

Examen Finale "Statistiques des Processus".

Aucun document n'est autorisé.

Durée : 1h30mn.

16 Janvier 2022.

Exercice 1.

Soit X_1, \dots, X_n une suite de variables aléatoires i.i.d de loi de Bernoulli de paramètre p , $p \in]0, 1[$.

1. Proposer un estimateur \hat{p} de p .
2. En utilisant l'inégalité de Tchébychef, construire un intervalle de confiance pour p de niveau $1 - \alpha$ pour $0 < \alpha < 1$.
3. Dédurre de l'inégalité de Hoeffding, un intervalle de confiance pour p de niveau $1 - \alpha$.
4. Comparer la précision de ces intervalles de confiances pour $n = 100$ et $\alpha = 5\%$.

► **Rappel** : Soit $\alpha \in]0, 1[$, un intervalle de confiance pour un paramètre θ de niveau $1 - \alpha$ est un couple d'estimateurs $(\underline{\theta}_n, \bar{\theta}_n)$ tel que $\mathbb{P}(\theta \in [\underline{\theta}_n, \bar{\theta}_n]) \geq 1 - \alpha$.

► **Indication** : Pour tout $x \in]0, 1[$, $x(1 - x) \leq \frac{1}{4}$.

Exercice 2.

On observe un n-échantillon $(X_1, Y_1), \dots, (X_n, Y_n)$ de même loi qu'un couple de variables aléatoires réelles (X, Y) . On suppose que le vecteur (X, Y) admet pour densité la fonction f telle que pour tout $(x, y) \in \mathbb{R}^2$, $f(x, y) > 0$. On considère \hat{f}_h l'estimateur de f défini pour tout $(x, y) \in \mathbb{R}^2$ par :

$$\hat{f}_h(x, y) = \frac{1}{nh_1h_2} \sum_{i=1}^n K\left(\frac{x - X_i}{h_1}, \frac{y - Y_i}{h_2}\right)$$

où $h = (h_1, h_2)$ avec $h_1 > 0$ et $h_2 > 0$. $K : \mathbb{R}^2 \rightarrow \mathbb{R}$ est un noyau tel que :

$$\iint |v| |K(u, v)| dudv < \infty, \quad \iint |u|^{1/2} |K(u, v)| dudv < \infty \quad \text{et} \quad \|K\|_2^2 < \infty$$

On pose pour tout $(x, y) \in \mathbb{R}^2$: $K_h(x, y) = \frac{1}{h_1h_2} K\left(\frac{x}{h_1}, \frac{y}{h_2}\right)$

1. Montrer que pour $x \in \mathbb{R}$ et $y \in \mathbb{R}$,

$$\mathbb{E}[\hat{f}_h(x, y)] = (K_h * f)(x, y).$$

2. Calculer pour $x \in \mathbb{R}$ et $y \in \mathbb{R}$, $Var\left(\hat{f}_h(x, y)\right)$ en fonction de n , $K_h * f$ et $K_h^2 * f$.
3. Montrer que si f est bornée par M , alors

$$Var\left(\hat{f}_h(x, y)\right) \leq \frac{M\|K\|_2^2}{nh_1h_2}.$$

4. Montrer que si f vérifie de plus la propriété suivante : pour tous $(u, v) \in \mathbb{R}^2$ et $(u', v') \in \mathbb{R}^2$

$$|f(u, v) - f(u', v')| \leq |u - u'|^{1/2} + |v - v'|,$$

alors le biais de $\hat{f}_h(x, y)$ vérifie pour tout $(x, y) \in \mathbb{R}^2$

$$\left| \text{Biais} \left(\hat{f}_h(x, y) \right) \right| \leq C_K (h_1^{1/2} + h_2)$$

où C_K est une constante ne dépendant que de K .

5. Dédurre une borne pour le risque quadratique ponctuel.
6. On note à présent f_X la densité de X . En utilisant \hat{f}_h , proposer un estimateur de f_X . Donner son expression si on suppose qu'il existe deux noyaux F et G telles que pour tous $x \in \mathbb{R}$ et $y \in \mathbb{R}$

$$K(x, y) = F(x)G(y).$$

- **Rappel** : Le produit de convolution des fonctions f_1 et f_2 de deux variables est défini par :

$$\forall (x, y) \in \mathbb{R}^2, (f_1 * f_2)(x, y) = \iint f_1(x - u, y - v) f_2(u, v) du dv$$

Bonne courage.

"Corrigé Epreuve finale"
"Statistique des Processus"

Exercice n°1:

$X_i \sim \mathcal{B}(p)$, $i=1, \dots, n$ avec X_i i.i.d. $p \in]0, 1[$.

1°/ $X_i \sim \mathcal{B}(p)$ donc $E(X_i) = p \quad \forall i=1, \dots, n$

Alors, $\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i$

2°/ Inégalité de Tchébychev:

ou $P(|\bar{X}_n - E(\bar{X}_n)| > \alpha) \leq \frac{\text{Var}(\bar{X}_n)}{\alpha^2}$ pour tout $\alpha \in \mathbb{R}^{++}$.

$E(\bar{X}_n) = p$.

$\text{Var}(\bar{X}_n) = \frac{1}{n^2} \sum_{i=1}^n \text{V}(X_i) = \frac{p(1-p)}{n}$

Donc;

ou $P(|\bar{X}_n - p| > \alpha) \leq \frac{p(1-p)}{4n\alpha^2} \leq \frac{1}{4n\alpha^2}$

$P(|\bar{X}_n - p| \leq \alpha) \geq 1 - \frac{1}{4n\alpha^2}$
 $P(\hat{p} - \alpha \leq p \leq \hat{p} + \alpha) \geq 1 - \frac{1}{4n\alpha^2} = 1 - \alpha$

ou Donc pour $\alpha = \frac{1}{2\sqrt{nd}}$ $\Leftrightarrow \alpha = \frac{1}{2\sqrt{nd}}$, on aura:

$P\left(\hat{p} - \frac{1}{2\sqrt{nd}} \leq p \leq \hat{p} + \frac{1}{2\sqrt{nd}}\right) \geq 1 - \alpha$

3°/ Inégalité de Hoeffding: $S_n = \sum_{i=1}^n X_i$

(0,1) $\forall n > 0, P(|S_n - E(S_n)| > \alpha) \leq 2 \exp\left(-\frac{2\alpha^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$

Comme $X_i \sim \mathcal{B}(p)$ alors $0 \leq X_i \leq 1$.

Donc: $P(|\hat{p} - p| > \frac{\alpha}{n}) \leq 2 e^{-\frac{2\alpha^2}{n}}$

(0,1) $P(\hat{p} - \frac{\alpha}{n} \leq p \leq \hat{p} + \frac{\alpha}{n}) \geq 1 - 2 e^{-\frac{2\alpha^2}{n}} = 1 - \alpha$

Alors pour $\alpha = 2 \exp\left(-\frac{2\alpha^2}{n}\right) \Leftrightarrow \alpha = \sqrt{\frac{n \ln(2/\alpha)}{2}}$

on a:

$P\left(\hat{p} - \sqrt{\frac{\ln(2/\alpha)}{2n}} \leq p \leq \hat{p} + \sqrt{\frac{\ln(2/\alpha)}{2n}}\right) \geq 1 - \alpha$

4°/ $n=100$ et $\alpha=1\%$.

1. C. Tchebychev:

$$P\left(\left|\hat{p} - \frac{1}{2}\right| \leq p \leq \hat{p} + \frac{1}{2}\sqrt{\frac{p(1-p)}{n}}\right) \geq 1 - 0.05 = 95\%$$

Donc la précision de l'intervalle est:

$$2 \cdot \frac{1}{2\sqrt{n}} \approx 0.447$$

(2)

1. C. Hoeffding:

$$2 \cdot \sqrt{\frac{\ln(2/0.05)}{2 \cdot 100}} \approx 0.271$$

On a bien l'inégalité de Hoeffding donne un intervalle de confiance plus précis pour la valeur p.

Exercice n°2:
$$\hat{f}_{h_1, h_2}(x, y) = \frac{1}{n h_1 h_2} \sum_{i=1}^n K\left(\frac{x-x_i}{h_1}, \frac{y-y_i}{h_2}\right)$$

1°) $(x, y) \in \mathbb{R}^2$, $E[\hat{f}_{h_1, h_2}(x, y)] = (K_h * f)(x, y)$

(3)
$$E(\hat{f}_{h_1, h_2}(x, y)) = \frac{1}{n h_1 h_2} \sum_{i=1}^n \int_{\mathbb{R}^2} K\left(\frac{x-u}{h_1}, \frac{y-v}{h_2}\right) f(u, v) du dv$$

$$= \int_{\mathbb{R}^2} K_h(x-u, y-v) f(u, v) du dv$$

$$= (K_h * f)(x, y)$$

2°) $Var(\hat{f}_{h_1, h_2}(x, y)) = \frac{1}{n^2} \sum_{i=1}^n Var(K_h(x-x_i, y-y_i))$

(4)
$$= \frac{1}{n} E(K_h^2(x-x, y-y)) - \frac{1}{n} E^2(K_h(x-x, y-y))$$

$$= \frac{1}{n} (K_h^2 * f)(x, y) - \frac{1}{n} (K_h * f)^2(x, y)$$

3°) $Var(\hat{f}_{h_1, h_2}(x, y)) \leq \frac{M \|K\|_2^2}{n h_1 h_2}$

01
$$Var(\hat{f}_{h_1, h_2}(x, y)) \leq \frac{1}{n} (K_h^2 * f)(x, y)$$

$$\leq \frac{M}{n h_1 h_2} \int_{\mathbb{R}^2} K^2\left(\frac{x-u}{h_1}, \frac{y-v}{h_2}\right) du dv$$

$$(u', v') = \left(\frac{x-u}{h_1}, \frac{y-v}{h_2}\right) \quad \langle du', dv' \rangle = \left(-\frac{du}{h_1}, -\frac{dv}{h_2}\right)$$

01
$$\text{Alors, } Var(\hat{f}_{h_1, h_2}(x, y)) \leq \frac{M}{n h_1 h_2} \int_{\mathbb{R}^2} K^2(u', v') du' dv'$$

$$= \frac{M \|K\|_2^2}{n h_1 h_2} \quad \text{c.q.f.d}$$

$$4^\circ) \quad | \text{biais}(\hat{f}_h(x,y)) | = | E(f_h(x,y)) - f(x,y) | \quad (0.7)$$

$$= | (K_h * f)(x,y) - f(x,y) |.$$

comme K est \rightarrow $= \frac{1}{h_1 h_2} \iint_{\mathbb{R}^2} K\left(\frac{x-u}{h_1}, \frac{y-v}{h_2}\right) f(u,v) du dv$
 le moyen \rightarrow $= \frac{1}{h_1 h_2} \iint_{\mathbb{R}^2} K\left(\frac{x-u}{h_1}, \frac{y-v}{h_2}\right) f(x,y) dudv$

$$\stackrel{(0.8)}{\leq} \frac{1}{h_1 h_2} \iint_{\mathbb{R}^2} |K\left(\frac{x-u}{h_1}, \frac{y-v}{h_2}\right)| |f(u,v) - f(x,y)| dudv$$

$$\leq \frac{1}{h_1 h_2} \iint_{\mathbb{R}^2} |K\left(\frac{x-u}{h_1}, \frac{y-v}{h_2}\right)| (|u-x|^{1/2} + |v-y|) dudv$$

Après un changement de variable: $(u,v) = \left(\frac{x-u}{h_1}, \frac{y-v}{h_2}\right)$.

$$(0.9) \quad | \text{biais}(\hat{f}_h(x,y)) | \leq h_1^{1/2} \iint_{\mathbb{R}^2} |K(u,v)| |u|^{1/2} du dv + h_2 \iint_{\mathbb{R}^2} |K(u,v)| |v| du dv$$

$$\stackrel{(0.9)}{=} C_1 h_1^{1/2} + C_2 h_2$$

$$\leq \max(C_1, C_2) (h_1^{1/2} + h_2)$$

$C_K = (-h_1^{1/2} + h_2) \quad \text{c.g. } f = \delta$

5) En utilisant successivement la décomposition "biais au carré - variance", les bornes obtenues dans les questions 4 et 5 et $(a+b)^2 \leq 2(a^2+b^2)$ par tout $a, b \geq 0$, nous obtenons:

$$(0.10) \quad \text{MSE}_{\hat{f}_h}(x,y) := \text{biais}^2 + \text{variance}$$

$$\leq \max(2C_K^2, M \|K\|_2^2) (h_1 + h_2^2 + \frac{1}{h_1 h_2})$$

6) f_x densité de X .

$$f_x(x) = \int_{\mathbb{R}} f(x,y) dy \quad \text{d'où}$$

$$\hat{f}_x(x) = \frac{1}{n h_1 h_2} \sum_{i=1}^n \int_{\mathbb{R}} K\left(\frac{x-x_i}{h_1}, \frac{y-y_i}{h_2}\right) dy$$

Si $H(x,y) \in \mathbb{R}^2$, $K(x,y) = F(x) \cdot G(y)$ alors:

$$\hat{f}_x(x) = \frac{1}{n h_1 h_2} \sum_{i=1}^n \int_{\mathbb{R}} F\left(\frac{x-x_i}{h_1}\right) G\left(\frac{y-y_i}{h_2}\right) dy$$

$$= \frac{1}{n h_1 h_2} \sum_{i=1}^n F\left(\frac{x-x_i}{h_1}\right) \int_{\mathbb{R}} G(u) \cdot h_2 du$$

$$\hat{f}_x(x) = \frac{1}{n h_1} \sum_{i=1}^n F\left(\frac{x-x_i}{h_1}\right)$$