

Troisième Année Licence M.I.A.S.H.S. 2024 – 2025

Statistique 2

Examen terminal, Mai 2025

Examen de 2h00. Tout document ou calculatrice est interdit.

Exercice 1 (Sur 10 points)

Soit X une variable aléatoire suivant une distribution uniforme sur $[0, 1]$. Pour $m \in [0, 1/2]$, on définit $Y = |X - m|$.

1. Prouver que pour tout $m \in [0, 1/2]$, $\mathbb{E}[Y] = m^2 - m + 1/2$ (**1pt**).
2. Déterminer la fonction de répartition de Y (**1 pt**) et en déduire que Y est une variable aléatoire absolument continue de densité $f_Y(y) = 2$ pour $y \in [0, m]$, $f_Y(y) = 1$ pour $y \in]m, 1 - m]$ et $f_Y(y) = 0$ ailleurs (**0.5pts**).
3. Supposons que m est inconnu et que (Y_1, \dots, Y_n) est une famille de variables aléatoires i.i.d. observées suivant la même distribution que Y . Avec \bar{Y}_n la moyenne empirique de (Y_1, \dots, Y_n) , montrer que $\bar{m}_n = \frac{1}{2}(1 - \sqrt{|4\bar{Y}_n - 1|})$ est un estimateur convergent de m (**1.5pts**). Montrer que $\text{var}(Y) = \frac{1}{12} - m^2(1 - m)^2$ (**1pt**). En déduire que pour $m \in [0, 1/2]$, \bar{m}_n vérifie un théorème central limite que l'on précisera (**2pts**).
4. On définit un second estimateur $\hat{m}_n = 1 - \max(Y_1, \dots, Y_n)$ de m . Montrer que $\mathbb{P}(m \leq \hat{m}_n \leq m + \varepsilon) = 1 - (1 - \varepsilon)^n$ pour $0 \leq \varepsilon \leq 1 - 2m$ (**1pt**). Déduisez un intervalle de confiance de niveau 95% pour m (**0.5pts**).
5. Démontrer que $n(\hat{m}_n - m) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{E}(1)$, loi exponentielle de paramètre 1 (**1.5pts**).

Proof. 1. Pour tout $m \in [0, 1/2]$, on a

$$\begin{aligned} \mathbb{E}[|X - m|] &= \int_0^m (m - x) dx + \int_m^1 (x - m) dx \\ &= \left[mx - \frac{1}{2}x^2 \right]_0^m + \left[-mx + \frac{1}{2}x^2 \right]_m^1 \\ &= m^2 - \frac{1}{2}m^2 - m + \frac{1}{2} + m^2 - \frac{1}{2}m^2 \\ &= m^2 - m + \frac{1}{2}. \end{aligned}$$

2. Puisque $m \in [0, 1/2]$, $|X - m| \in [0, 1 - m]$. Donc, pour $y < 0$, $F_Y(y) = 0$ et $y \geq 1 - m$, $F_Y(y) = 1$. Enfin, pour $0 \leq y \leq 1 - m$,

$$\begin{aligned} F_Y(y) &= \mathbb{P}(-y \leq X - m \leq y) \\ &= \mathbb{P}(m - y \leq X \leq y + m). \end{aligned}$$

Par conséquent, on a $F_Y(y) = \mathbb{P}(m - y \leq X \leq y + m) = 2y$ if $0 \leq y \leq m$ et $F_Y(y) = \mathbb{P}(0 \leq X \leq y + m) = y + m$ if $m \leq y \leq 1 - m$.

Il est clair que la fonction $y \in \mathbf{R} \mapsto F_Y(y)$ est une fonction de classe \mathcal{C}^1 par morceaux. Donc, Y est une variable absolument continue et sa densité de probabilité est donnée par sa dérivée (sauf aux points non différentiables: 0, m et $1 - m$). Par conséquent, on peut choisir que $f_Y(y) = 2$ pour $y \in [0, m]$, $f_Y(y) = 1$ pour $y \in (m, 1 - m]$ et $f_Y(y) = 0$ ailleurs.

3. Il est clair que (Y_1, \dots, Y_n) est une suite de v.a.i.i.d. telle que $\mathbb{E}[|Y_1|] < \infty$ donc la loi forte des grands nombres peut être appliquée et nous obtenons $\bar{Y}_n \xrightarrow[n \rightarrow +\infty]{\mathcal{P}} \mathbb{E}[Y_1]$. Considérons alors la fonction $g : y \in \mathbf{R} \mapsto \frac{1}{2}(1 - \sqrt{|4y - 1|})$ qui est continue. Alors $g(\bar{Y}_n) \xrightarrow[n \rightarrow +\infty]{\mathcal{P}} g(\mathbb{E}[Y_1])$. Comme $\mathbb{E}[Y_1] \geq 1/4$, $|4\mathbb{E}[Y_1] - 1| = 4(m^2 - m + \frac{1}{2}) - 1 = (1 - 2m)^2$ et donc $g(\mathbb{E}[Y_1]) = m$, ce qui implique que $g(\bar{Y}_n) = \bar{m}_n \xrightarrow[n \rightarrow +\infty]{\mathcal{P}} m$: l'estimateur est convergent.

Il suffit de calculer $\text{var}(Y_1) = \mathbb{E}[Y_1^2] - \mathbb{E}[Y_1]^2 = \int_0^1 (x - m)^2 dx - (m^2 - m + 1/2)^2 = 1/3 - m + m^2 - (m^2 - m + 1/2)^2 = \frac{1}{12} - m^2(1 - m)^2$. On utilise ensuite la méthode Delta. Pour ce faire, on vérifie que la fonction g est classe \mathcal{C}^1 sur $[0, 1/4] \cup [1/4, 1/2]$ et $(g'(x))^2 = 1/|1 - 4x|$ ne s'annule pas sur $[0, 1/4] \cup [1/4, 1/2]$. Or $\mathbb{E}[Y] = (m - 1/2)^2 + 1/4 > 1/4$ pour $m \in [0, 1/2]$. D'où pour $m \in [0, 1/2]$,

$$\sqrt{n}(\bar{m}_n - m) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}\left(0, \frac{1}{(1 - 2m)^2} \left(\frac{1}{12} - m^2(1 - m)^2\right)\right).$$

4. Puisque les Y_i prennent leurs valeurs dans $[0, 1 - m]$, alors $\max(Y_1, \dots, Y_n)$ prend aussi ses valeurs dans $[0, 1 - m]$, ce qui implique que \widehat{m}_n prend ses valeurs dans $[m, 1]$.
Soit $0 \leq \varepsilon \leq 1 - 2m$. Alors

$$\begin{aligned}\mathbb{P}(m \leq \widehat{m}_n \leq m + \varepsilon) &= \mathbb{P}(\widehat{m}_n \leq m + \varepsilon) \text{ puisque } \widehat{m}_n \text{ prend ses valeurs dans } [m, 1] \\ &= 1 - \mathbb{P}(1 - \max(Y_1, \dots, Y_n) > m + \varepsilon) \\ &= 1 - \mathbb{P}(\max(Y_1, \dots, Y_n) < 1 - m - \varepsilon) \\ &= 1 - \mathbb{P}(Y < 1 - m - \varepsilon)^n \text{ puisque } (Y_i) \text{ v.a.i.i.d.}\end{aligned}$$

Avec $0 \leq \varepsilon \leq 1 - 2m$ on a $m \leq 1 - m - \varepsilon \leq 1 - m$ et donc $\mathbb{P}(Y < 1 - m - \varepsilon) = F_Y(1 - m - \varepsilon) = 1 - m - \varepsilon + m = 1 - \varepsilon$. D'où le résultat.

De $\mathbb{P}(m \leq \widehat{m}_n \leq m + \varepsilon) = \mathbb{P}(\widehat{m}_n - \varepsilon \leq m \leq \widehat{m}_n) = 1 - (1 - \varepsilon)^n$, on déduit que si $1 - (1 - \varepsilon)^n = 0.95$ alors $[\widehat{m}_n - \varepsilon, \widehat{m}_n]$ est un intervalle de confiance de niveau 95% pour m . Cela implique que $(1 - \varepsilon)^n = 0,05$ et donc $\varepsilon = 1 - 0,05^{1/n}$. Par conséquent, $\varepsilon = 1 - 0,05^{1/n}$. $[\widehat{m}_n - 1 + 0,05^{1/n}, \widehat{m}_n]$ est un intervalle de confiance à 95%.

5. D'après la question 4, avec $x = n\varepsilon \geq 0$, soit $\varepsilon = x/n$,

$$\mathbb{P}(n(\widehat{m}_n - m) \leq n\varepsilon) = \mathbb{P}(n(\widehat{m}_n - m) \leq x) = 1 - (1 - x/n)^n = 1 - e^{n \log(1 - x/n)} \xrightarrow[n \rightarrow \infty]{} 1 - e^{-x},$$

puisque $\log(1 - x/n) \sim -x/n$ lorsque $n \rightarrow \infty$. Mais pour $x \geq 0$, $1 - e^{-x}$ est la fonction de distribution cumulative d'une distribution $\mathcal{E}(1)$. Par conséquent $n(\widehat{m}_n - m) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{E}(1)$. □

Exercice 2 (Sur 16 points)

Soit $(X_i)_{i \in \mathbb{N}}$ une suite de variables aléatoires identiquement distribuées, d'espérance m et de variance σ^2 , telle que pour tout $n \in \mathbb{N}^*$, le vecteur aléatoire (X_1, \dots, X_n) a une matrice de covariance $\Gamma_n = (\gamma_{ij}^{(n)})_{1 \leq i, j \leq n}$ définie positive.

1. Que vaut $\gamma_{11}^{(n)}$ (**0.5pts**)? Montrer que $|\gamma_{ij}^{(n)}| \leq \sigma^2$ pour tout $1 \leq i, j \leq n$ (**0.5pts**).
 2. On note $\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k$. Déterminer $\mathbb{E}[\bar{X}_n]$ (**0.5pts**) et montrer que $\text{var}(\bar{X}_n) = \frac{\sigma^2}{n} + \frac{2}{n^2} \sum_{1 \leq i < j \leq n} \gamma_{ij}^{(n)}$ (**0.5pts**).
 3. On suppose désormais qu'il existe une fonction $r : \mathbb{N} \rightarrow \mathbb{R}$ telle que $\gamma_{ij}^{(n)} = r(|j - i|)$ pour tout $1 \leq i, j \leq n$ et tout $n \in \mathbb{N}^*$ et on suppose que $\sum_{k=0}^{\infty} |r(k)| < \infty$. Montrer que $\text{var}(\bar{X}_n) = \frac{\sigma^2}{n} + \frac{2}{n} \sum_{k=1}^n \left(1 - \frac{k}{n}\right) r(k)$ (**1.5pts**). En écrivant que $\sum_{k=1}^n = \sum_{k=1}^{\lceil \sqrt{n} \rceil} + \sum_{k=\lceil \sqrt{n} \rceil + 1}^n$, montrer que $n \text{var}(\bar{X}_n) \xrightarrow[n \rightarrow \infty]{} \sigma^2 + 2 \sum_{k=1}^{\infty} r(k)$ (**2pts**). En déduire que $\bar{X}_n \xrightarrow[n \rightarrow +\infty]{\mathcal{P}} m$ (**0.5pts**).
 4. Soit $(Z_n)_{n \in \mathbb{N}}$ une suite de variables aléatoires gaussiennes. Montrer que si les suites $(\mathbb{E}[Z_n])_{n \in \mathbb{N}}$ et $(\text{var}(Z_n))_{n \in \mathbb{N}}$ convergent, alors $(Z_n)_{n \in \mathbb{N}}$ converge en loi vers une loi que l'on précisera (**1pt**).
 5. On suppose désormais que la suite $(X_i)_{i \in \mathbb{N}}$ est telle que le vecteur (X_1, \dots, X_n) est gaussien pour tout $n \in \mathbb{N}^*$. Montrer que $\sqrt{n}(\bar{X}_n - m)$ converge en loi vers une limite que l'on précisera (**1pt**).
 6. On suppose que (X_1, \dots, X_n) est observé avec Γ_n connue, mais m inconnue et on veut estimer m . Montrer que maximiser la vraisemblance du modèle en $x = {}^t(x_1, \dots, x_n) \in \mathbb{R}^n$ revient à minimiser ${}^t(x - m \mathbb{I}_n) \Gamma_n^{-1} (x - m \mathbb{I}_n)$, où $\mathbb{I}_n = {}^t(1, 1, \dots, 1)$ (**1pt**). En déduire que l'estimateur du maximum de vraisemblance \widehat{m}_n de m est unique et:
- $$\widehat{m}_n = ({}^t \mathbb{I}_n \Gamma_n^{-1} \mathbb{I}_n)^{-1} {}^t \mathbb{I}_n \Gamma_n^{-1} {}^t (X_1, \dots, X_n) \quad (\text{2pts}).$$
7. Montrer que \widehat{m}_n est un estimateur sans biais (**0.5pts**) et déterminer sous forme matricielle la variance de \widehat{m}_n (**1pt**).
 8. Ecrire également $\text{var}(\bar{X}_n)$ sous forme matricielle en utilisant Γ_n et \mathbb{I}_n (**1pt**). Montrer que pour tout vecteur $U \in \mathbb{R}^n$, $({}^t U U)^2 \leq ({}^t U \Gamma_n U) ({}^t U \Gamma_n^{-1} U)$ (**1.5pts**). En déduire que $\text{var}(\widehat{m}_n) \leq \text{var}(\bar{X}_n)$ (**1pt**).

Proof. 1. On a $\gamma_{11}^{(n)} = \text{var}(X_1) = \sigma^2$.

D'après l'Inégalité de Cauchy-Schwarz, $|\gamma_{ij}^{(n)}| = |\text{cov}(X_i, X_j)| \leq \sqrt{\text{var}(X_i) \text{var}(X_j)} = \sigma^2$.

2. On a $\mathbb{E}[\bar{X}_n] = m$.

Et d'une manière générale $\text{var}(\bar{X}_n) = \frac{1}{n^2} \left(\sum_{k=1}^n \gamma_{kk}^{(n)} + 2 \sum_{1 \leq i < j \leq n} \gamma_{ij}^{(n)} \right)$ d'où le résultat.

3. On a $\text{var}(\bar{X}_n) = \frac{\sigma^2}{n} + \frac{2}{n^2} \sum_{1 \leq i < j \leq n} r(j-i)$. Mais pour $k \in \{1, \dots, n-1\}$, il y a $n-k$ couples (i, j) avec $i < j$ tels que $j-i=k$. Donc $\text{var}(\bar{X}_n) = \frac{\sigma^2}{n} + \frac{2}{n^2} \sum_{1 \leq k \leq n-1} (n-k)r(k)$, d'où le résultat.

Pour $k \in \{1, \dots, [\sqrt{n}]\}$, alors $\frac{k}{n} \leq \frac{\sqrt{n}}{n}$, donc $\left| \sum_{1 \leq k \leq [\sqrt{n}]} \frac{k}{n} r(k) \right| \leq \frac{1}{\sqrt{n}} \sum_{1 \leq k \leq [\sqrt{n}]} |r(k)| \leq \frac{1}{\sqrt{n}} \sum_{1 \leq k \leq \infty} |r(k)| \xrightarrow[n \rightarrow \infty]{} 0$.

Pour $k \in \{[\sqrt{n}]+1, \dots, n\}$, alors $\frac{k}{n} \leq 1$, donc $\left| \sum_{k=[\sqrt{n}]+1}^n \frac{k}{n} r(k) \right| \leq \sum_{k=[\sqrt{n}]+1}^n |r(k)| \xrightarrow[n \rightarrow \infty]{} 0$ comme reste d'une série convergente.

De ces deux points, on en déduit que $\sum_{k=1}^n n \frac{k}{n} r(k) \xrightarrow[n \rightarrow \infty]{} 0$, donc $n \text{var}(\bar{X}_n) \xrightarrow[n \rightarrow \infty]{} \sigma^2 + 2 \sum_{k=1}^{\infty} r(k)$.

De ce qui précède $\text{var}(\bar{X}_n) \xrightarrow[n \rightarrow \infty]{} 0$, donc d'après l'Inégalité de Bienaymé-Tchebytchev, $\bar{X}_n \xrightarrow[n \rightarrow +\infty]{} m$.

4. Il suffit de considérer la fonction caractéristique de Z_n qui s'écrit $\phi_{Z_n}(u) = e^{-\frac{1}{2} \text{var}(Z_n) u^2 + i u \mathbb{E}[Z_n]}$. Si $\mathbb{E}[Z_n] \xrightarrow[n \rightarrow \infty]{} m_Z$ et $\text{var}(Z_n) \xrightarrow[n \rightarrow \infty]{} \sigma_Z^2$ alors il est clair que pour tout $u \in \mathbf{R}$, $\phi_{Z_n}(u) \xrightarrow[n \rightarrow \infty]{} e^{-\frac{1}{2} \sigma_Z^2 u^2 + i u m_Z}$, donc $Z_n \xrightarrow[n \rightarrow \infty]{} \mathcal{N}(m_Z, \sigma_Z^2)$.

5. Si on considère $Z_n = n(\bar{X}_n - m)$, on sait que \bar{X}_n est une v.a. gaussienne comme combinaison linéaire issu du vecteur gaussien (X_1, \dots, X_n) , donc Z_n est une v.a. gaussienne, comme transformation affine d'une v.a. gaussienne. De plus $\mathbb{E}[Z_n] = 0$ pour tout $n \in \mathbf{N}$, et $\text{var}(Z_n) \xrightarrow[n \rightarrow \infty]{} \sigma^2 + 2 \sum_{k=1}^{\infty} r(k)$ d'après la question 3. Donc $n(\bar{X}_n - m) \xrightarrow[n \rightarrow \infty]{} \mathcal{N}(0, \sigma^2 + 2 \sum_{k=1}^{\infty} r(k))$.

6. On a donc $(X_1, \dots, X_n) \xrightarrow{\mathcal{L}} \mathcal{N}(m \mathbb{I}_n, \Gamma_n)$. Comme Γ_n est définie positive, on peut donc écrire que la densité de (X_1, \dots, X_n) en x vaut:

$$\begin{aligned} f_{(X_1, \dots, X_n)}(x) &= \frac{1}{(2\pi)^{n/2} \det(\Gamma_n)^{1/2}} \exp\left(-\frac{1}{2} {}^t(x - m \mathbb{I}_n) \Gamma_n^{-1} (x - m \mathbb{I}_n)\right) \\ \implies \log(f_{(X_1, \dots, X_n)}(x)) &= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(\det(\Gamma_n)) - \frac{1}{2} {}^t(x - m \mathbb{I}_n) \Gamma_n^{-1} (x - m \mathbb{I}_n). \end{aligned}$$

Les deux premiers termes ne dépendent pas de m . On veut donc maximiser le troisième terme en maximisant la vraisemblance ce qui revient à minimiser ${}^t(x - m \mathbb{I}_n) \Gamma_n^{-1} (x - m \mathbb{I}_n)$.

En dérivant ${}^t(x - m \mathbb{I}_n) \Gamma_n^{-1} (x - m \mathbb{I}_n)$ par rapport à m , on obtient ${}^t \mathbb{I}_n \Gamma_n^{-1} (x - m \mathbb{I}_n) - {}^t(x - m \mathbb{I}_n) \Gamma_n^{-1} \mathbb{I}_n = -2 {}^t \mathbb{I}_n \Gamma_n^{-1} (x - m \mathbb{I}_n)$. En remplaçant x par ${}^t(X_1, \dots, X_n)$ et en cherchant les points critiques on obtient:

$${}^t \mathbb{I}_n \Gamma_n^{-1} (X_1, \dots, X_n) = m {}^t \mathbb{I}_n \Gamma_n^{-1} \mathbb{I}_n \implies m = ({}^t \mathbb{I}_n \Gamma_n^{-1} \mathbb{I}_n)^{-1} {}^t \mathbb{I}_n \Gamma_n^{-1} (X_1, \dots, X_n),$$

puisque la matrice est inversible du fait que Γ_n^{-1} est définie positive comme Γ_n .

Par ailleurs, si on dérive une seconde fois par rapport à m , on obtient $2 {}^t \mathbb{I}_n \Gamma_n^{-1} \mathbb{I}_n > 0$, donc la fonction est convexe et admet donc le point critique comme unique minimum local donc comme unique minimum absolu.

7. On a $(X_1, \dots, X_n) \xrightarrow{\mathcal{L}} \mathcal{N}(m \mathbb{I}_n, \Gamma_n)$. Par conséquent,

$$\begin{aligned} \bullet \mathbb{E}[\hat{m}_n] &= ({}^t \mathbb{I}_n \Gamma_n^{-1} \mathbb{I}_n)^{-1} {}^t \mathbb{I}_n \Gamma_n^{-1} \mathbb{E}[{}^t(X_1, \dots, X_n)] = m ({}^t \mathbb{I}_n \Gamma_n^{-1} \mathbb{I}_n)^{-1} {}^t \mathbb{I}_n \Gamma_n^{-1} \mathbb{I}_n = m: \text{ estimateur sans biais.} \\ \bullet \text{var}(\hat{m}_n) &= \text{cov}\left(({}^t \mathbb{I}_n \Gamma_n^{-1} \mathbb{I}_n)^{-1} {}^t \mathbb{I}_n \Gamma_n^{-1} (X_1, \dots, X_n)\right) = ({}^t \mathbb{I}_n \Gamma_n^{-1} \mathbb{I}_n)^{-1} {}^t \mathbb{I}_n \Gamma_n^{-1} \text{cov}\left((X_1, \dots, X_n)\right) {}^t \left(({}^t \mathbb{I}_n \Gamma_n^{-1} \mathbb{I}_n)^{-1} {}^t \mathbb{I}_n \Gamma_n^{-1}\right) \\ &= ({}^t \mathbb{I}_n \Gamma_n^{-1} \mathbb{I}_n)^{-1} {}^t \mathbb{I}_n \Gamma_n^{-1} \Gamma_n \Gamma_n^{-1} \mathbb{I}_n ({}^t \mathbb{I}_n \Gamma_n^{-1} \mathbb{I}_n)^{-1} = ({}^t \mathbb{I}_n \Gamma_n^{-1} \mathbb{I}_n)^{-1}. \end{aligned}$$

8. $\text{var}(\bar{X}_n) = \frac{1}{n^2} \text{cov}({}^t \mathbb{I}_n {}^t(X_1, \dots, X_n)) = \frac{1}{n^2} {}^t \mathbb{I}_n \text{cov}({}^t(X_1, \dots, X_n)) \mathbb{I}_n = \frac{1}{n^2} {}^t \mathbb{I}_n \Gamma_n \mathbb{I}_n$.

Considérons le produit scalaire $\langle U, V \rangle = {}^t U V$. Notons $\Gamma_n^{1/2}$ et $\Gamma_n^{-1/2}$ les matrices obtenues à partir de la diagonalisation de Γ_n en prenant les puissances $1/2$ et $-1/2$ des valeurs propres (toutes > 0). Alors en appliquant Cauchy-Schwarz:

$$\langle U, U \rangle^2 = \langle \Gamma_n^{1/2} U, \Gamma_n^{-1/2} U \rangle^2 \leq \langle \Gamma_n^{1/2} U, \Gamma_n^{1/2} U \rangle \langle \Gamma_n^{-1/2} U, \Gamma_n^{-1/2} U \rangle.$$

En réécrivant, on obtient le résultat demandé.

Il suffit d'appliquer l'inégalité précédente à $U = \mathbb{I}_n$. □