signaler que les variables sont indépendantes, et identiquement distribuées (c'est-à-dire de même loi).

En pratique, un échantillon de données x_1, \dots, x_n est la réalisation de n variables aléatoires X_1, \dots, X_n . L'objectif de l'estimation ponctuelle est alors de déterminer le paramètre θ (ou une fonction $g(\theta)$) qui « explique » au mieux les valeurs de l'échantillon.

DÉFINITION

estimateur

- On appelle **estimateur de θ** toute variable aléatoire de la forme $\varphi(X_1, \dots, X_n)$, où (X_1, \dots, X_n) est un n -échantillon et φ , une fonction de \mathbb{R}^n dans \mathbb{R} .
- Plus généralement, pour $g: \Theta \rightarrow \mathbb{R}$, une fonction, un **estimateur de $g(\theta)$** est une variable aléatoire de la forme $\varphi(X_1, \dots, X_n)$ où (X_1, \dots, X_n) est un n -échantillon.

Remarques.

- Un estimateur ne peut pas dépendre de θ puisque c'est la valeur que l'on souhaite déterminer.
- Estimer ponctuellement $g(\theta)$ par $\varphi(x_1, \dots, x_n)$ où $\varphi(X_1, X_2, \dots, X_n)$ est un estimateur de $g(\theta)$ et (x_1, \dots, x_n) est une réalisation de l'échantillon (X_1, \dots, X_n) , c'est décider d'accorder à $g(\theta)$ la valeur $\varphi(x_1, \dots, x_n)$.

Exemples.

- Avec les notations de la définition, soit (X_1, \dots, X_n) un n -échantillon de la loi de X . Il est souvent utile de considérer l'**estimateur de la moyenne empirique** donné par :

$$\overline{X}_n = \frac{X_1 + X_2 + \dots + X_n}{n} \quad (\bullet)$$

Mais on peut inventer toute une gamme d'estimateurs :

$$T_n = \max_{1 \leq k \leq n} X_k, \quad U_n = \min_{1 \leq k \leq n} X_k, \quad V_n = \ln(1 + |U_n|), \quad W_n = \max\{k \in [1, n] \mid X_k = T_n\}, \quad Y_n = 1, \quad \text{etc.}$$

```

theta=np.random.rand()
theta=0.4292081919392057 # le paramètre "inconnu"

n=500 # taille de l'échantillon
Ech=theta*np.random.rand(n)
# création de l'échantillon

np.round(Ech,2)
array([[0.18, 0.25, 0.28, 0.26, 0.05, 0.14, 0.15, 0.02,
        0.15, 0.35, 0.2 ,
        0.25, 0.41, 0.17, 0.39, 0.32, 0.22, 0.21, 0.12,
        0.14, 0.35, 0.18,
        0.16, 0.23, 0.1 , 0.19, 0.28, 0.3 , 0.31, 0.23, 0.4
        , 0.08, 0.18 ...

```

- Créons un échantillon d'une loi uniforme $[0; \theta]$.

On peut essayer de retrouver la valeur de θ à partir de deux estimateurs

$$2\overline{X}_n \quad \text{et} \quad T_n = \max_{1 \leq k \leq n} X_k.$$

On obtient deux estimations de θ par :

```
>>> 2*sum(Ech)/500
0.4386345472617729
```

```
>>> max(Ech)
0.4289769082767542
```

1.3 Biais, convergence et comparaison des estimateurs

Les définitions qui suivent permettent de quantifier la « qualité » d'un estimateur et de les comparer entre-eux.

DÉFINITION

biais d'un estimateur

Soit T_n un estimateur de $g(\theta)$ tel que tout $\theta \in \Theta$, T_n admet une espérance pour la probabilité \mathbf{P}_θ .

- On définit le **biais** de T_n en $g(\theta)$ par

$$b_\theta(T_n) = \mathbf{E}_\theta(T_n) - g(\theta).$$

- Si pour tout $\theta \in \Theta$, $b_\theta(T_n) = 0$, on dit l'estimateur est **sans biais**. Sinon, on dit que l'estimateur est **biaisé**.

Exemples.

- Avec les notations de la définition et si X admet une espérance θ , l'estimateur de la moyenne empirique \overline{X}_n est un estimateur sans biais de θ .

En effet, pour tout $i \in \llbracket 1, n \rrbracket$, on a $\mathbf{E}_\theta(X_i) = \theta$. Par linéarité de l'espérance \mathbf{E}_θ , il vient

$$\mathbf{E}_\theta(\overline{X}_n) = \frac{1}{n} \sum_{i=1}^n \mathbf{E}_\theta(X_i) = \theta, \quad \text{puis} \quad b_\theta(\overline{X}_n) = \mathbf{E}_\theta(\overline{X}_n) - \theta = 0.$$

- Considérons une variable X suivant une loi uniforme sur $[0, \theta]$. Pour tout $n \in \mathbb{N}^*$, posons $M_n = \max(X_1, \dots, X_n)$. On vérifie que M_n est une variable à densité, et une densité est donnée par

$$\forall x \in \mathbb{R}, \quad f_\theta(x) = \begin{cases} n \frac{x^{n-1}}{\theta^n} & \text{si } x \in [0; \theta] \\ 0 & \text{sinon.} \end{cases}$$

En tant que variable bornée, M_n admet une espérance avec

$$\mathbf{E}_\theta(M_n) = \int_0^\theta x f_\theta(x) dx = \int_0^\theta n \frac{x^n}{\theta^n} dx = \frac{n}{n+1} \theta.$$

On en déduit que M_n est un estimateur biaisé de θ et le biais de M_n est

$$b_\theta(M_n) = \frac{n}{n+1} \theta - \theta = -\frac{\theta}{n+1}.$$

Par contre, par linéarité de l'espérance, $\tilde{M}_n = \frac{n+1}{n} M_n$ est un estimateur sans biais de θ puisque

$$\mathbf{E}_\theta(\tilde{M}_n) = \frac{n+1}{n} \mathbf{E}_\theta(M_n) = \theta.$$

Remarque. On peut donner une définition moins contraignante : un estimateur est **asymptotiquement sans biais** si

$$b_{\theta}(T_n) \xrightarrow{n \rightarrow \infty} 0.$$

Par exemple, la variable M_n définie précédemment est un estimateur asymptotiquement sans biais.

Exercice 1



◆◆ Soit (X_1, \dots, X_n) un n -échantillon de la loi de Bernoulli de paramètre p . On pose

$$S_n = \sum_{i=1}^n X_i \quad \text{et} \quad T_n = \frac{S_n}{n} \left(1 - \frac{S_n}{n}\right).$$

p. ??

1. Déterminer $E(S_n)$ et $E(S_n^2)$. En déduire $E(T_n)$.
2. À l'aide de T_n , proposer un estimateur sans biais de la variance de cette loi de Bernoulli.

Exercice 2



◆◆ **Estimateur de la variance**

1. On suppose, dans cette question, que la variable X admet un espérance μ connue et une variance σ^2 , inconnue.

Montrer que la variable $T_n = \frac{1}{n} \sum_{k=1}^n (X_k - \mu)^2$ est un estimateur sans biais de σ^2 .

p. 24

2. On suppose, dans cette question, que μ est inconnu. On note $V_n = \frac{1}{n} \sum_{k=1}^n (X_k - \bar{X}_n)^2$.

- a) Montrer que V_n est un estimateur asymptotiquement sans biais de σ^2 et calculer le biais de cet estimateur.
- b) Construire, à partir de V_n , un estimateur sans biais de σ^2 .

DÉFINITION

estimateur convergent

Soit $(T_n)_{n \in \mathbb{N}^*}$, une suite d'estimateurs de $g(\theta)$.

On dit que la suite $(T_n)_{n \in \mathbb{N}}$ est **convergente** si pour tout $\theta \in \Theta$, la suite (T_n) converge en probabilité vers la variable aléatoire presque sûrement constante $g(\theta)$. Autrement dit

$$\forall \theta \in \Theta, \quad \forall \varepsilon > 0, \quad \mathbf{P}_{\theta}(|T_n - g(\theta)| \geq \varepsilon) \xrightarrow{n \rightarrow \infty} 0.$$

Vocabulaire. Par abus de langage, on dit simplement que T_n est un estimateur convergent de $g(\theta)$.

Exercice 3



◆ **Exemples**

1. Considérons X une variable aléatoire dont la loi est uniforme sur $[0, \theta]$, où le paramètre θ est inconnu. Vérifier que $T_n = 2\bar{X}_n$ est un estimateur convergent de θ .

2. Soit X une variable à densité donnée pour $\theta \in \mathbb{R}_*^+$ par

$$\forall x \in \mathbb{R}, \quad f_{\theta}(x) = \frac{2}{\theta} \left(1 - \frac{x}{\theta}\right) \mathbf{1}_{[0, \theta]}(x).$$

p. 24

- a) Pour un échantillon (X_1, \dots, X_n) , montrer que $3\bar{X}_n$ est un estimateur sans biais de θ .
- b) Vérifier que cet estimateur est convergent.

Remarque. La notion de convergence des estimateurs ne donne aucune assurance pratique que la valeur prise par un estimateur à partir de l'échantillon de données sera assez proche de la vraie valeur du paramètre. On quantifie la qualité des estimateurs par la notion de **risque quadratique**. Cette notion est maintenant hors-programme mais on pourra consulter l'exercice 11, page 15, pour des précisions.

PROPOSITION

Soit $(T_n)_{n \in \mathbb{N}}$, une suite d'estimateurs de $g(\theta)$.

- Si** | \rightarrow La suite $(T_n)_{n \in \mathbb{N}}$ est une suite d'estimateurs convergente de $g(\theta)$.
 | \rightarrow La fonction f est continue sur \mathbb{R} .

Alors $(f(T_n))_{n \in \mathbb{N}}$ est une suite d'estimateurs convergente de $f(g(\theta))$.

Preuve. C'est une conséquence directe de l'énoncé : Soit $(X_n)_{n \in \mathbb{N}}$ une suite de variables aléatoires sur $(\Omega, \mathcal{A}, \mathbf{P})$.

- Si** | \rightarrow La suite $(X_n)_{n \in \mathbb{N}}$ converge en probabilité vers X .
 | \rightarrow La fonction f est continue sur \mathbb{R} à valeurs réelles.

Alors $f(X_n) \xrightarrow[n \rightarrow +\infty]{\mathbf{P}} f(X)$.

Exemples. Nous en donnons deux.

- Reprenons le cas de X une variable aléatoire dont la loi est uniforme sur $[0; \theta]$ avec θ inconnu. À partir du théorème de transfert, on vérifie que $\ln(X)$ a une espérance avec

$$\mathbf{E}_\theta(\ln(X)) = \int_0^\theta \ln(t) dt = \ln(\theta) - 1.$$

On vérifie de même que $\ln(X)$ admet un moment d'ordre 2 (et donc une variance). Toujours d'après la loi faible des grands nombres, $\left(\frac{\sum_{i=1}^n \ln(X_i)}{n}\right) / n$ est un estimateur convergent de $\ln(\theta) - 1$. La fonction $f : t \in \mathbb{R} \mapsto \exp(t + 1)$ est continue, donc l'estimateur suivant est un estimateur convergent de θ :

$$e \cdot \sqrt[n]{\prod_{i=1}^n X_i} = \exp\left(\frac{\sum_{i=1}^n \ln(X_i)}{n} + 1\right) = f\left(\frac{1}{n} \sum_{i=1}^n \ln(X_i)\right) \xrightarrow[n \rightarrow +\infty]{\mathbf{P}} f(\ln(\theta) - 1) = \theta.$$

- Soit $(X_n)_{n \in \mathbb{N}^*}$ des variables aléatoires indépendantes de même loi géométrique $\mathcal{G}(p)$. Comme la loi géométrique est d'espérance $1/p$ (et admettant une variance), la loi faible justifie que $\frac{\sum_{i=1}^n X_i}{n}$ est un estimateur de $1/p$. Comme la fonction inverse est continue sur \mathbb{R}_*^+ , $n / \left(\frac{\sum_{i=1}^n X_i}{n}\right)$ est un estimateur convergent de p .

Notons que la conclusion du théorème demeure si f est continue sur un intervalle I contenant $X(\Omega)$.

Exercice 4

◆ Soit (X_1, X_2, \dots, X_n) un n -échantillon d'une variable aléatoire X suivant une loi normale $\mathcal{N}(0, \sigma^2)$.

1. Pour tout i , calculer l'espérance de la variable aléatoire $|X_i|$.
2. En déduire un estimateur sans biais et convergent de σ .

p. 25

Exercice 5**◆◆ Biases et composition**

Soit (X_1, X_2, \dots, X_n) un n -échantillon d'une variable aléatoire X suivant une loi de Poisson de paramètre λ inconnu. On sait que la moyenne empirique $\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k$ est un estimateur sans biais et convergent de λ . On cherche à estimer $e^{-\lambda}$.

Est-ce que l'estimateur $T_n = e^{-\bar{X}_n}$ est un estimateur sans biais de $e^{-\lambda}$? asymptotiquement sans biais? convergent?

p. 25

PROPOSITION

condition suffisante de convergence

Soit $(T_n)_{n \in \mathbb{N}}$, une suite d'estimateurs de $g(\theta)$.Si $\mathbf{E}(T_n) \xrightarrow[n \rightarrow \infty]{} g(\theta)$ et $\mathbf{V}(T_n) \xrightarrow[n \rightarrow \infty]{} 0$.Alors $(T_n)_{n \in \mathbb{N}}$ est une suite d'estimateurs convergente de $g(\theta)$.**Exercice 6**

◇ Prouver cet énoncé.

p. 25

Remarque. Un estimateur sans biais Y_n est meilleur qu'un autre estimateur sans biais Z_n si $\mathbf{V}(Y_n) \leq \mathbf{V}(Z_n)$ pour tout entier n . Le fait d'être meilleur se matérialise dans l'inégalité de Bienaymé-Tchebychev :

$$\forall \varepsilon \in \mathbb{R}_*^+, \quad \mathbf{P}_\theta (|T_n - g(\theta)| \geq \varepsilon) \leq \frac{\mathbf{V}(Y_n)}{\varepsilon^2} \leq \frac{\mathbf{V}(Z_n)}{\varepsilon^2}.$$

Exemple. L'estimateur de la moyenne empirique est convergent si X admet une variance.**Exercice 7**◇ Soit $(X_n)_{n \geq 1}$ une suite de variables aléatoires mutuellement indépendantes et suivant toutes la loi $\mathcal{B}(p)$, où p est un paramètre inconnu que l'on cherche à estimer. On pose

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{et} \quad T_n = \frac{2}{n(n+1)} \sum_{i=1}^n i X_i.$$

p. 25

1. Montrer que \bar{X}_n et T_n sont deux estimateurs sans biais de p .
2. Calculer et comparer les variances de \bar{X}_n et de T_n .
3. Montrer que \bar{X}_n et T_n sont deux estimateurs convergents de p .

Exemple. Questions sensibles lors d'un sondage d'opinion.Certains sujets abordés dans les sondages d'opinion peuvent être sensibles et les personnes interrogées peuvent refuser de répondre honnêtement. Détaillons une procédure qui permet aux sondés de répondre plus librement. Considérons n personnes sondées et une question fermée à deux réponses possibles dont on veut estimer la probabilité p de réponses positives dans la population générale. On demande à chaque sondé de lancer un dé.

- S'il obtient 6, la personne doit donner sa réponse sans mentir.
- Sinon, elle donne la réponse contraire à la sienne.

Si le sondeur ignore le résultat du dé, il ne pourra pas savoir si la réponse est franche ou non, et on peut espérer que la personne sondée acceptera plus facilement de répondre honnêtement à la question.

Généralisons la procédure en fixant t , la probabilité que la personne réponde sans mentir. Le réel t est connu, il vaut $1/6$ dans l'exemple du dé. Posons pour tout $i \in \llbracket 1; n \rrbracket$, la variable aléatoire X_i valant 1 le i -ème sondé répond positivement et 0 sinon. X_i suit une loi de Bernoulli. Donnons son paramètre à l'aide de la formule de probabilités totales avec le système complet d'événements constitué de A : « Le i sondé ne ment pas » et \bar{A} .

$$\mathbf{P}(X_i = 1) = \mathbf{P}(A)\mathbf{P}_A(X_i = 1) + \mathbf{P}(\bar{A})\mathbf{P}_{\bar{A}}(X_i = 1) = tp + (1-t)(1-p).$$

Nous avons vu que l'estimateur de la moyenne empirique $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ est un estimateur sans biais et convergent de $r := tp + (1-t)(1-p)$. Or, pour $t \neq 1/2$, on peut inverser la relation

$$r = tp + (1-t)(1-p) \iff p = \frac{1-r-t}{2t-1} \iff p = f(r) \quad \text{où} \quad f: x \in \mathbb{R} \mapsto \frac{1-x-t}{2t-1}.$$

Posons alors :

$$T_n = f(\bar{X}_n) = \frac{1-t-\bar{X}_n}{1-2t}.$$

 T_n est alors est estimateur convergent de p . Donnons deux arguments :

- La fonction f est continue et \bar{X}_n est un estimateur convergent de r . Donc $T_n = f(\bar{X}_n)$ est un estimateur convergent de $f(r) = p$.
- Par linéarité de l'espérance, on vérifie que l'espérance de T_n est p . T_n est un estimateur sans biais de p . De plus, la variance de T_n vaut :

$$V(T_n) = \frac{r(1-r)}{n(2t-1)^2} \xrightarrow{n \rightarrow \infty} 0.$$

D'après la proposition précédente, on retrouve le fait que l'estimateur T_n est convergent de p .

Comparaison des estimateurs sans biais

Simulation Python.

Reprenons le cas d'une variable X suivant une loi uniforme sur $[0; \theta]$ et des deux estimateurs sans biais :

$$\bar{X}_n = \frac{2}{n} \sum_{i=1}^n X_i \quad \text{et} \quad \tilde{M}_n = \frac{n+1}{n} \max(X_1, \dots, X_n).$$

Simulons un grand nombre d'échantillons de données et affichons les réalisations de ces deux estimateurs à l'aide des histogrammes.

Editeur

```
# Importation des bibliothèques
# et définition des deux estimateurs

import numpy as np
import random as rd
import matplotlib.pyplot as plt

theta=0.4292081919392057
# choix du paramètre "inconnu"

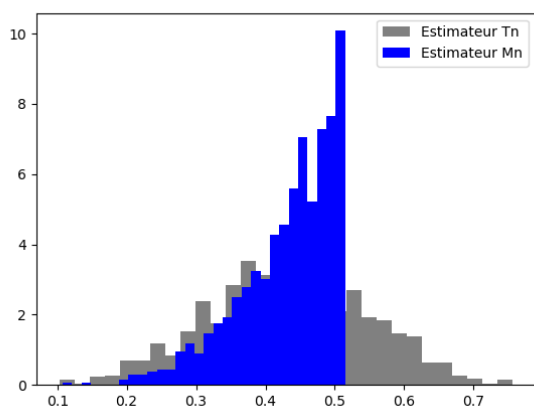
def estimateurT(n):
    return 2*np.sum(theta*np.random.rand(n))/n

def estimateurM(n):
    return ((n+1)/n)*max(theta*np.random.rand(n))

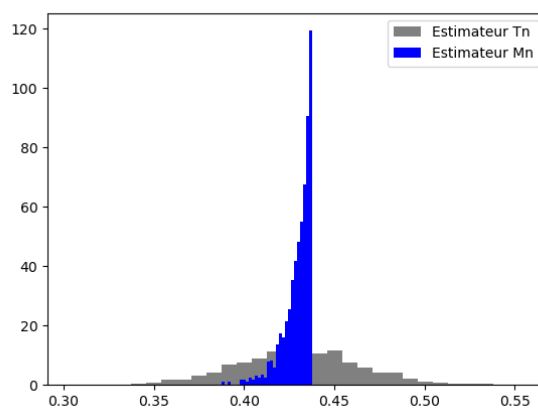
# affichage des histogrammes

def histogrammes(n):
    LM=[]
    LT=[]
    for i in range(1000):
        LM.append(estimateurM(n))
        LT.append(estimateurT(n))
    plt.clf()
    plt.hist(LM,30)
    plt.hist(LT,30)
    plt.show()
```

>>> histogrammes(5)



>>> histogrammes(50)

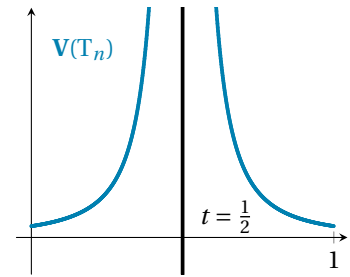


Le meilleur estimateur sans biais est celui qui a la variance la plus faible. En effet, c'est celui où le risque que l'échantillon de données donnent une valeur loin de l'espérance (qui est dans ce cas le paramètre à estimer).

Exemple. On a vu que dans le cas de questions sensibles lors d'un sondage d'opinion, la variance de l'estimateur est

$$V(T_n) = \frac{r(1-r)}{n(2t-1)^2}.$$

On retrouve bien le cas que l'estimateur est meilleur lorsque $t = 0$ (le sondé ne ment jamais) ou $t = 1$ (le sondé ment systématiquement).



2

Estimation par intervalle de confiance

S'il existe des critères pour juger des qualités d'un estimateur ponctuel T_n de $g(\theta)$, aucune certitude ne peut jamais être apportée quant au fait que l'estimation donnée par l'échantillon de données soit une « bonne » valeur du paramètre $g(\theta)$. L'estimation par intervalle de confiance permet de trouver un intervalle aléatoire qui contienne $g(\theta)$ avec une probabilité minimale donnée. Dans tout ce paragraphe, $(U_n)_{n \in \mathbb{N}^*}$ et $(V_n)_{n \in \mathbb{N}^*}$ désigneront deux suites d'estimateurs de $g(\theta)$ telles que pour tous $\theta \in \Theta$ et $n \in \mathbb{N}^*$,

$$\mathbf{P}_\theta ([U_n \leq V_n]) = 1.$$

2.1 Intervalle de confiance

Définition

DÉFINITION

intervalle de confiance

Soient $\alpha \in]0, 1[$, U_n et V_n deux estimateurs de $g(\theta)$ tels que pour tout $\theta \in \Theta$

$$\mathbf{P}_\theta (U_n \leq g(\theta) \leq V_n) \geq 1 - \alpha.$$

On dit que l'intervalle $[U_n, V_n]$ est un **intervalle de confiance** de $g(\theta)$ avec un risque d'au plus α ou au niveau de confiance au moins égal à $1 - \alpha$.

Remarque. En pratique, on part d'un échantillon de données x_1, x_2, \dots, x_n . On en déduit les valeurs

$$u_n = U_n(x_1, x_2, \dots, x_n) \quad \text{et} \quad v_n = V_n(x_1, x_2, \dots, x_n).$$

Ainsi on construit un intervalle aléatoire $[u_n, v_n]$ dans lequel $g(\theta)$ à une probabilité supérieure à $1 - \alpha$ de s'y trouver.

Remarque. Dans les conditions usuelles, on considère des niveaux de confiance de 95% ou 99% (soit $\alpha = 0,05$ ou $\alpha = 0,01$).

Estimation par intervalle de confiance en utilisant l'inégalité de Bienaymé-Tchebychev

Méthode

Comment obtenir un intervalle de confiance à partir de l'inégalité de Bienaymé-Tchebychev?

Soit T_n , une variable aléatoire admettant une variance, l'inégalité s'écrit

$$\mathbf{P}_\theta \left(|T_n - \mathbf{E}_\theta(T_n)| \geq \varepsilon \right) \leq \frac{\mathbf{V}_\theta(T_n)}{\varepsilon^2}.$$

Comme $\left[|T_n - \mathbf{E}_\theta(T_n)| < \varepsilon \right] \subset \left[|T_n - \mathbf{E}_\theta(T_n)| \leq \varepsilon \right]$, on obtient par passage au complémentaire

$$\mathbf{P}_\theta \left(|T_n - \mathbf{E}_\theta(T_n)| \leq \varepsilon \right) \geq 1 - \frac{\mathbf{V}_\theta(T_n)}{\varepsilon^2}.$$

Si T_n est un estimateur sans biais de $g(\theta)$ (c'est-à-dire $\mathbf{E}_\theta(T_n) = g(\theta)$) et si on peut trouver un entier n tel que $\frac{\mathbf{V}_\theta(T_n)}{\varepsilon^2} \leq \alpha$, on peut récrire l'inégalité

$$\mathbf{P}_\theta \left(T_n - \varepsilon \leq g(\theta) \leq T_n + \varepsilon \right) \geq 1 - \alpha.$$

L'intervalle de confiance est alors $[T_n - \varepsilon; T_n + \varepsilon]$, il est de longueur 2ε .

Exemple. Estimation du paramètre p d'une loi de Bernoulli.

On réalise un sondage sur n personnes avec une unique question. On suppose que les réponses des personnes sont indépendantes et on veut déterminer un intervalle de confiance d'au moins 0.95 de la probabilité p de répondre positivement à l'unique question posée. D'après la loi faible des grands nombres, cet intervalle sera d'autant plus petit (au sens de l'inclusion) que le nombre de personnes interrogées sera grand. Désignons, pour un entier naturel non nul quelconque n , par X_n la variable aléatoire égale à 1 si la $n^{\text{ième}}$ personne répond positivement et 0 sinon. La variable X_n suit une loi de Bernoulli, de paramètre p . En particulier

$$\mathbf{E}(X_n) = p \quad \text{et} \quad \mathbf{V}(X_n) = p(1-p).$$

Soit \bar{X}_n la moyenne empirique

$$\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k, \quad \mathbf{E}(\bar{X}_n) = p \quad \text{et} \quad \mathbf{V}(\bar{X}_n) = \frac{p(1-p)}{n} \quad (\text{indépendances}).$$

D'après l'inégalité de Bienaymé-Tchebychev :

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

Par une étude de fonction, on peut majorer $p(1-p)$ par $\frac{1}{4}$. On obtient

$$\mathbf{P} \left(\left| \bar{X}_n - p \right| \leq \varepsilon \right) \geq 1 - \alpha \quad \text{et} \quad \alpha = \frac{1}{4n\varepsilon^2} \iff \varepsilon = \frac{1}{2\sqrt{n\alpha}}.$$

Dans ce cas, on obtient :

$$\mathbf{P} \left(p \in \left[\bar{X}_n - \frac{1}{2\sqrt{n\alpha}}, \bar{X}_n + \frac{1}{2\sqrt{n\alpha}} \right] \right) \geq 1 - \alpha$$

On obtient ainsi un intervalle de confiance de p à un niveau de confiance $1 - \alpha$ avec

$$\left[\bar{X}_n - \frac{1}{2\sqrt{n\alpha}}, \bar{X}_n + \frac{1}{2\sqrt{n\alpha}} \right].$$

Pour $\alpha = 0.05$, on a en particulier

$$\mathbf{P} \left(p \in \left[\bar{X}_n - \frac{1}{\sqrt{0.05n}}, \bar{X}_n + \frac{1}{\sqrt{0.05n}} \right] \right) \geq 0.95$$

et si on prend maintenant $n = 100$, on obtient $\varepsilon \approx 0.22$. Donc $[\bar{X}_n - 0.22, \bar{X}_n + 0.22]$ est un intervalle de confiance de p au niveau de risque 0.05. Autrement dit, il y a plus de 95% de chances que p soit compris entre $\bar{X}_n - 0.22$ et $\bar{X}_n + 0.22$. L'étendue de l'intervalle est énorme (0.44) lorsque qu'on se rappelle que $p \in [0; 1]$.

À l'inverse, on peut fixer ε , et s'intéresser à la taille de l'échantillon nécessaire pour garantir que l'étendue de l'intervalle $[\bar{X}_n - \varepsilon, \bar{X}_n + \varepsilon]$ soit inférieure à 1%, il faut

$$n \geq \frac{1}{4 \times 0.05 \times \varepsilon^2} = 50000.$$

On constate qu'il faut interroger 50 000 personnes!

Remarques.

- Notons aussi que pour diviser par 2 la longueur de l'intervalle, il faut quadrupler le nombre de sondés.
- On retrouve le fait que l'inégalité de Bienaymé-Tchebychev est trop générale pour donner des résultats précis : elle ne prend pas suffisamment en compte la loi de \bar{X}_n , seules son espérance et sa variance sont importantes. Nous verrons dans la suite comment construire des intervalles de confiance plus petits en utilisant plus finement la loi de \bar{X}_n .

Simulation Python. Le code suivant crée plusieurs séries de données et construit l'intervalle aléatoire associé à chaque série de données.

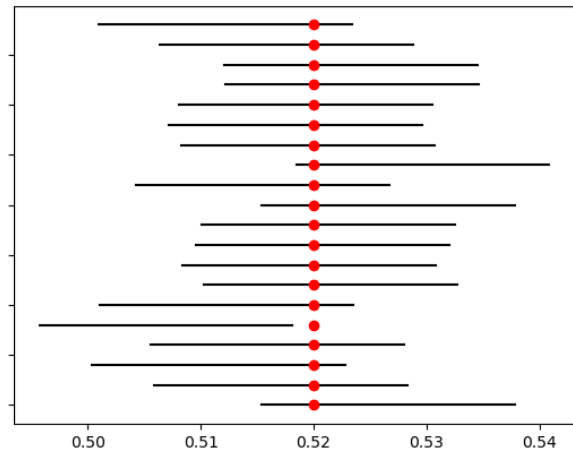
Editeur

```
import numpy as np
import matplotlib.pyplot as plt

p=0.52 # le paramètre "inconnu" à estimer

n=10000 # Nombre de sondés
alpha=0.20 # niveau de risque

NbreTest=20
for i in range(NbreTest):
    # Création de l'échantillon de données
    Ech=np.random.rand(n)<p
    # Calcul des bords de l'intervalle
    u=np.mean(Ech)-1/(2*(n*alpha)**(1/2))
    v=np.mean(Ech)+1/(2*(n*alpha)**(1/2))
    plt.plot([u,v],[i,i],'k')
    plt.plot([p],[i],'ro')
plt.show()
```



Estimation par intervalle de confiance de la moyenne d'une loi normale dont l'écart type est connu

Soit (X_1, X_2, \dots, X_n) un n -échantillon issu d'une loi normale $\mathcal{N}(\mu, \sigma^2)$. On suppose σ connu mais l'espérance μ est inconnue et on cherche à l'estimer. On considère \bar{X}_n la moyenne empirique de l'échantillon. Nous avons vu que c'est un estimateur sans biais et convergent de μ . De plus, par les règles de stabilités des lois normales indépendantes :

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

On a aussi $Y_n = \sqrt{n} \cdot \frac{\bar{X}_n - \mu}{\sigma} \hookrightarrow \mathcal{N}(0, 1)$.

Posons t_α positif tel que $t_\alpha = \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)$. Avec ce choix

$$\mathbf{P}(-t_\alpha \leq Y_n \leq t_\alpha) = \Phi(t_\alpha) - \Phi(-t_\alpha) = 2\Phi(t_\alpha) - 1 = 1 - \alpha.$$

En revenant à \bar{X}_n , on trouve

$$\mathbf{P}\left(\bar{X}_n - t_\alpha \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + t_\alpha \frac{\sigma}{\sqrt{n}}\right) = \mathbf{P}\left(-t_\alpha \leq \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \leq t_\alpha\right) = \mathbf{P}(-t_\alpha \leq Y_n \leq t_\alpha) = 1 - \alpha$$

On vient de montrer 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 μ avec un niveau de confiance égal à 0,95.

Applications numériques

- Pour un risque $\alpha = 5\%$, $1 - \frac{\alpha}{2} = 0,975$ et $t_{0,05} = \Phi^{-1}(0,975) \simeq 1,96$.
- Pour un risque de $\alpha = 1\%$, $1 - \frac{\alpha}{2} = 0,995$ et $t_{0,01} = \Phi^{-1}(0,995) \simeq 2,58$.

◆◆ Comparaison des méthodes précédentes

Dans une population donnée, une étude statistique faite sur un groupe de 100 personnes donne lieu à la série statistique suivante.

Poids	48	49	50	51	52	53	54	55
Effectif	3	5	2	6	6	10	12	10
Poids	56	57	58	59	60	61	62	63
Effectif	9	8	8	6	5	4	3	3

Exercice 8



On suppose que le poids d'un individu du groupe est une variable aléatoire X qui suit une loi normale d'écart-type $\sigma = 3,5$. Dans chaque groupe de 100 personnes étudié, on désigne par X_i la variable aléatoire égale au poids du i -ème individu, pour tout $i \in \llbracket 1, 100 \rrbracket$.

1. En utilisant les deux méthodes précédentes, déterminer une valeur approchée d'un intervalle de confiance à 95%, de la moyenne des poids des individus.
2. Comparer les deux méthodes.

p. 26

2.2 Intervalle de confiance asymptotique

DÉFINITION

intervalle de confiance asymptotique

Soient $\alpha \in]0, 1[$, U_n et V_n deux estimateurs de $g(\theta)$ tels que pour tout $\theta \in \Theta$

$$\mathbf{P}_\theta (U_n \leq g(\theta) \leq V_n) \geq 1 - \alpha_n \quad \text{et} \quad \alpha_n \xrightarrow[n \rightarrow \infty]{} \alpha.$$

On dit que l'intervalle $[U_n, V_n]$ est un **intervalle de confiance asymptotique** de $g(\theta)$ avec un risque d'au plus α ou au niveau de confiance au moins égal à $1 - \alpha$.

Intervalle de confiance asymptotique du paramètre d'une loi de Bernoulli

Soit (X_1, X_2, \dots, X_n) un n -échantillon issu d'une loi de Bernoulli de paramètre p . Par indépendance, on sait que

$$n\bar{X}_n = \sum_{i=1}^n X_i \hookrightarrow \mathcal{B}(n, p).$$

Par le théorème central limite, on sait de plus que

$$\bar{X}_n^* = \frac{\bar{X}_n - \mathbf{E}(\bar{X}_n)}{\sqrt{\mathbf{V}(\bar{X}_n)}} = \sqrt{n} \left(\frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \right) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} Z \quad \text{où} \quad Z \hookrightarrow \mathcal{N}(0, 1).$$

Autrement dit, pour tous $a, b \in \mathbb{R}$ avec $a < b$

$$\mathbf{P} \left(a < \bar{X}_n^* \leq b \right) \xrightarrow[n \rightarrow \infty]{} \mathbf{P}(a < Z \leq b) = \Phi(b) - \Phi(a).$$

En particulier, pour un intervalle centré en 0, $a = -b$ et en reprenant les notations précédentes : soit t_α tel que $\Phi(t_\alpha) = 1 - \frac{\alpha}{2}$.

$$\begin{aligned} \mathbf{P}\left(\bar{X}_n - t_\alpha \frac{\sqrt{p(1-p)}}{\sqrt{n}} \leq p \leq \bar{X}_n + t_\alpha \frac{\sqrt{p(1-p)}}{\sqrt{n}}\right) &= \mathbf{P}\left(-t_\alpha \frac{\sqrt{p(1-p)}}{\sqrt{n}} \leq \bar{X}_n - p \leq t_\alpha \frac{\sqrt{p(1-p)}}{\sqrt{n}}\right) \\ &= \mathbf{P}\left(-t_\alpha < \sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \leq t_\alpha\right) \\ &= \mathbf{P}(t_\alpha < \bar{X}_n^* \leq t_\alpha) \xrightarrow[n \rightarrow \infty]{} \Phi(t_\alpha) - \Phi(-t_\alpha) = 2\Phi(t_\alpha) - 1 = 1 - \alpha. \end{aligned}$$

Ensuite, à partir de l'encadrement $0 \leq p(1-p) \leq \frac{1}{4}$

$$\left[\bar{X}_n - t_\alpha \frac{\sqrt{p(1-p)}}{\sqrt{n}} \leq p \leq \bar{X}_n + t_\alpha \frac{\sqrt{p(1-p)}}{\sqrt{n}}\right] \subset \left[\bar{X}_n - \frac{t_\alpha}{2\sqrt{n}} \leq p \leq \bar{X}_n + \frac{t_\alpha}{2\sqrt{n}}\right].$$

Par croissance de la probabilité

$$\mathbf{P}\left(\bar{X}_n - t_\alpha \frac{\sqrt{p(1-p)}}{\sqrt{n}} \leq p \leq \bar{X}_n + t_\alpha \frac{\sqrt{p(1-p)}}{\sqrt{n}}\right) \leq \mathbf{P}\left(\bar{X}_n - \frac{t_\alpha}{2\sqrt{n}} \leq p \leq \bar{X}_n + \frac{t_\alpha}{2\sqrt{n}}\right).$$

Si, pour tout $n \in \mathbb{N}$, on pose α_n de sorte que $1 - \alpha_n$ soit le terme de gauche dans l'inégalité précédente, on a bien trouvé une suite $(\alpha_n)_{n \in \mathbb{N}^*}$ telle que

$$\mathbf{P}\left(\bar{X}_n - \frac{t_\alpha}{2\sqrt{n}} \leq p \leq \bar{X}_n + \frac{t_\alpha}{2\sqrt{n}}\right) \geq 1 - \alpha_n \quad \text{et} \quad \alpha_n \xrightarrow[n \rightarrow \infty]{} \alpha.$$

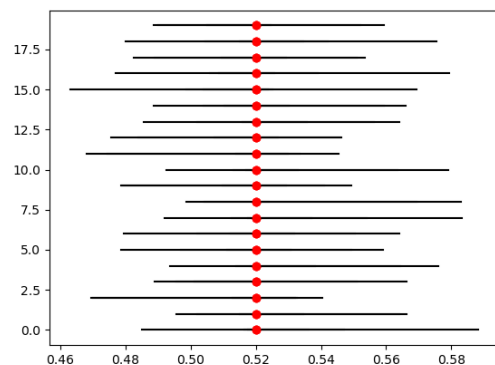
En reprenant la définition, un intervalle de confiance asymptotique de p au niveau de confiance $1 - \alpha$ est

$$\left[\bar{X}_n - \frac{t_\alpha}{2\sqrt{n}}, \bar{X}_n + \frac{t_\alpha}{2\sqrt{n}}\right].$$

• Simulation Python

Editeur

```
p=0.52 # le paramètre "inconnu" à estimer
n=10000 # Nombre de sondés
alpha=0.05
talpha=1.96
NbreTest=20
for i in range(NbreTest):
    Ech=np.random.rand(n)<p
    u=np.mean(Ech)-talpha/(2*(n)**(1/2))
    v=np.mean(Ech)+talpha/(2*(n)**(1/2))
    plt.plot([u,v],[i,i],'k')
    plt.plot([p],[i],'ro')
plt.show()
```

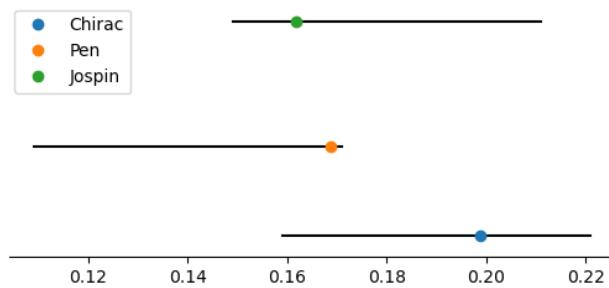


Remarque. Les intervalles obtenus sont bien plus précis que ceux obtenus par l'inégalité de Bienaymé-Tchebychev. En effet, les longueurs des intervalles obtenus par le théorème central limite est plus petite que celles obtenues par l'inégalité de Bienaymé-Tchebychev.

Exemple. Premier tour de l'élection présidentielle de 2002 (21 avril)

Voici quelques résultats des sondages d'opinion quelques jours avant l'élection.

	Jacques CHIRAC	Jean Marie LE PEN	Lionel JOSPIN
CSA - 11 avril	21%	12%	19%
IFOP - 12 avril	19%	11,5%	17%
SOFRES - 13 avril	20%	13%	18%
CSA - 18 avril	19,15%	14%	18%
IPSOS - 18 avril	20%	14%	18%
BVA - 19 avril	19%	14%	18%
Sondage confidentiel - 21 avril	18%	14,5%	17%
Résultat final	19,88%	16,88%	16,18%



Extrait du site du conseil constitutionnel : "Dans le cas précis de l'élection présidentielle, la publication d'un sondage d'intention de vote prévoyant un faible écart entre les candidats (51-49 par exemple) rend nécessaire une telle précaution si l'on considère que, dans le cas des sondages portant sur un échantillon de 1 000 personnes (qui est l'échantillon standard en France pour les élections nationales et notamment l'élection présidentielle), la marge d'erreur est estimée à 3%.

L'obligation pour les instituts de sondage de rendre publique leur marge d'erreur est prévue par la proposition de loi sénatoriale."



Exercices



Exercice 9. ♦ Estimation d'une variance

Soient $n \in \mathbb{N}^*$ et $p \in]0; 1[$. Une variable aléatoire $X \hookrightarrow \mathcal{B}(n, p)$. Démontrer qu'il existe $\alpha_n, \beta_n \in \mathbb{N}^*$ qui ne dépendent que de n , tels que la variable aléatoire $T = \alpha_n X + \beta_n X^2$ soit un estimateur sans biais de $V(X)$.

Exercice 10. ♦ Estimation d'un paramètre d'une loi uniforme

Soit $\theta \in \mathbb{R}_*^+$. Soit X une variable aléatoire suivant une loi uniforme sur $]0; \theta[$ et (X_1, \dots, X_n) un n -échantillon de X .

1. Montrer que pour tout $n \in \mathbb{N}^*$, $T_n = \frac{2}{n}(X_1 + \dots + X_n)$ est un estimateur sans biais de θ .
2. Considérons Y_{\min} et Y_{\max} définies par :

$$Y_{\min} = \min_{i \in \llbracket 1; n \rrbracket} X_i \quad \text{et} \quad Y_{\max} = \max_{i \in \llbracket 1; n \rrbracket} X_i.$$

Pour chacune de ces deux variables, préciser :

- a) La fonction de répartition,
 - b) La densité de probabilité,
 - c) L'espérance et la variance. *On pourra remarquer que Y_{\min} et $\theta - Y_{\max}$ ont même loi.*
3. Posons $T'_n = \frac{n+1}{n} Y_{\max}$. Justifier que T'_n est un estimateur sans biais de θ .
 4. Quel est le meilleur estimateur de θ entre T'_n et T_n ?
Un estimateur sans biais est d'autant meilleur que sa variance est faible.
 5. Posons $T''_n = Y_{\min} + Y_{\max}$. Montrer sans calculs superflus que $V(T''_n) \leq 4V(Y_{\max})$.
 6. En déduire que T''_n est meilleur estimateur de θ que T_n .
 7. Est-ce que les estimateurs T_n, T'_n et T''_n sont convergents?

» Solution p. 26

Exercice 11. ♦ Risque quadratique

Soit T_n un estimateur de $g(\theta)$. On suppose que pour tout $\theta \in \Theta$, T_n admet un moment d'ordre 2 pour la probabilité \mathbf{P}_θ . Pour tout $\theta \in \Theta$, on appelle risque quadratique de T_n en $g(\theta)$ et on note $r_\theta(T_n)$ le réel :

$$r_\theta(T_n) = \mathbf{E}_\theta \left((T_n - g(\theta))^2 \right).$$

1. Justifier que $r_\theta(T_n) = [b_\theta(T_n)]^2 + V_\theta(T_n)$.
2.
 - a) Justifier que si le risque quadratique de T_n tend vers 0 quand n tend vers l'infini, alors T_n est un estimateur convergent de $g(\theta)$.
 - b) En déduire que, si T_n est un estimateur asymptotiquement sans biais de $g(\theta)$ et si la variance de T_n tend vers 0 lorsque n tend vers l'infini, alors T_n est un estimateur convergent de $g(\theta)$.
3. Pour comparer deux estimateurs T_n et U_n de $g(\theta)$, on peut calculer leur risque quadratique. Si on a :

$$\forall \theta \in \Theta, \quad r_\theta(T_n) \leq r_\theta(U_n),$$

alors on dira que T_n est un meilleur estimateur de $g(\theta)$ que U_n .

a) Exemple 1.

Soit $(X_n)_{n \geq 1}$ une suite de variables aléatoires indépendantes suivant la loi exponentielle $\mathcal{E}(\lambda)$. On pose

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{et} \quad T_n = \frac{1}{n+1} \sum_{i=1}^n X_i.$$

Comparer ces deux estimateurs de $1/\lambda$.

b) Exemple 2.

Soient T_0 et T_1 deux estimateurs de θ , sans biais et indépendants. Pour tout a réel, on pose $T_a = aT_1 + (1-a)T_0$.

- i) T_a est-il un estimateur sans biais de θ ?
- ii) Parmi tous les estimateurs T_a , lequel est le meilleur?

>> Solution p. ??

Exercice 12. ♦♦ Estimateur sans biais de l'écart-type σ d'une loi normale centrée

extrait HEC 2005

Soit X une variable aléatoire qui suit une loi normale centrée et d'écart-type σ , le paramètre réel inconnu σ , est strictement positif.

1. Montrer que la variable aléatoire $T = \frac{X^2}{2\sigma^2}$ suit une loi γ de paramètre $1/2$. En déduire la valeur de $\Gamma(1/2)$.
2. Pour n entier naturel non nul, on considère un n -échantillon (X_1, X_2, \dots, X_n) constitué de variables indépendantes et de même loi que X .
 - a) On désigne par S_n la variable aléatoire $S_n = \sum_{i=1}^n \frac{X_i^2}{2\sigma^2}$. Quelle est la loi de probabilité de S_n ?
 - b) En déduire que la variable aléatoire Y_n définie par $Y_n = \frac{1}{n} \sum_{i=1}^n X_i^2$ est un estimateur sans biais de σ^2 .

>> Solution p. ??

Exercice 13. ♦♦

d'après EDHEC 2000

Un sondage consiste à proposer l'affirmation « A » à certaines personnes d'une population donnée. Le sujet abordé étant délicat, le stratagème suivant est mis en place afin de mettre en confiance les personnes sondées pour qu'elles ne mentent pas...

L'enquêteur dispose d'un paquet de 20 cartes, numérotées de 1 à 20, qu'il remet à la personne sondée. Celle-ci tire une carte au hasard et ne la montre pas à l'enquêteur. La règle est alors la suivante :

- si la carte porte le numéro 1, la personne sondée répond "vrai" si elle est d'accord avec l'affirmation « A » et "faux" sinon.
- si la carte porte un autre numéro, la personne sondée répond "vrai" si elle n'est pas d'accord avec l'affirmation « A » et "faux" sinon.

Le but de l'enquête est d'évaluer la proportion p de gens de cette population qui sont réellement d'accord avec l'affirmation « A ».

1. On interroge une personne selon ce procédé et on considère l'événement suivant, noté V : « la personne répond "vrai" ». On note $\theta = \mathbf{P}(V)$. En utilisant la formule des probabilités totales, exprimer θ en fonction de p , puis en déduire p en fonction de θ .
2. Certaines considérations théoriques laissent penser que $p = \frac{17}{18}$.
 - a) Vérifier que $\theta = \frac{1}{10}$.
 - b) Calculer la probabilité pour qu'une personne ayant répondu "vrai" soit d'accord avec l'affirmation « A ».

On revient au cas général où l'on ne connaît ni p , ni θ .

3. On considère un échantillon aléatoire, de taille n , extrait de la population considérée et on note S_n le nombre de réponses "vrai" obtenues. On suppose n assez grand pour pouvoir considérer que cet échantillonnage est assimilable à un tirage avec remise.
 - a) Donner la loi de S_n ainsi que son espérance et sa variance.
 - b) Montrer que $\frac{S_n}{n}$ est un estimateur sans biais et convergent de θ .
4. Dans cette question, on suppose que l'on a réalisé un échantillon de 100 personnes et on constate que 23 personnes ont répondu "vrai".
 - a) Donner une estimation ponctuelle de θ et de p .
 - b) Donner un intervalle de confiance à 95% de θ puis de p .
On rappelle que, si Φ désigne la fonction de répartition d'une variable X suivant la loi normale $\mathcal{N}(0, 1)$, alors $\Phi(1,96) = 0,975$

>> Solution p. ??

Exercice 14. ♦♦ Estimation des paramètres d'une loi de Pareto

Soient $(b, \theta) \in \mathbb{R}^{+*} \times \mathbb{R}$. On suppose que X est une va à densité qui prend ses valeurs dans $[\theta; +\infty[$, dont une densité f est définie par :

$$\forall x \in \mathbb{R}, \quad f(x) = \begin{cases} 0 & \text{si } x < \theta \\ \frac{1}{b} \exp\left(-\frac{x-\theta}{b}\right) & \text{si } x \geq \theta. \end{cases}$$

Soient $n \in \mathbb{N} \setminus \{0; 1\}$ et (X_1, \dots, X_n) un n -échantillon de même loi que X . On pose :

$$\overline{X}_n = \frac{X_1 + \dots + X_n}{n} \quad \text{et} \quad T_n = \min(X_1, \dots, X_n).$$

1. Soit $k \in \llbracket 1; n \rrbracket$. Reconnaitre la loi de $Y_k = (X_k - \theta)/b$.
2. On pose $\overline{Y}_n = \frac{Y_1 + \dots + Y_n}{n}$ et $U_n = \min(Y_1, \dots, Y_n)$. Calculer $\mathbf{E}(\overline{Y}_n)$ et $\mathbf{E}(U_n)$.
3. Exprimer les variables aléatoires \overline{X}_n et T_n en fonction des variables \overline{Y}_n et U_n . En déduire $\mathbf{E}(\overline{X}_n)$ et $\mathbf{E}(T_n)$.

4. Déterminer un estimateur sans biais $\hat{\theta}_n$ de θ et un estimateur sans biais \hat{b}_n de b sous la forme de combinaisons linéaires de \bar{X}_n et T_n .

>> Solution p. 27

Exercice 15. ♦♦ Estimations des paramètres d'une loi de Pareto

Une variable X suit une loi de Pareto $VP(\alpha, 1, x_0)$ si la fonction de répartition est donnée par

$$\forall x \in \mathbb{R}, \quad F(x) = \begin{cases} 0 & \text{si } x < 1 + x_0 \\ 1 - \frac{1}{(x-x_0)^\alpha} & \text{sinon.} \end{cases}$$

1. Dans la suite, on souhaite estimer le paramètre α d'une loi $VP(\alpha, 1, 0)$.
- Donner la loi suivie par $Z = \ln(X)$.
 - Soient $(X_n)_{n \in \mathbb{N}^*}$ une suite de variables mutuellement indépendantes qui suivent toutes une loi de Pareto $VP(\alpha, 1, 0)$. Pour tout $n \in \mathbb{N}^*$, on pose

$$Z_n = \ln(X_1 X_2 \dots X_n).$$

Donner la loi de Z_n .

- Préciser l'espérance et la variance de Z_n , en déduire que $T_n = \frac{1}{n} Z_n$ est un estimateur de α^{-1} .
 - Démontrer que $W_n = \frac{n}{Z_n}$ est un estimateur de α , calculer son espérance, sa variance.
2. On revient au cas général. On veut maintenant estimer le paramètre x_0 d'une loi de Pareto de paramètre $(\alpha, 1, x_0)$.
- Démontrer que X suit une loi de Pareto $VP(\alpha, 1, x_0)$ si et seulement si $X_0 = X - x_0$ suit un loi de Pareto $VP(\alpha, 1, 0)$.
 - Calculer $E(X_0)$ et $V(X_0)$ pour $\alpha > 2$, et en déduire

$$E(X) = x_0 + \frac{\alpha}{\alpha - 1} \quad \text{et} \quad V(X) = \frac{\alpha}{(\alpha - 2)(\alpha - 1)^2}.$$

- Soient $(X_n)_{n \in \mathbb{N}^*}$ une suite de variables mutuellement indépendantes qui suivent toutes une loi de Pareto $VP(\alpha, 1, x_0)$, on pose pour tout $n \in \mathbb{N}^*$,

$$Y_n = \inf(X_1, X_2, \dots, X_n).$$

Préciser la loi, l'espérance la variance de Y_n . En déduire que $\tilde{Y}_n = Y_n - 1$ est un estimateur de rang n de x_0 . Quel est le risque lié à l'estimateur \tilde{Y}_n ?

Exercice 16. ♦♦ Le boulanger

d'après oral ESCP 2022, sujet 23

Un boulanger vend du pain chaque jour.

- La quantité de pain produite chaque jour est une quantité fixée Q choisie par le boulanger, Q étant exprimée en kilogramme.
- La demande de pain de la part des clients est une variable aléatoire X strictement positive, toujours exprimée en kilogramme.
- On suppose que la variable X admet une densité f strictement positive sur \mathbb{R}_+^* , nulle sur \mathbb{R}_- , continue sur \mathbb{R} et on note F la fonction de répartition de X .
- Le coût de fabrication par kilogramme est c euros et le prix de vente est v euros par kilogramme.
- On note B la variable aléatoire égale au bénéfice quotidien.
- La variable indicatrice d'un événement A est notée $\mathbb{1}_A$.
- On suppose que $0 < c < v$.

Si la demande de pain X est inférieure à l'offre Q , le boulanger ne vend que la quantité X (le pain invendu un jour donné n'est pas remis en vente le lendemain!); si la demande est supérieure à l'offre, il ne vend que la quantité Q . Dans ces conditions, on cherche la quantité optimale à produire, c'est-à-dire la quantité Q_0 qui maximise l'espérance de B .

1. Établir la relation suivante :

$$B = v [Q + (X - Q) \mathbb{1}_{[X < Q]}] - cQ.$$

2. Montrer que la variable $X \mathbb{1}_{[X < Q]}$ admet une espérance et donner son expression sous forme d'intégrale.
3. En déduire l'égalité suivante :

$$E(B) = (v - c)Q + v \left[\int_0^Q t f(t) dt - QF(Q) \right].$$

4. Exprimer Q_0 à l'aide de F , de v et de c . Le boulanger cherche à prévoir sa demande journalière. La demande aléatoire X_n qui va s'exprimer le jour n n'est pas connu à l'avance mais le boulanger fait l'hypothèse que la demande ne variera pas beaucoup d'un jour à l'autre et que :

$$X_{n+1} = X_n + U_{n+1}$$

où :

- X_0 est une constante strictement positive fixée.
 - Les U_k sont des variables aléatoires indépendantes, de même loi, d'espérance nulle et de variance σ^2 non nulle.
5. a) Exprimer X_n en fonction des U_i et de X_0 .
- b) Montrer que la suite $\left(\frac{X_n}{n}\right)$ converge en probabilité vers 0.
- c) Démontrer que si deux suites de variables aléatoires (A_n) et (B_n) convergent en probabilité respectivement vers des variables aléatoires a et b , alors la suite $(A_n + B_n)$ converge en probabilité vers $a + b$.
- d) En déduire que $\left(\frac{X_n}{n}\right)$ converge en probabilité vers une variable que l'on précisera.
- e) Montrer que la suite $\left(\frac{X_n}{n}\right)$ converge en loi vers 0.
- f) Montrer que la suite $\left(\frac{X_n}{n}\right)$ converge en loi et préciser la loi limite.
On pourra utiliser le Théorème de Slutsky : Si X_n converge en loi vers X , et si Y_n converge en probabilité vers une constante c , alors $X_n Y_n$ converge en loi vers cX .

>> Solution p. ??

Exercice 17. ♦♦ Calcul d'un intervalle de confiance

Soit (X_1, \dots, X_n) un échantillon de la loi exponentielle de paramètre $\lambda > 0$ inconnu. On pose

$$Y_n = \min(X_1, \dots, X_n).$$

1. Montrer que Y_n suit une loi exponentielle dont on précisera le paramètre.
2. Soit $\alpha \in]0; 2/3[$. Déterminer deux réels a_n et b_n tels que

$$\mathbf{P}(n\lambda Y_n \leq a_n) = \frac{\alpha}{2} = \mathbf{P}(n\lambda Y_n \geq b_n).$$

3. Justifier que $\left[\frac{nY_n}{a_n}, \frac{nY_n}{b_n}\right]$ est un intervalle de confiance de $\frac{1}{\lambda}$ au niveau de confiance d'au moins $1 - \alpha$.

>> Solution p. 28

Exercice 18. ♦♦ Calcul d'un intervalle de confiance II

Soit $(X_i)_{i \in \mathbb{N}^*}$ une suite de variables aléatoires indépendantes et de même loi exponentielle $\mathcal{E}(\lambda)$ où le paramètre $\lambda > 0$ est inconnu. On pose $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$.

1. Justifier que $\lambda\sqrt{n}\bar{X}_n - \sqrt{n}$ converge en loi vers une variable aléatoire de loi $\mathcal{N}(0, 1)$.
2. Soit $\alpha \in]0, 1[$. Notons t_α l'unique réel tel que $\Phi(t_\alpha) = 1 - \frac{\alpha}{2}$. Montrer que

$$\left[\left(1 - \frac{t_\alpha}{\sqrt{n}}\right) \frac{1}{\bar{X}_n}, \left(1 + \frac{t_\alpha}{\sqrt{n}}\right) \frac{1}{\bar{X}_n} \right]$$

est un intervalle de confiance asymptotique de λ au niveau de risque α .

>> Solution p. ??

Exercice 19. ♦♦ Estimation et loi de Poisson

d'après oral ESCP 2022, sujet 35

Soit X_1, \dots, X_n un n -échantillon de variables aléatoires indépendantes qui suivent la loi de Poisson de paramètre $\lambda \in \mathbb{R}_*^+$ inconnu. On cherche dans cet exercice à estimer $e^{-\lambda}$.

Pour tout $k \in \llbracket 1, n \rrbracket$, on note Y_k la fonction indicatrice de l'événement $[X_k = 0]$. Pour tout $n \in \mathbb{N}^*$, on pose :

$$\bar{Y}_n = \frac{1}{n} \sum_{k=1}^n Y_k \quad \text{et} \quad S_n = \sum_{k=1}^n X_k$$

1. a) Déterminer la loi de Y_k .
 b) Montrer que \bar{Y}_n est un estimateur sans biais de $e^{-\lambda}$.
 c) \bar{Y}_n est-il un estimateur convergent de $e^{-\lambda}$?
2. Pour $j \in \mathbb{N}$, on pose $\varphi(j) = \mathbf{P}_{[S_n=j]}(X_1 = 0)$. Calculer $\varphi(j)$.
3. On pose à présent $T_n = \varphi(S_n)$.
 a) Montrer que T_n est un estimateur sans biais de $e^{-\lambda}$.
 b) T_n est-il un estimateur convergent de $e^{-\lambda}$?
4. Comparer les variances des deux estimateurs T_n et \bar{Y}_n .

>> Solution p. ??

Lien entre estimation et optimisation

Exercice 20. ♦

Soit $(X_1, X_2, \dots, X_n, Y_1, Y_2, \dots, Y_m)$ un échantillon de variables aléatoires indépendantes, de même loi de Bernoulli de paramètre inconnu $p \in]0; 1[$. On pose :

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i, \quad \bar{Y}_m = \frac{1}{m} \sum_{i=1}^m Y_i \quad \text{et} \quad Z = a\bar{X}_n + b\bar{Y}_m, \quad \text{où} \quad (a, b) \in \mathbb{R}^2.$$

Pour quelles valeurs de a et b , Z est-il le meilleur estimateur sans biais de p ?

» Solution p. 28

Exercice 21. ♦♦ Soit X une variable aléatoire discrète d'espérance $\mathbf{E}(X) = \theta \neq 0$ et de variance $\mathbf{V}(X) = 1$. On considère un n -échantillon (X_1, \dots, X_n) de la loi de X . On pose classiquement

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

1. Montrer que \bar{X}_n est un estimateur sans biais de θ et calculer son risque quadratique défini par

$$r_\theta = \mathbf{E}_\theta((\bar{X}_n - \theta)^2).$$

2. Soit $(\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n$. On note $Y_n = \sum_{i=1}^n \alpha_i X_i$.

- Donner une condition nécessaire et suffisante sur $\alpha_1, \dots, \alpha_n$ pour que Y_n soit un estimateur sans biais de θ . On suppose dans la suite que cette condition est vérifiée.
- Calculer $\text{cov}(\bar{X}_n, Y_n)$. En déduire $\mathbf{V}(\bar{X}_n) \leq \mathbf{V}(Y_n)$. Que dire si $\mathbf{V}(\bar{X}_n) = \mathbf{V}(Y_n)$?
- Que peut-on en déduire sur les estimateurs \bar{X}_n et Y_n de θ .

» Solution p. 28

Exercice 22. ♦♦ Maximum de vraisemblance

D'après EDHEC 2014, voie E

Dans cet exercice, θ désigne un réel strictement positif et n un entier naturel supérieur ou égal à 2. Pour tout k de \mathbb{N} , on pose :

$$u_k = \frac{1}{1+\theta} \left(\frac{\theta}{1+\theta} \right)^k.$$

1. Montrer que la suite $(u_k)_{k \in \mathbb{N}}$ définit bien une loi de probabilité. On considère maintenant une variable aléatoire X prenant ses valeurs dans \mathbb{N} et dont la loi est donnée par :

$$\forall k \in \mathbb{N}, \quad \mathbf{P}(X = k) = u_k.$$

- On pose $Y = X + 1$. Reconnaître la loi de Y , puis en déduire l'espérance et la variance de X .
 - Écrire une fonction Python qui prend en argument θ et simule la loi d'une variable aléatoire X .
- Dans cette question, on souhaite estimer le paramètre θ par la méthode du maximum de vraisemblance. Pour ce faire, on considère un échantillon (X_1, X_2, \dots, X_n) composé de variables aléatoires indépendantes ayant toutes la même loi que X et on introduit la fonction L , de \mathbb{R}_+^* dans \mathbb{R} , définie par :

$$\forall \theta \in \mathbb{R}_+^*, \quad L(\theta) = \prod_{k=1}^n \mathbf{P}(X_k = x_k)$$

où x_1, x_2, \dots, x_n désignent des entiers naturels éléments de $X(\Omega)$.

L'objectif est de choisir la valeur de θ qui rend $L(\theta)$ maximale.

- Écrire $\ln(L(\theta))$ en fonction de θ et de $S_n = \sum_{k=1}^n x_k$.
 - On considère la fonction φ , définie par :

$$\forall \theta \in]0; +\infty[, \quad \varphi(\theta) = S_n \ln \theta - (S_n + n) \ln(1 + \theta).$$

Montrer que la fonction φ admet un maximum, atteint en un seul réel que l'on notera $\hat{\theta}_n$ et que l'on exprimera en fonction de S_n . Que représente $\hat{\theta}_n$ pour la fonction L ?

On pose dorénavant : $T_n = \frac{1}{n} \sum_{i=1}^n X_i$. La variable T_n est appelée estimateur du maximum de vraisemblance pour θ .

- Vérifier que T_n est un estimateur sans biais de θ .

d) On définit le risque quadratique par

$$r_\theta(T_n) = \mathbf{E}_\theta \left((T_n - g(\theta))^2 \right).$$

Calculer $r_\theta(T_n)$ de T_n et vérifier que $\lim_{n \rightarrow +\infty} r_\theta(T_n) = 0$.

>> Solution p. ??

Exercice 23. ♦♦♦ Maximum de vraisemblance

D'après HEC, 2020

Lorsque l'on cherche à estimer un paramètre inconnu à partir d'un échantillon de données, on appelle statistique exhaustive toute fonction de ces données qui résume à elle seule l'information que ces données fournissent sur le paramètre. On donne ici une définition précise de cette notion d'exhaustives dans le cas des échantillons de variables aléatoires discrètes, illustrée de plusieurs exemples qui en montrent l'intérêt. On s'intéressera dans ce problème à l'estimation d'un paramètre réel inconnu θ appartenant à un intervalle Θ .

On dispose pour cela de plusieurs observations x_1, \dots, x_n considérées comme les réalisations de variables aléatoire discrètes X_1, \dots, X_n définie sur le même espace probabilisable (Ω, \mathcal{A}) , à valeurs dans une partie B de \mathbb{N} .

L'espace probabilisable (Ω, \mathcal{A}) est muni d'une famille $(\mathbf{P}_\theta)_{\theta \in \Theta}$ de probabilités indexées par le paramètre θ .

On fait, pour toutes les valeurs du paramètre θ , les trois hypothèses suivantes.

- Les variables aléatoires X_1, \dots, X_n sont mutuellement indépendantes, c'est-à-dire :

$$\forall (x_1, \dots, x_n) \in B^n, \quad \mathbf{P}_\theta \left(\bigcap_{i=1}^n [X_i = x_i] \right) = \prod_{i=1}^n \mathbf{P}_\theta([X_i = x_i]) \quad (1)$$

- Les variables aléatoires X_1, \dots, X_n suivent toutes la même loi qu'une variable aléatoire de référence, dotée X, à valeurs dans B, c'est-à-dire :

$$\forall i \in [1, n], \quad \forall x \in B, \quad \mathbf{P}_\theta([X_i = x]) = \mathbf{P}_\theta([X = x]) \quad (2)$$

- Tous les éléments de B sont des valeurs effectivement possibles de X, c'est-à-dire :

$$\forall x \in B, \quad \mathbf{P}_\theta([X = x]) > 0 \quad (3)$$

On appelle statistique toute variable aléatoire S de la forme $\omega \mapsto s(X_1(\omega), \dots, X_n(\omega))$, où s désigne une application définie sur B^n et à valeurs réelles. On note alors $S = s(X_1, \dots, X_n)$.

Pour tout $\theta \in \Theta$, on note $\mathbf{E}_\theta(S)$ l'espérance de S lorsque (Ω, \mathcal{A}) est muni de la probabilité \mathbf{P}_θ (si cette espérance existe). On note de même $\mathbf{V}_\theta(S)$ la variance de S (si elle existe).

Partie 1 : développements en série

1. Dans cette question, x désigne un nombre réel strictement compris entre 0 et 1.

a) Justifier la convergence de la série $\sum_{k \geq 1} x^k/k$.

b) Vérifier, pour tout $m \in \mathbb{N}^*$ et tout $t \in]0, 1[$, l'égalité : $\frac{1}{1-t} = \frac{t^m}{1-t} + \sum_{k=0}^{m-1} t^k$.

c) Démontrer que l'intégrale $\int_0^x \frac{t^m}{1-t} dt$ tend vers 0 quand l'entier m tend vers l'infini.

d) En déduire la somme de la série $\sum_{k \geq 1} x^k/k$.

2. Dans cette question, indépendante de la précédente, $(a_k)_{k \in \mathbb{N}}$ désigne une suite de nombres réels telle que la série $\sum_{k \geq 0} a_k c^k$ est absolument convergente pour un réel strictement positif c.

a) Justifier que la fonction $f : x \mapsto a_0 + \sum_{k=1}^{+\infty} a_k x^k$ est bien définie sur le segment $[-c, c]$.

b) Pour un entier naturel m, on pose : $M_m = \sum_{k=m+1}^{+\infty} |a_k| c^{k-m-1}$. Justifier, pour tout $x \in [-c, c]$, l'inégalité :

$$\left| \sum_{k=m+1}^{+\infty} a_k x^k \right| \leq M_m |x|^{m+1}.$$

c) Justifier, pour tout $m \in \mathbb{N}^*$, le développement limite au voisinage de 0 :

$$f(x) = a_0 + \sum_{k=1}^m a_k x^k + o(x^m)$$

d) Démontrer que si la fonction f est nulle sur l'intervalle $]0, c[$, alors $(a_k)_{k \in \mathbb{N}}$ est la suite nulle.

- Dans toute la suite du problème, pour tout $\theta \in \Theta$ et tout $(x_1, \dots, x_n) \in B^n$, on note :

$$L(x_1, \dots, x_n, \theta) = \prod_{i=1}^n \mathbf{P}_\theta([X_i = x_i]) \quad (4)$$

Cette quantité, qui s'écrit aussi $\prod_{i=1}^n \mathbf{P}_\theta([X = x_i])$ d'après (2), est appelée *la vraisemblance* de la valeur θ du paramètre au vu des observations x_1, \dots, x_n .

Partie II : estimateur du maximum de vraisemblance, un exemple

Dans cette partie, Θ est l'intervalle ouvert $]0, 1[$, B est égal à \mathbb{N}^* et on a :

$$\forall x \in B, \quad \mathbf{P}_\theta([X = x]) = (1 - \theta)^{x-1} \theta.$$

On note \bar{X} la variable aléatoire $\frac{1}{n} \sum_{i=1}^n X_i$.

3. Soit $\theta \in \Theta$

- Reconnaître la loi de X lorsque (Ω, \mathcal{A}) est muni de la probabilité \mathbf{P}_θ .
- En déduire que \bar{X} est un estimateur sans biais du paramètre $1/\theta$
- Quel est le risque quadratique de cet estimateur? Rappelons que le risque quadratique est défini par

$$r_\theta(\bar{X}) = \mathbf{E}_\theta \left((\bar{X} - 1/\theta)^2 \right).$$

4. On note T la variable aléatoire $\frac{1}{n} \sum_{i=1}^n \frac{1}{X_i}$.

- En utilisant le résultat de la question 1.d, justifier que :

$$\forall \theta \in \Theta, \quad \mathbf{E}_\theta(T) = \frac{\theta \ln(\theta)}{\theta - 1}.$$

- En déduire que T est un estimateur de θ dont le biais $b_\theta(T)$ est strictement positif.

5. Soit $(x_1, \dots, x_n) \in B^n$.

- Justifier, pour tout $\theta \in \Theta$, l'égalité :

$$\ln(L(x_1, \dots, x_n, \theta)) = n \ln(\theta) - \left(n - \sum_{i=1}^n x_i \right) \ln(1 - \theta).$$

- En déduire que, lorsque les x_i ne sont pas tous égaux à 1, le nombre $n / \left(\sum_{i=1}^n x_i \right)$ est l'unique valeur de θ qui maximise la vraisemblance $L(x_1, \dots, x_n, \theta)$

6. On note U la variable aléatoire $n / \left(\sum_{i=1}^n X_i \right)$.

- établir, pour tout $\theta \in \Theta$ et tout entier $k \geq n$, l'égalité :

$$\frac{n}{k} = \theta - \theta^2 \left(\frac{k}{n} - \frac{1}{\theta} \right) + \int_{1/\theta}^{k/n} \left(\frac{k}{n} - t \right) \frac{2}{t^3} dt.$$

- En déduire que U est un estimateur de θ dont le biais $b_\theta(U)$ est donné par :

$$\forall \theta \in \Theta, \quad b_\theta(U) = \sum_{k=n}^{+\infty} \mathbf{P} \left(\left[\sum_{i=1}^n X_i = k \right] \right) \int_{1/\theta}^{k/n} \left(\frac{k}{n} - t \right) \frac{2}{t^3} dt.$$

- Justifier que $b_\theta(U)$ est strictement positif, quelle que soit la valeur du paramètre θ

7. Dans cette question, on suppose que le nombre des observations est illimité. On dispose donc, pour estimer le paramètre θ , d'une suite $(X_n)_{n \in \mathbb{N}^*}$ de variables aléatoires mutuellement indépendantes et de même loi. Pour tout entier $n \in \mathbb{N}^*$, on note

$$T_n = \frac{1}{n} \sum_{i=1}^n \frac{1}{X_i} \quad \text{et} \quad U_n = \frac{n}{\sum_{i=1}^n X_i}.$$

Étudier la convergence des deux suites d'estimateurs $(T_n)_{n \in \mathbb{N}^*}$ et $(U_n)_{n \in \mathbb{N}^*}$ du paramètre θ .

>> to be continued..

>> Solution p. ??

Python

Exercice 24. ♦ Retour sur la méthode de Monte-Carlo

Soient $(U_i)_{i \in \mathbb{N}}$ des variables aléatoires réelles indépendantes de loi uniforme sur $[0, 1]$ et N une variable aléatoire de loi géométrique de paramètre p indépendante de la suite $(U_i)_{i \in \mathbb{N}}$. On pose

$$X = \max_{1 \leq i \leq N} U_i.$$

1. Déterminer la fonction de répartition de la variable aléatoire X .
2. Calculer l'espérance de X .
3. Simuler la variable et vérifier votre résultat.

>> Solution p. 29

Exercice 25. ♦ Soit (X_1, \dots, X_n) un échantillon suivant la loi $\mathcal{E}(\lambda)$. Nous disposons des deux estimateurs de $\theta = \frac{1}{\lambda}$:

- La moyenne empirique $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$;
- L'écart-type empirique $\sigma_n = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X}_n)^2}$.

On souhaite comparer ces deux estimateurs.

1. Écrire un programme qui prend en argument n et λ et simule \bar{X}_n , puis σ_n .
2. Afficher des histogrammes et comparer les estimateurs.

>> Solution p. 30

Exercice 26. ♦

Soient $n \in \mathbb{N}^*$ et (X_1, X_2, \dots, X_n) un échantillon d'une loi de Poisson $\mathcal{P}(\theta)$. On cherche à estimer $\exp(-\theta)$. Pour cela, on pose :

$$A_n = \left(1 - \frac{1}{n}\right)^{\sum_{i=1}^n X_i} \quad \text{et} \quad B_n = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{X_i=0\}}.$$

1. Écrire deux programmes qui prend en arguments n , θ et simulent respectivement A_n et B_n .
2. On suppose que $\theta = 1$.
En déduire un programme qui prend en argument n et renvoie une approximation de l'espérance de A_n et B_n .
Que peut-on conjecturer sur le biais de chacun de ces estimateurs?
3. Afficher l'histogramme de 10000 réalisations de T_{100} et U_{100} .
Que peut-on en déduire sur la qualité de ces estimateurs?

>> Solution p. ??

Exercice 27. Dans cet exercice toutes les variables aléatoires sont définies sur un espace probabilisé $(\Omega, \mathcal{A}, \mathbb{P})$.

- 1) Question de cours : Rappeler la définition de la convergence en probabilité d'une suite de variables aléatoires.
- 2) Soit $(u, v) \in \mathbb{R}^2$. On considère deux suites $(U_n)_{n \in \mathbb{N}^*}$ et $(V_n)_{n \in \mathbb{N}^*}$ de variables aléatoires convergeant en probabilité, la première vers u , la seconde vers v .

Démontrer que la suite de variables aléatoires $(U_n + V_n)_{n \in \mathbb{N}^*}$ converge en probabilité vers $u + v$.

- 3) Dans cette question, on suppose les réels u et v supérieurs ou égaux à 1.

a) Établir l'inégalité : $|\ln(u) - \ln(v)| \leq |u - v|$.

b) En déduire que si $(U_n)_{n \in \mathbb{N}^*}$ est une suite de variables aléatoires ne prenant que des valeurs supérieures ou égales à 1, convergeant en probabilité vers u , la suite de variables aléatoires $(\ln(U_n))_{n \in \mathbb{N}^*}$ converge en probabilité vers $\ln(u)$.

- 4) Dans cette question, x désigne un nombre réel strictement compris entre 0 et 1.

a) Établir, pour tout $t \in [0, x]$, l'inégalité : $\frac{x-t}{1-t} \leq x$.

b) En déduire la limite de l'intégrale $\int_0^x \left(\frac{x-t}{1-t}\right)^n dt$ quand l'entier n tend vers l'infini.

c) Démontrer l'égalité : $\ln\left(\frac{1}{1-x}\right) = \sum_{k=1}^{+\infty} \frac{x^k}{k}$.

5) Dans cette question, p désigne un paramètre strictement compris entre 0 et 1. Soit X une variable aléatoire discrète telle que, pour tout entier k strictement positif, la probabilité que la variable X prenne la valeur k soit donnée par : $\mathbb{P}[X = k] = \frac{(1-p)^k}{k \ln\left(\frac{1}{p}\right)}$. Le résultat de la question précédente, appliqué à $x = 1 - p$, prouve que X ne peut pas prendre d'autres valeurs.

- a) Établir l'égalité : $\frac{\mathbb{E}[X]}{\mathbb{E}[X^2]} = p$, où \mathbb{E} désigne l'espérance.

b) Démontrer que si $(X_n)_{n \in \mathbb{N}^*}$ est une suite de variables aléatoires indépendantes et de même loi que X , alors

$$\widehat{p}_n(X_1, X_2, \dots, X_n) = \frac{\sum_{k=1}^n X_k}{\sum_{k=1}^n (X_k)^2} \text{ est un estimateur convergent du paramètre } p.$$



Indications et solutions



Exercice 2

p. 5

1. Soit $n \in \mathbb{N}^*$. Par linéarité de l'espérance

$$\begin{aligned} \mathbf{E}_\theta(T_n) &= \frac{1}{n} \sum_{k=1}^n \mathbf{E}_\theta \left((X_k - \mu)^2 \right) \\ &= \frac{1}{n} \sum_{k=1}^n \mathbf{E}_\theta \left((X_k - \mathbf{E}_\theta(X_k))^2 \right) \\ &= \frac{1}{n} \sum_{k=1}^n \sigma^2 = \sigma^2. \end{aligned}$$

On a donc bien un estimateur sans biais.

2.(a) Par linéarité de l'espérance

$$\mathbf{E}_\theta(V_n) = \frac{1}{n} \sum_{k=1}^n \mathbf{E}_\theta \left((X_k - \bar{X}_n)^2 \right).$$

Soit $k \in \llbracket 1; n \rrbracket$

$$(X_k - \bar{X}_n)^2 = X_k^2 - 2X_k\bar{X}_n + \bar{X}_n^2.$$

Calculons l'espérance de chacun des termes. Par la formule de Koenig-Huygens

$$\mathbf{E}_\theta(X_k^2) = \mathbf{V}_\theta(X_k) + \mathbf{E}_\theta(X_k)^2 = \sigma^2 + \mu^2$$

$$\mathbf{E}_\theta(\bar{X}_n^2) = \mathbf{V}_\theta(\bar{X}_n) + \mathbf{E}_\theta(\bar{X}_n)^2 = \frac{\sigma^2}{n} + \mu^2.$$

De plus
$$X_k\bar{X}_n = \frac{1}{n}X_k^2 + \frac{1}{n} \sum_{\substack{i=1 \\ i \neq k}}^n X_k X_i.$$

Et

$$\begin{aligned} \mathbf{E}_\theta(X_k\bar{X}_n) &= \frac{1}{n} \mathbf{E}_\theta(X_k^2) + \frac{1}{n} \sum_{\substack{i=1 \\ i \neq k}}^n \mathbf{E}_\theta(X_k X_i) \\ &= \frac{1}{n} \mathbf{E}_\theta(X_k^2) + \frac{1}{n} \sum_{\substack{i=1 \\ i \neq k}}^n \mathbf{E}_\theta(X_k) \mathbf{E}_\theta(X_i) \quad (\text{indépendance}) \\ &= \frac{1}{n} (\sigma^2 + \mu^2) + \frac{1}{n} \sum_{\substack{i=1 \\ i \neq k}}^n \mu^2 \\ &= \frac{1}{n} (\sigma^2 + \mu^2) + \frac{n-1}{n} \mu^2 \\ &= \frac{1}{n} \sigma^2 + \mu^2. \end{aligned}$$

Résumons

$$\begin{aligned} \mathbf{E}_\theta \left((X_k - \bar{X}_n)^2 \right) &= \sigma^2 + \mu^2 - 2 \left(\frac{1}{n} \sigma^2 + \mu^2 \right) + \frac{\sigma^2}{n} + \mu^2 \\ &= \sigma^2 \left(1 - \frac{1}{n} \right) = \sigma^2 \frac{(n-1)}{n} \end{aligned}$$

Finalement

$$b_\theta(V_n) = -\frac{1}{n} \sigma^2 \xrightarrow{n \rightarrow +\infty} 0.$$

Ce qui conclut.

2.(b) Il suffit de poser

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

de sorte que $b_\theta(\tilde{V}_n) = 0$.

Exercice 3

p. 5

1. La variable X admet un moment d'ordre et les variables constituant \bar{X}_n sont mutuellement indépendantes. D'après la loi faible des grands nombres,

$$\bar{X}_n \xrightarrow[n \rightarrow +\infty]{\mathbf{P}} \mathbf{E}_\theta(X) = \frac{\theta}{2}.$$

Puis, par composition avec $t \mapsto 2t$, continue

$$T_n = 2\bar{X}_n \xrightarrow[n \rightarrow +\infty]{\mathbf{P}} \theta.$$

2.(a) En tant que variable bornée, X admet une espérance (pour toute probabilité \mathbf{P}_θ). De plus

$$\begin{aligned} \mathbf{E}_\theta(X) &= \int_0^\theta t f_\theta(t) dt = \frac{2}{\theta} \int_0^\theta \left(t - \frac{t^2}{\theta} \right) dt \\ &= \frac{2}{\theta} \left[\frac{t^2}{2} - \frac{t^3}{3\theta} \right]_0^\theta = \frac{2}{\theta} \left(\frac{\theta^2}{2} - \frac{\theta^3}{3\theta} \right) = \frac{\theta}{3}. \end{aligned}$$

Ensuite, par linéarité de l'espérance

$$\mathbf{E}_\theta(3\bar{X}_n) = \frac{3}{n} \sum_{i=1}^n \mathbf{E}_\theta(X_i) = \frac{3}{n} \sum_{i=1}^n \frac{\theta}{3} = \theta.$$

On a bien un estimateur sans biais.

2.(b) La variable X admet une espérance. Par la loi faible des grands nombres

$$\bar{X}_n \xrightarrow[n \rightarrow +\infty]{\mathbf{P}} \mathbf{E}_\theta(X) = \frac{\theta}{3}.$$

Puis, par composition avec $t \mapsto 3t$, continue

$$3\bar{X}_n \xrightarrow[n \rightarrow +\infty]{\mathbf{P}} \theta.$$

L'estimateur $3\bar{X}_n$ est bien convergent.

Exercice 4

p. 6

1. Avec le changement de variable \mathcal{C}^1 et bijectif sur \mathbb{R}_*^+ , $u = \frac{t^2}{2\sigma^2}$, on obtient

$$\begin{aligned} \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{+\infty} |t| e^{-\frac{t^2}{2\sigma^2}} dt &= \frac{2}{\sigma\sqrt{2\pi}} \int_0^{+\infty} |t| e^{-\frac{t^2}{2\sigma^2}} dt \\ &= 2 \frac{\sigma}{\sqrt{2\pi}} \int_0^{+\infty} e^{-u} du = \sigma \sqrt{\frac{2}{\pi}}. \end{aligned}$$

Ce calcul prouve la convergence (absolue). Pour tout indice i , en utilisant le théorème de transfert, $|X_i|$ admet une espérance avec

$$\mathbf{E}(|X_i|) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{+\infty} |t| e^{-\frac{t^2}{2\sigma^2}} dt = \sigma \sqrt{\frac{2}{\pi}}.$$

2. Par linéarité de l'espérance,

$$\mathbf{E}\left(\sqrt{\frac{\pi}{2}} |X_i|\right) = \sigma.$$

Il en résulte, compte tenu des propriétés de la moyenne empirique, que

$$T_n = \sqrt{\frac{\pi}{2}} \left(\frac{1}{n} \sum_{i=1}^n |X_i| \right) = \sqrt{\frac{\pi}{2}} \bar{X}_n.$$

On vérifie que X_i a une variance et par la loi faible des grands nombres, $(T_n)_{n \in \mathbb{N}^*}$ converge en probabilité vers σ . C'est donc bien une suite d'estimateurs convergente.

Exercice 5

p. 6

• Calculons $\mathbf{E}(T_n)$

Méthode 1.

Remarquons que par indépendance

$$\bar{X}_n = \frac{1}{n} Y_n, \quad \text{où } Y_n = \sum_{k=1}^n X_k$$

suit la loi de Poisson de paramètre $n\lambda$.

Ensuite,

$$\begin{aligned} \sum \exp\left(-\frac{k}{n}\right) \mathbf{P}(Y_n = k) &= \sum \exp\left(-\frac{k}{n}\right) \frac{(n\lambda)^k e^{-n\lambda}}{k!} \\ &= e^{-n\lambda} \sum \frac{(n\lambda e^{-\frac{1}{n}})^k}{k!}. \end{aligned}$$

On reconnaît une série exponentielle. La série est absolument convergente. Le théorème de transfert justifie l'existence de l'espérance et le calcul :

$$\begin{aligned} \mathbf{E}(T_n) &= \mathbf{E}\left(\exp\left(-\frac{1}{n} Y_n\right)\right) \\ &= e^{-n\lambda} \sum_{k=0}^{+\infty} \frac{(n\lambda e^{-\frac{1}{n}})^k}{k!} = e^{-n\lambda} e^{n\lambda e^{-\frac{1}{n}}} = e^{-n\lambda(1-e^{-\frac{1}{n}})}. \end{aligned}$$

Cela prouve que T_n n'est pas un estimateur sans biais de $e^{-\lambda}$ puisque

$$\mathbf{E}(T_n) \neq e^{-\lambda}.$$

Méthode 2.

On vérifie par le théorème de transfert que $e^{-X_1/n}$ admet une espérance donnée par

$$\mathbf{E}\left(e^{-X_1/n}\right) = \sum_{k=0}^{+\infty} \left(e^{-k/n}\right) \cdot e^{-\lambda} \frac{\lambda^k}{k!} = e^{-\lambda} \cdot e^{\lambda e^{-1/n}}.$$

Par le lemme des coalitions, les variables $e^{-X_i/n}$ sont indépendantes. La fonction $t \mapsto e^{-\lambda t}$ étant continue, elles sont aussi même loi. Dès lors,

$$\mathbf{E}(T_n) = \prod_{i=1}^n \mathbf{E}\left(e^{-X_i/n}\right).$$

On retrouve le résultat de la première méthode.

• Mais on sait que

$$\left(1 - e^{-\frac{1}{n}}\right) \sim \frac{1}{n}$$

et donc

$$-n\lambda \left(1 - e^{-\frac{1}{n}}\right) \xrightarrow{n \rightarrow +\infty} -\lambda,$$

et donc que

$$\mathbf{E}(T_n) \xrightarrow{n \rightarrow +\infty} e^{-\lambda}.$$

• Enfin, l'application $f : t \in \mathbb{R} \mapsto e^{-t}$ est continue et l'estimateur de la moyenne empirique \bar{X}_n est un estimateur convergent de $\mathbf{E}(X) = \lambda$. Donc $t_n = f(\bar{X}_n)$ est encore un estimateur convergent de $e^{-\lambda}$.

Exercice 6

p. 7

Soit $\epsilon \in \mathbb{R}_*^+$. À partir de l'inégalité de Markov

$$\mathbf{P}_\theta\left(|T_n - g(\theta)| \geq \epsilon\right) = \mathbf{P}_\theta\left((T_n - g(\theta))^2 \geq \epsilon^2\right) \leq \frac{\mathbf{E}\left((T_n - g(\theta))^2\right)}{\epsilon^2}.$$

Or

$$\begin{aligned} \mathbf{E}\left((T_n - g(\theta))^2\right) &= \mathbf{E}\left((T_n - \mathbf{E}(T_n) + \mathbf{E}(T_n) - g(\theta))^2\right) \\ &= \mathbf{E}\left((T_n - \mathbf{E}(T_n))^2\right) \\ &\quad + 2\mathbf{E}(T_n - \mathbf{E}(T_n)) (\mathbf{E}(T_n) - g(\theta)) \\ &\quad + (\mathbf{E}(T_n) - g(\theta))^2 \\ &= \mathbf{V}(T_n) + (\mathbf{E}(T_n) - g(\theta))^2. \end{aligned}$$

Par hypothèse

$$(\mathbf{E}(T_n) - g(\theta))^2 \xrightarrow{n \rightarrow +\infty} 0 \quad \text{et} \quad \mathbf{V}(T_n) \xrightarrow{n \rightarrow +\infty} 0.$$

D'où $\mathbf{E}\left((T_n - g(\theta))^2\right) \xrightarrow{n \rightarrow +\infty} 0$ et par encadrement

$$\mathbf{P}_\theta\left(|T_n - g(\theta)| \geq \epsilon\right) \xrightarrow{n \rightarrow +\infty} 0.$$

On retrouve la définition de la convergence de l'estimateur.

Exercice 7

p. 7

1. On vérifie que

$$\mathbf{E}_\theta(\bar{X}_n) = p \quad \text{et} \quad \mathbf{E}_\theta(T_n) = p.$$

2. On a par indépendance

$$V_{\theta}(\bar{X}_n) = \frac{1}{n} V(X) = \frac{p(1-p)}{n}$$

et

$$V_{\theta}(T_n) = \frac{4}{n^2(n+1)^2} \sum_{i=1}^n i^2 V(X_i) = \frac{2(2n+1)}{3(n+1)} \frac{p(1-p)}{n}.$$

On a

$$V_{\theta}(\bar{X}_n) \geq V_{\theta}(\bar{T}_n).$$

On dira que \bar{T}_n est un "meilleur" estimateur que \bar{X}_n .

3. On a

$$V_{\theta}(\bar{X}_n) \xrightarrow{n \rightarrow +\infty} 0 \quad \text{et} \quad V_{\theta}(T_n) \xrightarrow{n \rightarrow +\infty} 0.$$

On est dans les conditions de la proposition précédente et les estimateurs sont convergents.

Exercice 8

p. 12

1. Notons m , la moyenne.

- Via l'inégalité de Bienaymé-Tchebychev.

On a vu que

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

est un estimateur sans biais et convergent de m . On va donc utiliser cet estimateur, via l'inégalité de Bienaymé-Tchebychev, pour déterminer un intervalle de confiance de m , au niveau de risque $\alpha = 0,05$.

$$\forall \varepsilon \in \mathbb{R}^+, \quad \mathbf{P}\left(|\bar{X}_n - m| \geq \varepsilon\right) \leq \frac{V(\bar{X}_n)}{\varepsilon^2}$$

$$\Leftrightarrow \mathbf{P}\left(|\bar{X}_n - m| < \varepsilon\right) \geq 1 - \frac{V(\bar{X}_n)}{\varepsilon^2}.$$

Or :

$$\mathbf{P}\left(\bar{X}_n - \varepsilon \leq m \leq \bar{X}_n + \varepsilon\right) \geq \mathbf{P}\left(\bar{X}_n - \varepsilon < m < \bar{X}_n + \varepsilon\right),$$

et si on a

$$\frac{V(\bar{X}_n)}{\varepsilon^2} \leq 0,05,$$

alors :

$$\mathbf{P}\left(\bar{X}_n - \varepsilon \leq m \leq \bar{X}_n + \varepsilon\right) \geq 0,95.$$

Donc pour des valeurs de ε tels que $\frac{\sigma^2}{n\varepsilon^2} \leq 0,05$ i.e. :

$$\varepsilon \geq \sqrt{\frac{\sigma^2}{n0,05}},$$

on obtient l'intervalle de confiance

$$\left[\bar{X}_n - \varepsilon, \bar{X}_n + \varepsilon\right]$$

est un intervalle de confiance de m à 95%.

Application numérique. On trouve

$$[53,75; 56,89].$$

- Via la loi normale.

Les variables X_1, \dots, X_n sont indépendantes de même loi normale de paramètres (m, σ^2) . Par stabilité de la loi normale par combinaison linéaire, \bar{X}_n^* suit la loi $\mathcal{N}(0, 1)$. Par symétrie de la densité de la loi normale centrée réduite, il existe un unique réel strictement positif t_{α} tel que :

$$\Phi(t_{\alpha}) = 1 - \frac{\alpha}{2} = 0,975.$$

$$\text{Alors} \quad \mathbf{P}\left(-t_{\alpha} \leq \bar{X}_n^* \leq t_{\alpha}\right) = 1 - \alpha$$

c'est-à-dire,

$$\mathbf{P}\left(\bar{X}_n - t_{\alpha} \frac{\sigma}{\sqrt{n}} \leq m \leq \bar{X}_n + t_{\alpha} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha.$$

Application numérique. On trouve comme intervalle de confiance de m à 95% :

$$[54,634; 56,006]$$

2. La meilleure méthode est celle dont l'intervalle de confiance possède la plus petite amplitude. La 1^{re} méthode donne un intervalle de confiance d'amplitude égale à 3,14 et la 2^e méthode donne un intervalle de confiance d'amplitude égale à 1,372. Donc la deuxième méthode est plus performante.

Exercice 10

p. 15

1. On a par linéarité de l'espérance

$$\mathbf{E}(T_n) = \frac{2}{n} \sum_{i=1}^n \mathbf{E}(X_i) = \frac{2}{n} \sum_{i=1}^n \frac{\theta}{2} = \theta.$$

Ainsi

$$b_{\theta}(T_n) = \mathbf{E}(T_n) - \theta = 0.$$

L'estimateur est sans biais.

2. En reprenant le cours sur les variables à densité

$$\forall t \in \mathbb{R}, \quad F_{\min}(t) = 1 - (1 - F_X(t))^n$$

$$F_{\max}(t) = F_X(t)^n.$$

On vérifie que F_{\min} et F_{\max} sont continues et \mathcal{C}^1 sauf éventuellement en 0 et 1. Y_{\min} et Y_{\max} sont à densité avec

$$f_{\min}(t) = n f_X(t) (1 - F_X(t))^{n-1}$$

$$f_{\max}(t) = n f_X(t) F_X(t)^{n-1}$$

La variable Y_{\max} est une variable bornée, elle admet donc une espérance et une variance avec

$$\mathbf{E}(Y_{\max}) = \int_0^{\theta} t f_{\max}(t) dt = n \int_0^{\theta} \frac{t}{\theta} \cdot \left(\frac{t}{\theta}\right)^{n-1} dt = \frac{\theta n}{n+1}.$$

Par la formule de Koenig- Huygens

$$V(Y_{\max}) = \mathbf{E}(Y_{\max}^2) - \mathbf{E}(Y_{\max})^2.$$

Or

$$\begin{aligned} \mathbf{E}(Y_{\max}^2) &= \int_0^\theta t^2 f_{\max}(t) dt \\ &= n \int_0^\theta \frac{t^2}{\theta} \cdot \left(\frac{t}{\theta}\right)^{n-1} dt \\ &= \frac{n}{\theta^{n-1}} \int_0^\theta t^{n+1} dt \\ &= \frac{n}{\theta^n} \left[\frac{1}{n+2} t^{n+2} \right]_0^\theta \\ &= \frac{n}{n+2} \frac{\theta^{n+2}}{\theta^n} = \frac{n}{n+2} \theta^2. \end{aligned}$$

D'où

$$\begin{aligned} \mathbf{V}(Y_{\max}) &= \frac{n}{n+2} \theta^2 - \left(\frac{n\theta}{n+1}\right)^2 \\ &= n\theta^2 \left(\frac{1}{n+2} - \frac{n}{(n+1)^2} \right) \\ &= n\theta^2 \left(\frac{(n+1)^2 - n(n+2)}{(n+2)(n+1)^2} \right) \\ &= \frac{n}{(n+2)(n+1)^2} \theta^2. \end{aligned}$$

Par symétrie Y_{\min} et $\theta - Y_{\max}$ ont même loi donc

$$\mathbf{E}(Y_{\min}) = \theta - \mathbf{E}(Y_{\max}) = \theta - \frac{n}{n+1} \theta = \frac{\theta}{n+1}.$$

Et

$$\mathbf{V}(Y_{\min}) = \mathbf{V}(\theta - Y_{\max}) = \mathbf{V}(Y_{\max}).$$

3. On a

$$\mathbf{E}(T'_n) = \frac{n+1}{n} \mathbf{E}(Y_{\max}) = \frac{n+1}{n} \cdot \frac{\theta n}{n+1} = \theta.$$

On a bien un estimateur sans biais.

4. On a deux estimateurs sans biais avec

$$\mathbf{V}(T'_n) = \left(\frac{n+1}{n}\right)^2 \mathbf{V}(Y_{\max}) = \frac{1}{n(n+2)} \theta^2$$

$$\text{et } \mathbf{V}(T_n) = \frac{\theta^2}{3n}.$$

On constate que

$$\mathbf{V}(T'_n) \leq \mathbf{V}(T_n).$$

T'_n est un meilleur estimateur que T_n .

5. On a

$$\begin{aligned} \mathbf{V}(T''_n) &= \mathbf{V}(Y_{\min} + Y_{\max}) \\ &= \mathbf{V}(Y_{\min}) + \mathbf{V}(Y_{\max}) + 2 \text{Cov}(Y_{\min}, Y_{\max}). \end{aligned}$$

Or

$$\mathbf{V}(Y_{\min}) = \mathbf{V}(Y_{\max})$$

et en remarquant que

$$0 \leq \mathbf{V}(Y_{\min} - Y_{\max}) = \mathbf{V}(Y_{\min}) + \mathbf{V}(Y_{\max}) - 2 \text{Cov}(Y_{\min}, Y_{\max})$$

on a

$$2 \text{Cov}(Y_{\min}, Y_{\max}) \leq \mathbf{V}(Y_{\min}) + \mathbf{V}(Y_{\max})$$

On obtient

$$\mathbf{V}(T''_n) \leq 4 \mathbf{V}(Y_{\max}).$$

6. Dans un premier temps. T''_n est sans biais

$$\mathbf{E}(T''_n) = \mathbf{E}(Y_{\min}) + \mathbf{E}(Y_{\max}) = \theta.$$

On sait par indépendance que

$$\mathbf{V}(T_n) = \frac{4}{n^2} \sum_{i=1}^n \mathbf{V}(X_i) = \frac{4}{n} \cdot \frac{\theta^2}{12} = \frac{\theta^2}{3n}.$$

Et on a vu que

$$\mathbf{V}(Y_{\max}) = \frac{n}{(n+2)(n+1)^2} \theta^2.$$

D'après ce qui précède

$$\begin{aligned} \mathbf{V}(T''_n) &\leq \frac{4n}{(n+2)(n+1)^2} \sigma^2 = \frac{12n^2}{(n+2)(n+1)^2} \mathbf{V}(T_n) \\ &\leq \frac{12}{n+2} \mathbf{V}(T_n). \end{aligned}$$

Dès que $n \geq 10$, on a

$$\mathbf{V}(T''_n) \leq \mathbf{V}(T_n)$$

et T_n est un meilleur estimateur de T_n .

7. Dans chaque cas, les estimateurs sont sans biais avec une variance qui tend vers 0. Par conséquent, on sait que chaque estimateur est convergent.

Exercice 14

p. 16

1. La variable Y_k est à valeurs dans \mathbb{R}^+ et pour $t \in \mathbb{R}^+$

$$\begin{aligned} F_{Y_k}(t) &= \mathbf{P}(Y_k \leq t) \\ &= \mathbf{P}(X_k \leq bt + \theta) = F_{X_k}(bt + \theta). \end{aligned}$$

Par composition, la fonction de répartition de Y_k est continue et de classe \mathcal{C}^1 sauf éventuellement en un nombre fini de point. On en déduit que Y_k est une variable à densité. Une densité est donnée par

$$\forall t \in \mathbb{R} \quad f(t) = \begin{cases} 0 & \text{si } t \in \mathbb{R}_*^- \\ e^{-t} & \text{si } t \in \mathbb{R}^+ \end{cases}$$

On reconnaît une loi exponentielle de paramètre 1.

2. Par linéarité et l'égalité en loi des Y_i

$$\mathbf{E}(\overline{Y}_n) = \mathbf{E}(Y_1) = 1.$$

On montre que U_n suit une loi exponentielle de paramètre n . Dès lors

$$\mathbf{E}(U_n) = \frac{1}{n}.$$

3. On a

$$\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} \sum_{i=1}^n (bY_i + \theta) = b\overline{Y}_n + \theta$$

et

$$\begin{aligned} T_n &= \min_{i \in \llbracket 1; n \rrbracket} X_i = \min_{i \in \llbracket 1; n \rrbracket} bY_i + \theta \\ &= b \left(\min_{i \in \llbracket 1; n \rrbracket} Y_i \right) + \theta \\ &= bU_n + \theta. \end{aligned}$$

Par croissance de la fonction affine $t \in \mathbb{R} \mapsto bt + \theta$. Par linéarité

$$\mathbf{E}(\overline{X}_n) = b + \theta, \quad \mathbf{E}(T_n) = \frac{b}{n} + \theta.$$

4. On constate que

$$\mathbf{E}(\overline{X}_n - T_n) = \left(1 - \frac{1}{n}\right)b = \frac{n-1}{n}b.$$

Un estimateur sans biais de b est donc donné par

$$\hat{b}_n = \frac{n}{n-1}(\overline{X}_n - T_n).$$

Dans ce cas

$$\begin{aligned}\mathbf{E}(\overline{X}_n - \hat{b}_n) &= \mathbf{E}(\overline{X}_n) - \mathbf{E}(\hat{b}_n) \\ &= b + \theta - b = \theta.\end{aligned}$$

et un estimateur sans biais de θ est

$$\hat{\theta}_n = \overline{X}_n - \hat{b}_n = \frac{n}{n-1}T_n - \frac{1}{n-1}\overline{X}_n.$$

Exercice 17

p. 18

1. Vérifier que

$$Y_n \hookrightarrow \mathcal{E}(n\lambda).$$

2. Par les règles sur les transformations linéaires, $n\lambda Y_n \hookrightarrow \mathcal{E}(1)$.

Par conséquent

$$\mathbf{P}(n\lambda Y_n \leq a_n) = F(a_n)$$

où F désigne la fonction de répartition d'une loi exponentielle de paramètre 1. Soit

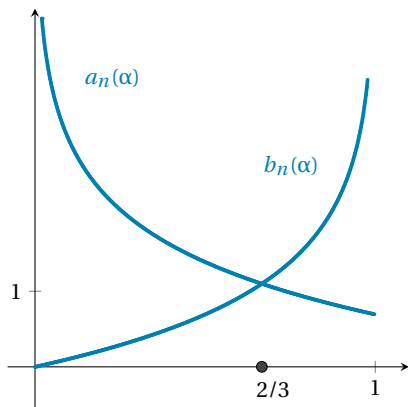
$$\forall x \in \mathbb{R}, \quad F(x) = \begin{cases} 0 & \text{si } x < 0 \\ 1 - e^{-x} & \text{sinon.} \end{cases}$$

Il vient pour $a_n \geq 0$,

$$\begin{aligned}\frac{\alpha}{2} &= F(a_n) = 1 - e^{-a_n} \\ \Leftrightarrow e^{-a_n} &= 1 - \frac{\alpha}{2} \\ \Leftrightarrow a_n &= -\ln\left(1 - \frac{\alpha}{2}\right).\end{aligned}$$

Sachant que Y_n est une variable à densité :

$$\begin{aligned}\mathbf{P}(n\lambda Y_n \geq b_n) &= 1 - F(b_n) \\ \Leftrightarrow 1 - F(b_n) &= \frac{\alpha}{2} \\ \Leftrightarrow e^{-b_n} &= \frac{\alpha}{2} \\ \Leftrightarrow b_n &= -\ln\left(\frac{\alpha}{2}\right).\end{aligned}$$



Notons que pour $\alpha \in [0; 2/3]$,

$$0 < b_n \leq a_n.$$

3. Posons

$$A_n = [nY_n\lambda \leq a_n]$$

$$B_n = [nY_n\lambda \geq b_n].$$

Ainsi

$$\begin{aligned}\mathbf{P}\left(\frac{1}{\lambda} \in \left[\frac{nY_n}{a_n}, \frac{nY_n}{b_n}\right]\right) \\ = \mathbf{P}\left(\frac{nY_n}{a_n} \leq \frac{1}{\lambda} \leq \frac{nY_n}{b_n}\right) \\ = \mathbf{P}([nY_n\lambda \leq a_n] \cap [b_n \leq nY_n\lambda]) \\ = \mathbf{P}(A_n \cap B_n) = 1 - \mathbf{P}(\bar{A}_n \cup \bar{B}_n).\end{aligned}$$

En utilisant la formule du crible, on a aussi

$$\mathbf{P}(\bar{A}_n \cup \bar{B}_n) \leq \mathbf{P}(\bar{A}_n) + \mathbf{P}(\bar{B}_n).$$

Comme Y_n est une variable à densité, le membre de droite vaut α . Finalement

$$\mathbf{P}\left(\frac{1}{\lambda} \in \left[\frac{nY_n}{a_n}, \frac{nY_n}{b_n}\right]\right) \geq 1 - \alpha$$

Exercice 21

p. 19

Il est clair que

$$\mathbf{E}(\overline{X}_n) = p, \quad \mathbf{E}(\overline{Y}_m) = p.$$

D'où

$$\mathbf{E}(Z) = (a + b)p.$$

Z est donc un estimateur sans biais si et seulement si

$$a + b = 1.$$

De plus, par le lemme des coalitions, \overline{X}_n et \overline{Y}_m sont indépendantes. Ainsi

$$\begin{aligned}\mathbf{V}(Z) &= a^2 \mathbf{V}(\overline{X}_n) + b^2 \mathbf{V}(\overline{Y}_m) \\ &= \left(\frac{a^2}{n} + \frac{b^2}{m}\right)p(1-p) \\ &= \left(\frac{a^2}{n} + \frac{(1-a)^2}{m}\right)p(1-p).\end{aligned}$$

Le meilleur estimateur sans biais étant celui de variance minimale, on cherche donc à minimiser la fonction polynomiale de degré 2

$$a \in \mathbb{R} \mapsto \frac{a^2}{n} + \frac{(1-a)^2}{m}.$$

Une étude de variation donne un minimum atteint pour

$$a = \frac{n}{n+m}.$$

Dans ce cas

$$Z = \frac{n}{n+m}\overline{X}_n + \frac{m}{n+m}\overline{Y}_m = \frac{1}{n+m}\left(\sum_{i=1}^n X_i + \sum_{i=1}^m Y_i\right).$$

C'est tout simplement l'estimateur de la moyenne empirique construit à partir de l'échantillon

$$(X_1, \dots, X_n, Y_1, \dots, Y_m).$$

Exercice 24

p. 22

1. Soit $x \in [0; 1]$. Appliquons la formule des probabilités totales avec le système complet d'événements $(N = n)_{n \in \mathbb{N}}$

$$\begin{aligned} F(x) &= \mathbf{P}(X \leq x) \\ &= \sum_{n=1}^{+\infty} \mathbf{P}(N = n) \mathbf{P}_{[N=n]}(X \leq x) \\ &= \sum_{n=1}^{+\infty} \mathbf{P}(N = n) \mathbf{P}_{[N=n]}(\max\{U_1 \dots U_n\} \leq x) \\ &= \sum_{n=1}^{+\infty} \mathbf{P}(N = n) \mathbf{P}(\max\{U_1, \dots, U_n\} \leq x). \end{aligned}$$

La dernière égalité s'obtenant par indépendance de N avec $\max\{U_1, \dots, U_n\}$. Ensuite

$$\begin{aligned} F(x) &= \sum_{n=1}^{+\infty} \mathbf{P}(N = n) \mathbf{P}([U_1 \leq x] \cap \dots \cap [U_n \leq x]) \\ &= \sum_{n=1}^{+\infty} \mathbf{P}(N = n) \mathbf{P}(U_1 \leq x)^n \end{aligned}$$

par indépendance de $(U_i)_{i \in \mathbb{N}^*}$ et même loi. En remplaçant par les expressions des lois :

$$\begin{aligned} F(x) &= \sum_{n=1}^{+\infty} p(1-p)^{n-1} \cdot x^n \\ &= px \sum_{n=1}^{+\infty} ((1-p)x)^{n-1} = \frac{px}{1 - (1-p)x}. \end{aligned}$$

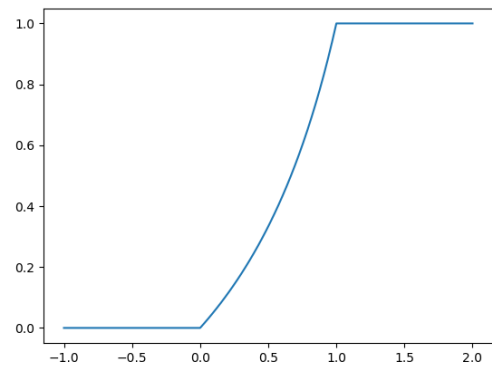
Pour $x \leq 0$, $F(x) = 0$, et $x \geq 1$, $F(x) = 1$.

Voici un code pour tracer la courbe représentative de la fonction de répartition.

```
def FctionRep(x, p):
    if x < 0 :
        return 0
    if x > 1:
        return 1
    else :
        return p*x / (1 - (1-p)*x)

def trace(p):
    x=np.linspace(-1, 2, 100)
    y=np.zeros(len(x))
    for i in range(len(x)):
        y[i]=FctionRep(x[i], p)
    plt.clf()
    plt.plot(x, y)
    plt.show()

# test
trace(1/2)
```



2. Appliquons la formule de l'espérance totale :

- $(N = n)_{n \in \mathbb{N}}$ est un système complet d'événements.
- Soit $n \in \mathbb{N}^*$. Posons $\bar{U}_n = \max(U_1, \dots, U_n)$. En reprenant le calcul précédent. Pour $x \in [0; 1]$

$$F_{\bar{U}_n}(x) = x^n.$$

\bar{U}_n est à valeurs dans $[0, 1]$, et une densité est :

$$\forall x \in \mathbb{R}, \quad f_n(x) = \begin{cases} nx^{n-1} & \text{si } x \in [0; 1] \\ 0 & \text{sinon.} \end{cases}$$

\bar{U}_n est bornée, elle admet une espérance avec

$$\mathbf{E}(\bar{U}_n) = \int_0^1 x \cdot nx^{n-1} dx = \frac{n}{n+1}.$$

- Ainsi \bar{U}_n est positive et

$$\mathbf{E}(\bar{U}_n) = \mathbf{E}(|X| | [N = n]).$$

De plus

$$\sum \mathbf{P}(N = n) \mathbf{E}(|X| | [N = n]) = \sum pq^{n-1} \cdot \frac{n}{n+1}.$$

La série est convergente.

D'après la formule de l'espérance totale, X admet une espérance et :

$$\begin{aligned} \mathbf{E}(X) &= \sum_{n=1}^{+\infty} \mathbf{P}(N = n) \mathbf{E}(X | [N = n]) \\ &= \sum_{n=1}^{+\infty} pq^{n-1} \cdot \frac{n}{n+1} \\ &= \sum_{n=1}^{+\infty} pq^{n-1} \left(\frac{n+1-1}{n+1} \right) \\ &= \sum_{n=1}^{+\infty} pq^{n-1} - pq^{-2} \sum_{n=1}^{+\infty} \frac{q^{n+1}}{n+1}. \end{aligned}$$

Or on montre que pour $x \in]-1; 1[$

$$\ln(1-x) = - \sum_{n=1}^{+\infty} \frac{x^n}{n}.$$

D'où

$$\begin{aligned} \mathbf{E}(X) &= 1 - pq^{-2} \sum_{n=2}^{+\infty} \frac{q^n}{n^2} \\ &= 1 - pq^{-2} \left(\sum_{n=1}^{+\infty} \frac{q^n}{n} - q \right) \\ &= 1 + pq^{-2} (\ln(1-q) + q) \\ &= 1 + \frac{p}{(1-p)^2} (\ln(p) + 1 - p). \end{aligned}$$

3.

```
def simuGeometrique(p):
    n=1
    while np.random.rand()>p:
        n=n+1
    return n

def simu(p):
    n=simuGeometrique(p)
    U=np.random.rand(n)
    return max(U)

def approxEsperance(p):
    s=0
    m=500
    for i in range(m):
        s+=simu(p)

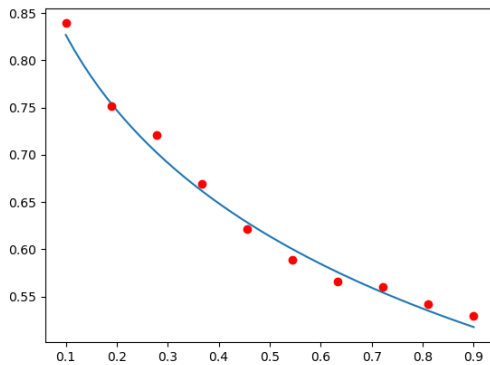
    ValTheo=1+p*(1-p)**(-2)*(np.log(p)+1-p)
    #print(s/1000,'à comparer à',ValTheo)
    return s/m

def Valtheorique(p):
    return 1+p*(1-p)**(-2)*(np.log(p)+1-p)

def comparaison():
    plt.clf()
    x=np.linspace(0.1,0.9,50)
    plt.plot(x,Valtheorique(x))

    p=np.linspace(0.1,0.9,10)
    for i in range(len(p)):
        plt.plot([p[i]], [approxEsperance(p[i])], 'ro')
    plt.show()

comparaison()
```



Exercice 25

1.

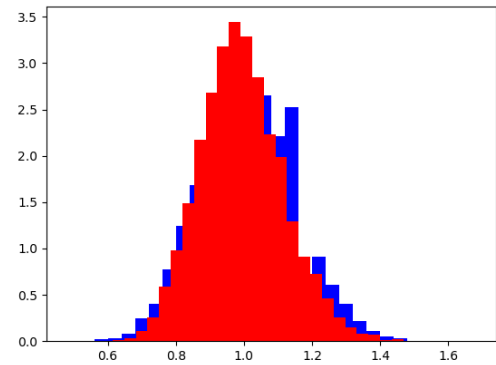
```
def Estimateurs(n, lbda):
    X=np.random.poisson(lbda, n)
    Xbar=np.mean(X)
    Sigma=np.std(X)*(n/(n-1))**(1/2)
    return [Xbar, Sigma]
```

2.

```
def comparaisonPoisson(lbda, n):
    m=10000
    Xbar=np.zeros(m)
    Sigma=np.zeros(m)
    for i in range(m):
        S=Estimateurs(n, lbda)
        Xbar[i]=float(S[0])
        Sigma[i]=float(S[1])

    plt.clf()
    plt.hist(Xbar, 30, density=True, color='blue')
    plt.hist(Sigma, 30, density=True, color='red')
    plt.show()
```

```
comparaisonPoisson(1, 50)
```



Après plusieurs essais, on constate que la dispersion par rapport à la moyenne est plus faible pour σ_n (en rouge). On conjecture que

$$V(\sigma_n) \leq V(\bar{X}_n).$$

L'estimateur σ_n serait meilleur que \bar{X}_n .

1.