
DS7 (vB) - Concours blanc

La présentation, la lisibilité, l'orthographe, la qualité de la rédaction, la clarté et la précision des raisonnements entreront pour une part importante dans l'appréciation des copies.

*Les candidat·es sont invité·es à **encadrer** dans la mesure du possible les résultats de leurs calculs.*

*Aucun document n'est autorisé. **L'utilisation de toute calculatrice et de tout matériel électronique est interdite.** Seule l'utilisation d'une règle graduée est autorisée.*

Si au cours de l'épreuve, un candidat ou une candidate repère ce qui lui semble être une erreur d'énoncé, il la signalera sur sa copie et poursuivra sa composition en expliquant les raisons des initiatives qu'il sera amené à prendre.

On suppose, et c'est valable pour toute l'épreuve, que la librairie suivante est importée sous son alias habituel :

- `import numpy as np`
- `import numpy.linalg as al`
- `import numpy.random as rd`
- `import matplotlib.pyplot as plt`
- `import pandas as pd`

On s'intéresse dans ce sujet à la méthode de Stein, introduite par Charles Stein (1920/2016) en 1972, dont les développements et applications sont nombreux.

Les parties 1 et 2 concernent la justification de la méthode, elles sont indépendantes.

Dans la partie 3, on s'intéresse à l'estimation en un point d'une densité d'une loi de probabilité. Cette partie peut être traitée indépendamment des deux premières parties.

Dans la partie 4, on met en œuvre la méthode de Stein, vue dans les parties 1 et 2, pour établir des convergences « uniformes » en loi et on démontre le résultat admis dans la partie 3. Cette partie est indépendante de la partie 3 à l'exception de sa dernière question.

Dans tout le problème :

- les variables aléatoires considérées sont définies sur le même espace probabilisé $(\Omega, \mathcal{A}, \mathbb{P})$.
- si X est une variable aléatoire, $\mathbb{E}(X)$ et $\mathbb{V}(X)$ désignent respectivement, lorsqu'elles existent, l'espérance et la variance de X .
- W désigne l'ensemble des fonctions h de classe \mathcal{C}^1 sur \mathbb{R} telles que :

$$\forall x \in \mathbb{R}, |h'(x)| \leq 1$$

- N est une variable aléatoire qui suit la loi normale $(0, 1)$.
- on admet que si X est une variable aléatoire possédant une espérance et $h \in W$, $\mathbb{E}(h(X))$ existe. On note en particulier c_h l'espérance de $h(N)$.

- On note Φ la fonction de répartition de la loi normale $(0, 1)$ définie par, pour tout $x \in \mathbb{R}$,

$$\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt. \text{ On rappelle que c'est la primitive sur } \mathbb{R}, \text{ qui vaut } \frac{1}{2} \text{ en } 0, \text{ de la fonction}$$

$$\varphi : t \mapsto \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}}.$$

Partie 1 - Transformation de Stein

Soit $h \in W$. On définit sur \mathbb{R} , la fonction $\theta : x \mapsto \frac{\Phi(x)}{\varphi(x)}$ et la fonction f_h par,

$$f_h : x \mapsto \theta(-x) \int_{-\infty}^x h'(t)\Phi(t) dt + \theta(x) \int_x^{+\infty} h'(t)(1 - \Phi(t)) dt$$

lorsque ces intégrales convergent.

L'objectif principal de cette partie est d'obtenir, pour X une variable aléatoire admettant une espérance, une expression de $\mathbb{E}(h(X)) - \mathbb{E}(h(N))$ qui ne fait pas intervenir N directement.

- 1. a)** Montrer que pour tout $x \geq 0$ et $t \in [x, +\infty[$, $0 \leq x\varphi(t) \leq t\varphi(t)$. En déduire que :

$$\forall x \geq 0, 0 \leq x(1 - \Phi(x)) \leq \varphi(x)$$

(on remarquera que pour tout $t \in \mathbb{R}$, $\varphi'(t) = -t\varphi(t)$)

- b)** Procéder de façon analogue pour montrer que : $\forall x \leq 0, -\varphi(x) \leq x\Phi(x) \leq 0$.

- c)** En déduire à l'aide d'une intégration par parties, pour tout x réel, la convergence des intégrales qui suivent et montrer que :

$$\int_{-\infty}^x \Phi(t) dt = x\Phi(x) + \varphi(x) \quad \text{et} \quad \int_x^{+\infty} (1 - \Phi(t)) dt = -x(1 - \Phi(x)) + \varphi(x) \quad (R_1)$$

2. a) Montrer que pour tous réels x et y ,

$$|h(x) - h(y)| \leq |x - y|, \quad \text{puis que } |h(x)| \leq |x| + |h(0)|$$

b) Pour tout x réel, justifier la convergence de $\int_{-\infty}^x h'(t)\Phi(t) dt$ et montrer que :

$$\int_{-\infty}^x h'(t)\Phi(t) dt = h(x)\Phi(x) - \int_{-\infty}^x h(t)\varphi(t) dt$$

On admet de même que, $\int_x^{+\infty} h'(t)(1 - \Phi(t)) dt$ converge et que,

$$\int_x^{+\infty} h'(t)(1 - \Phi(t)) dt = -h(x)(1 - \Phi(x)) + \int_x^{+\infty} h(t)\varphi(t) dt$$

c) En déduire que, pour tout x réel :

$$-\int_{-\infty}^x h'(t)\Phi(t) dt + \int_x^{+\infty} h'(t)(1 - \Phi(t)) dt = c_h - h(x)$$

3. a) Établir que pour tout $x \in \mathbb{R}$,

$$\begin{aligned} \theta'(x) &= 1 + x\theta(x) \\ \theta''(x) &= x + (1 + x^2)\theta(x) \\ \theta(-x)\Phi(x) &= \theta(x)(1 - \Phi(x)) \end{aligned}$$

b) En déduire que f_h est une fonction de classe \mathcal{C}^1 sur \mathbb{R} qui vérifie, pour tout x réel :

$$f_h'(x) - xf_h(x) = c_h - h(x)$$

Pourquoi peut-on alors affirmer que f_h est de classe \mathcal{C}^2 sur \mathbb{R} ?

c) En conclure que, si X est une variable aléatoire admettant une espérance,

$$|\mathbb{E}(h(X)) - \mathbb{E}(h(N))| = |\mathbb{E}(f_h'(X) - Xf_h(X))|$$

4. Majoration de $|f_h|$.

a) Montrer, en utilisant les égalités (R_1), que pour tout x réel :

$$\theta(-x) \int_{-\infty}^x \Phi(t) dt + \theta(x) \int_x^{+\infty} (1 - \Phi(t)) dt = 1$$

b) En déduire que pour tout x réel : $|f_h(x)| \leq 1$.

5. Majoration de $|f_h''|$.

a) Montrer que pour tout x réel :

$$\theta''(-x) \int_{-\infty}^x \Phi(t) dt + \theta''(x) \int_x^{+\infty} (1 - \Phi(t)) dt = 1$$

b) Établir pour tout x réel l'égalité :

$$f_h''(x) = -h'(x) + \theta''(-x) \int_{-\infty}^x h'(t)\Phi(t) dt + \theta''(x) \int_x^{+\infty} h'(t)(1 - \Phi(t)) dt$$

c) Étudier les variations sur \mathbb{R} de la fonction $x \mapsto \Phi(x) + \frac{x}{1+x^2}\varphi(x)$. En déduire son signe et le signe de θ'' .

En conclure que, pour tout x réel : $|f_h''(x)| \leq 2$.

Partie 2 - Majoration uniforme de la distance de Kolmogorov

Dans la suite du problème, si X est une variable aléatoire de fonction de répartition F_X , on définit, pour tout x réel, $d_X(x)$ la distance de Kolmogorov au point x entre la loi de X et la loi normale centrée réduite par :

$$d_X(x) = |F_X(x) - \Phi(x)|$$

On définit, pour tout x réel, la fonction h_x sur \mathbb{R} par $h_x(t) = \begin{cases} 1 & \text{si } t \leq x \\ 0 & \text{si } t > x \end{cases}$.

On définit aussi la fonction γ sur \mathbb{R} par :

$$\gamma(t) = \begin{cases} 1 & \text{si } t < 0 \\ 0 & \text{si } t > 1 \\ 1 - 3t^2 + 2t^3 & \text{si } t \in [0, 1] \end{cases}$$

Soit X une variable aléatoire.

6. Pour tout x réel, quelle est la loi de la variable aléatoire $h_x(X)$? En déduire que $\mathbb{E}(h_x(X))$ existe et vaut $F_X(x)$.

7. a) Écrire une fonction **Python** `gamma(t)` qui calcule et renvoie la valeur de $\gamma(t)$, t étant donné.

b) Utiliser la fonction précédente pour écrire un script qui affiche le graphe de γ sur le segment $[-1, 2]$ dans un repère.

8. a) Montrer que γ est continue sur \mathbb{R} , de classe \mathcal{C}^1 sur \mathbb{R} privé de 0 et 1.

b) Étudier les variations de γ sur $[0, 1]$ et montrer que pour tout $t \in \mathbb{R}$, $\gamma(t) \in [0, 1]$.

c) Établir que γ est dérivable en 1 et que $\gamma'(1) = 0$.

On montrerait de même que γ est dérivable en 0 et que $\gamma'(0) = 0$. On l'admet.

d) Justifier que γ est de classe \mathcal{C}^1 sur \mathbb{R} et que pour tout t réel $|\gamma'(t)| \leq \frac{3}{2}$.

On suppose dans la suite de cette partie que X admet une espérance et on considère un réel M_X qui vérifie, pour tout $h \in W$, $|\mathbb{E}(h(X)) - \mathbb{E}(h(N))| \leq M_X$.

9. Soit $t > 0$ et x un réel. Pour tout $y \in \mathbb{R}$, on pose $k_x(y) = \gamma\left(\frac{y-x}{t}\right)$.

a) Montrer que pour tout y réel, $h_x(y) \leq k_x(y)$.

b) On admet l'existence de $\mathbb{E}(k_x(X))$ et de $\mathbb{E}(k_x(N))$. Justifier l'inégalité suivante :

$$\mathbb{E}(h_x(X)) - \mathbb{E}(h_x(N)) \leq \mathbb{E}(k_x(X)) - \mathbb{E}(k_x(N)) + \mathbb{E}(k_x(N)) - \mathbb{E}(h_x(N))$$

c) Montrer que $\mathbb{E}(k_x(N)) - \mathbb{E}(h_x(N)) = \int_x^{x+t} k_x(u)\varphi(u)du$.

d) Établir que la fonction g , définie sur \mathbb{R} par $g : u \mapsto \frac{2t}{3}k_x(u)$, appartient à W . En déduire que :

$$\mathbb{E}(h_x(X)) - \mathbb{E}(h_x(N)) \leq \frac{3}{2t}M_X + \frac{t}{\sqrt{2\pi}} \leq \frac{3}{2t}M_X + \frac{t}{2}$$

On admet de même, qu'en utilisant la fonction k_{x-t} , on a :

$$\mathbb{E}(h_x(N)) - \mathbb{E}(h_x(X)) \leq \frac{3}{2t}M_X + \frac{t}{2}$$

10. En étudiant la fonction $t \mapsto \frac{3}{2t}M_X + \frac{t}{2}$ sur $]0, +\infty[$, en déduire que, pour tout x réel,

$$|\mathbb{E}(h_x(X)) - \mathbb{E}(h_x(N))| \leq \sqrt{3M_X}, \text{ puis que } d_X(x) \leq \sqrt{3M_X} \quad (R_2)$$

Partie 3 - Estimation d'une densité

On considère X une variable aléatoire à densité de fonction de répartition F et de densité de probabilité f qui dépendent d'un paramètre inconnu θ , où $\theta \in \Theta$, Θ un intervalle de \mathbb{R} .

Soit a un point de continuité de f , fixé. On souhaite estimer $f(a)$.

Par exemple, si X suit la loi exponentielle de paramètre θ et $a > 0$, on souhaite estimer $\theta e^{-\theta a}$.

On dispose pour tout $\theta \in \Theta$, d'une suite de variables aléatoires $(X_i)_{i \geq 1}$ indépendantes de même loi que X .

On choisit une suite $(h_n)_{n \geq 1}$ de réels strictement positifs tels que :

$$\lim_{n \rightarrow +\infty} h_n = 0 \text{ et } \lim_{n \rightarrow +\infty} nh_n = +\infty$$

Pour tout $n \in \mathbb{N}^*$, et $\omega \in \Omega$, on définit :

$C_n(\omega)$ comme le nombre d'indices $i \in \llbracket 1, n \rrbracket$ tels que $X_i(\omega) \in]a - h_n, a + h_n]$

et $f_n(\omega) = \frac{1}{2nh_n} C_n(\omega)$.

11. On suppose que l'on dispose d'un fichier `stats.csv` qui comporte une colonne nommée `salaire`. On considère que les valeurs de cette colonne constituent la réalisation d'un échantillon de la loi de X dont la taille dépasse 10000.

a) Après avoir exécuté `import pandas as pd`, quelle(s) instruction(s) permet(tent) de lire dans le fichier `stats.csv` les valeurs de la colonne `salaire` et d'affecter cette série `pandas` obtenue à une variable échantillon ?

On supposera que le fichier `stats.csv` se trouve dans le répertoire de travail.

b) On souhaite calculer et afficher $f_n(\omega)$ pour a donné, lorsque la réalisation d'un échantillon $(X_1(\omega), \dots, X_n(\omega))$ de la loi de X est représentée en **Python** par `échantillon` et, pour tout $n \in \mathbb{N}^*$, $h_n = \frac{1}{\sqrt{n}}$.

Compléter le script suivant pour qu'il réalise cette tâche.

```

1  a = float(input('a='))
2  n = échantillon.count()
3  h = 1 / np.sqrt(n)
4  C = 0
5  for i in range(n):
6      if ... and ...:
7          ... += 1
8  print(C / ...)
```

12. Montrer que C_n suit une loi binomiale de paramètres (n, p_n) en précisant l'expression de p_n en fonction de a et h_n .

En déduire que $\mathbb{E}(f_n)$ existe et vaut : $\frac{F(a + h_n) - F(a - h_n)}{2h_n}$.

13. a) En utilisant la dérivabilité de F en a , montrer que $\lim_{n \rightarrow +\infty} \mathbb{E}(f_n) = f(a)$.

b) Montrer que $\mathbb{V}(f_n)$ existe et que $\lim_{n \rightarrow +\infty} \mathbb{V}(f_n) = 0$.

On suppose désormais, que $f(a) > 0$, que pour tout $n \in \mathbb{N}^*$, $p_n \in]0, 1[$, que F est de classe \mathcal{C}^2 au voisinage de a et que $\lim_{n \rightarrow +\infty} nh_n^3 = 0$.

On note pour tout $n \geq 1$, $\sigma_n = \sqrt{np_n(1-p_n)}$ et $\theta_n = \sqrt{2h_n f(a)}$.

On définit les variables aléatoires : $D_n = \frac{C_n - np_n}{\sigma_n}$ et $\hat{f}_n = \frac{\theta_n \sqrt{n}}{f(a)}(f_n - f(a))$.

14. a) En utilisant le développement limité de F à l'ordre 2 au point a , montrer que :

$$p_n \underset{n \rightarrow +\infty}{=} 2h_n f(a) + o(h_n^2) \underset{n \rightarrow +\infty}{=} \theta_n^2 + o(h_n^2)$$

b) En déduire que : $p_n \underset{n \rightarrow +\infty}{\sim} 2h_n f(a)$, puis que $\lim_{n \rightarrow +\infty} np_n = +\infty$.

c) Montrer que $\hat{f}_n = \frac{\sigma_n}{\theta_n \sqrt{n}} D_n + \sqrt{n} \left(\frac{p_n}{\theta_n} - \theta_n \right)$ et que l'on a :

$$\lim_{n \rightarrow +\infty} \frac{\sigma_n}{\theta_n \sqrt{n}} = 1 \quad ; \quad \lim_{n \rightarrow +\infty} \sqrt{n} \left(\frac{p_n}{\theta_n} - \theta_n \right) = 0$$

On admet, dans la suite de cette partie, que $(\hat{f}_n)_{n \geq 1}$ converge en loi vers N ce qui implique que pour tout $(x, y) \in \mathbb{R}^2$, avec $x \leq y$, on a :

$$\lim_{n \rightarrow +\infty} \mathbb{P}(x \leq \hat{f}_n \leq y) = \Phi(y) - \Phi(x)$$

15. Soit $\alpha \in]0, 1[$. On pose $\eta_\alpha = t_\alpha^2$ où t_α est le quantile d'ordre $1 - \frac{\alpha}{2}$ de la loi normale $(0, 1)$.

a) Montrer que $\lim_{n \rightarrow +\infty} \mathbb{P} \left((f(a))^2 - \left(2f_n + \frac{\eta_\alpha}{2nh_n} \right) f(a) + f_n^2 \leq 0 \right) = 1 - \alpha$.

b) On note, pour $n \geq 1$, $\Delta_n = \sqrt{\left(f_n + \frac{\eta_\alpha}{4nh_n} \right)^2 - f_n^2}$.

Montrer que :

$$\lim_{n \rightarrow +\infty} \mathbb{P} \left(f(a) \in \left[f_n + \frac{\eta_\alpha}{4nh_n} - \Delta_n, f_n + \frac{\eta_\alpha}{4nh_n} + \Delta_n \right] \right) = 1 - \alpha$$

Partie 4 - Convergence « uniforme » en loi vers la loi normale

Soit $n \in \mathbb{N}^*$. On considère n variables aléatoires X_1, \dots, X_n indépendantes centrées qui possèdent un moment d'ordre 3. On admet alors que ces variables aléatoires possèdent une variance.

On pose, pour tout $k \in \llbracket 1, n \rrbracket$, $\mathbb{E}(X_k^2) = v_k$, $S_n = \sum_{k=1}^n X_k$, $Y_k = S_n - X_k$ et on suppose que $\sum_{k=1}^n v_k = 1$.

Soit f une fonction de classe \mathcal{C}^2 sur \mathbb{R} telle que pour tout x réel, $|f(x)| \leq 1$ et $|f''(x)| \leq 2$.

On admet que si X et Y sont des variables aléatoires possédant une espérance alors $\mathbb{E}(Yf(X))$ et $\mathbb{E}(f'(X))$ existent.

16. a) Montrer que $\sum_{k=1}^n v_k \mathbb{E}(f'(S_n) - f'(Y_k)) + \sum_{k=1}^n \mathbb{E}(X_k^2 f'(Y_k)) = \mathbb{E}(f'(S_n))$.

b) Montrer que $\sum_{k=1}^n \mathbb{E}(X_k [f(S_n) - f(Y_k)]) = \sum_{k=1}^n \mathbb{E}(X_k f(S_n)) = \mathbb{E}(S_n f(S_n))$.

c) En déduire que :

$$\mathbb{E}(f'(S_n) - S_n f(S_n)) = \sum_{k=1}^n v_k \mathbb{E}(f'(S_n) - f'(Y_k)) + \sum_{k=1}^n \mathbb{E}(X_k [X_k f'(Y_k) - (f(S_n) - f(Y_k))])$$

17. Soit a et b deux réels.

a) Montrer que :

$$bf'(a) - (f(a+b) - f(a)) = \int_0^1 b(f'(a) - f'(a+tb)) dt$$

b) En déduire que :

$$|bf'(a) - (f(a+b) - f(a))| \leq b^2$$

c) En conclure que :

$$|\mathbb{E}(f'(S_n) - S_n f(S_n))| \leq 2 \sum_{k=1}^n v_k \mathbb{E}(|X_k|) + \sum_{k=1}^n \mathbb{E}(|X_k|^3)$$

puis, grâce à l'inégalité (R_2), que, pour tout x réel :

$$d_{S_n}(x) \leq \sqrt{3 \left(2 \sum_{k=1}^n v_k \mathbb{E}(|X_k|) + \sum_{k=1}^n \mathbb{E}(|X_k|^3) \right)} \quad (R_3)$$

Une définition - Dans la suite du sujet, si $(X_n)_{n \geq 1}$ est une suite de variables aléatoires réelles et $(\delta_n)_{n \geq 1}$ une suite réelle de limite nulle qui vérifient,

$$\forall n \in \mathbb{N}^*, \forall x \in \mathbb{R}, d_{X_n}(x) \leq \delta_n$$

on dira alors que $(X_n)_{n \geq 1}$ converge uniformément en loi vers N .

On remarque, et on l'admet pour la suite, que si $(X_n)_{n \geq 1}$ converge uniformément en loi vers N alors $(X_n)_{n \geq 1}$ converge en loi vers N .

18. *Une première application.* On suppose dans cette question que $(Z_k)_{k \geq 1}$ est une suite de variables aléatoires indépendantes, suivant la même loi et admettant des moments d'ordre 1 à 3.

On note pour $i \in \llbracket 1, 3 \rrbracket$, $s_i = \mathbb{E}(|Z_k - \mathbb{E}(Z_k)|^i)$, $\sigma = \sqrt{s_2}$ et $X_k = \frac{Z_k - \mathbb{E}(Z_k)}{\sigma \sqrt{n}}$ pour tout $k \in \mathbb{N}^*$.

On suppose que $\sigma \neq 0$.

On utilise les notations de la question précédente.

a) Montrer que l'on peut appliquer l'inégalité (R_3) qui donne ici :

$$d_{S_n}(x) \leq \sqrt{3 \frac{2\sigma^2 s_1 + s_3}{\sigma^3 \sqrt{n}}}$$

b) En déduire que $(S_n)_{n \geq 1}$ converge uniformément en loi vers N , donc converge en loi vers N . Quel résultat du cours nous aurait permis d'obtenir cette dernière convergence directement ?

19. *Une deuxième application.* On suppose dans cette question que Z_1, \dots, Z_n sont des variables aléatoires indépendantes suivant la même loi de Bernoulli de paramètre $p_n \in]0, 1[$.

On pose $\sigma_n = \sqrt{np_n(1-p_n)}$ et $X_k = \frac{Z_k - p_n}{\sigma_n}$ pour tout $k \in \llbracket 1, n \rrbracket$.

a) Montrer que $\mathbb{E}(|X_k|) = \frac{2\sigma_n}{n}$ et $\mathbb{E}(|X_k|^3) \leq \frac{2}{n\sigma_n}$.

b) En déduire que, pour tout x réel :

$$d_{S_n}(x) \leq 2 \sqrt{3 \left(\frac{\sigma_n}{n} + \frac{1}{2\sigma_n} \right)}$$

c) Justifier le résultat suivant :

si pour tout $n \in \mathbb{N}^*$, T_n est une variable aléatoire qui suit la loi binomiale (n, p_n) avec $\lim_{n \rightarrow +\infty} p_n = 0$ et $\lim_{n \rightarrow +\infty} np_n = +\infty$ alors, $\left(\frac{T_n - np_n}{\sqrt{np_n(1-p_n)}} \right)_{n \geq 1}$ converge uniformément en loi vers N .

20. *Un petit lemme.* Soit $(V_n)_{n \geq 1}$ une suite de variables aléatoires réelles qui converge uniformément en loi vers N . Soit $(a_n)_{n \geq 1}$ une suite de réels strictement positifs qui converge vers a , tel que $a > 0$, et $(b_n)_{n \geq 1}$ une suite de réels qui converge vers b .

a) Soit X une variable aléatoire et (α, β) un couple de réels avec $\alpha > 0$. On note $F_{\alpha X + \beta}$ et $F_{\alpha N + \beta}$ les fonctions de répartition respectives de $\alpha X + \beta$ et $\alpha N + \beta$.

Montrer que, pour tout x réel,

$$|F_{\alpha X + \beta}(x) - F_{\alpha N + \beta}(x)| = d_X \left(\frac{x - \beta}{\alpha} \right)$$

b) Montrer que pour tout x réel,

$$\lim_{n \rightarrow +\infty} (\mathbb{P}(a_n V_n + b_n \leq x) - \mathbb{P}(a_n N + b_n \leq x)) = 0$$

c) Établir que $\lim_{n \rightarrow +\infty} \mathbb{P}(a_n N + b_n \leq x) = \mathbb{P}(aN + b \leq x)$ puis en déduire que $(a_n V_n + b_n)_{n \geq 1}$ converge en loi vers $aN + b$. Quelle est la loi de la variable aléatoire $aN + b$?

21. On reprend les notations de la partie 3.

a) Justifier que $(D_n)_{n \geq 1}$ converge uniformément en loi vers N .

b) En utilisant les résultats des questions 14. et 20., en déduire que la suite $(\hat{f}_n)_{n \geq 1}$ converge en loi vers N .