
DS8 - 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 les librairies suivantes sont importées sous leur 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`

Notations

- Tout au long du sujet $(\Omega, \mathcal{F}, \mathbb{P})$ désignera un espace probabilisé et les variables aléatoires utilisées seront toutes définies sur cet espace probabilisé. Sous réserve d'existence, l'espérance mathématique d'une variable aléatoire réelle X sera notée $\mathbb{E}(X)$ et sa variance sera notée $\mathbb{V}(X)$.
- Pour un événement A , on notera $\mathbb{P}_B(A)$ la probabilité conditionnelle de A sachant B où B est un événement non négligeable.

Si T est un estimateur de θ (respectivement de $g(\theta)$), on définit :

- le biais de T par

$$b_\theta(T) = \mathbb{E}(T) - \theta \quad (\text{respectivement,} \quad b_{g(\theta)}(T) = \mathbb{E}(T) - g(\theta))$$

sous réserve de l'existence de $\mathbb{E}(T)$,

- et le risque quadratique de T par

$$r_\theta(T) = \mathbb{E}((T - \theta)^2) \quad (\text{respectivement,} \quad r_{g(\theta)}(T) = \mathbb{E}((T - g(\theta))^2))$$

sous réserve de l'existence de $\mathbb{V}(T)$.

On dit que T est *sans biais* lorsque son biais est nul.

On considérera qu'un estimateur T est meilleur (au sens large) qu'un autre estimateur T' en termes de risque quadratique si le risque quadratique de T est inférieur ou égal à celui de T' .

Le sujet est composé de quatre parties. Les parties I, II, III et IV.1 sont **indépendantes**. Il s'agit de variations autour de la notion de risque quadratique en théorie de l'estimation.

I. Premier problème d'estimation

Dans ce premier problème d'estimation, on dispose d'une seule observation notée X .

On suppose que X admet pour densité la fonction f_θ définie sur \mathbb{R} par :

$$f_\theta : x \mapsto \begin{cases} \frac{k+1}{\theta^{k+1}} x^k & \text{si } x \in [0, \theta] \\ 0 & \text{sinon} \end{cases}$$

où k est un entier naturel non nul et θ un paramètre réel inconnu strictement positif que l'on souhaite estimer.

1. Montrer que f_θ est bien une densité de probabilité.
2. Calculer $\mathbb{E}(X)$.
3. Déterminer λ_0 un réel dépendant uniquement de k tel que $\lambda_0 X$ soit un estimateur de θ sans biais.
4. Calculer $\mathbb{V}(X)$.
5. Démontrer que pour tout T estimateur de θ admettant une variance :

$$r_\theta(T) = (\mathbb{E}(T) - \theta)^2 + \mathbb{V}(T)$$

6. Donner la valeur de $r_\theta(\lambda_0 X)$.

Le but de la fin de cette partie I est de déterminer un estimateur de θ ayant un plus petit risque quadratique que celui de $\lambda_0 X$.

7. En utilisant I.5 montrer que pour tout λ réel

$$r_\theta(\lambda X) = \theta^2 Q(\lambda)$$

où Q est un polynôme de degré 2 dont les coefficients ne dépendent que de k .

8. Montrer que la fonction $\lambda \mapsto Q(\lambda)$ atteint son minimum en un unique réel noté λ^* que l'on exprimera en fonction de k .
9. Conclure sur le but recherché.

II. Second problème d'estimation

Dans ce second problème d'estimation, on dispose de n observations indépendantes ($n \geq 2$) notées X_1, \dots, X_n de même loi de Poisson de paramètre θ inconnu ($\theta \in]0, +\infty[$). On souhaite estimer le paramètre $\exp(-\theta)$. On définit pour tout i élément de $\llbracket 1, n \rrbracket$ la variable aléatoire Y_i par :

$$Y_i : \omega \mapsto \begin{cases} 1 & \text{si } X_i(\omega) = 0 \\ 0 & \text{sinon} \end{cases}$$

Puis on note :

$$\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$$

10. Pour tout i élément de $\llbracket 1, n \rrbracket$, donner la loi de Y_i .

11. Donner la loi de $\sum_{i=1}^n Y_i$, puis montrer que $\mathbb{E}(\bar{Y}_n) = \exp(-\theta)$.

On dira dans ce cas que \bar{Y}_n est un estimateur sans biais de $\exp(-\theta)$.

12. Calculer $\mathbb{V}(\bar{Y}_n)$.

Pour tout k élément de $\llbracket 1, n \rrbracket$ on définit $S_k = \sum_{i=1}^k X_i$.

13. Rappeler sans démonstration la loi de S_k pour tout k élément de $\llbracket 1, n \rrbracket$.

On définit jusqu'à la fin de cette partie II pour tout j entier naturel :

$$\varphi(j) = \mathbb{P}_{[S_n=j]}([X_1 = 0])$$

14. Montrer que pour tout j entier naturel :

$$\varphi(j) = \left(1 - \frac{1}{n}\right)^j$$

On a donc $\varphi(j)$ indépendant du paramètre θ inconnu.

D'après la question II.14, on peut définir l'estimateur :

$$\varphi(S_n) = \left(1 - \frac{1}{n}\right)^{S_n}$$

15. Montrer que $\varphi(S_n)$ admet une espérance et que $\mathbb{E}(\varphi(S_n)) = \exp(-\theta)$.

On dira dans ce cas que $\varphi(S_n)$ est un estimateur sans biais de $\exp(-\theta)$.

16. Montrer que $\varphi(S_n)$ admet une variance vérifiant :

$$\mathbb{V}(\varphi(S_n)) = \exp(-2\theta) \left(\exp\left(\frac{\theta}{n}\right) - 1 \right)$$

17. On souhaite comparer les performances de \bar{Y}_n et $\varphi(S_n)$ en tant qu'estimateurs de $\exp(-\theta)$.

a) Démontrer :

$$1 \leq \frac{\exp(\theta) - 1}{\theta} \leq \exp(\theta)$$

b) Soit la fonction $h : [0, 1] \rightarrow \mathbb{R}$ définie par :

$$h(t) = t \exp(\theta) + (1 - t) - \exp(t\theta)$$

pour tout $t \in [0, 1]$. Étudier les variations de h .

c) En déduire :

$$\exp\left(\frac{\theta}{n}\right) \leq \frac{\exp(\theta)}{n} + \frac{n-1}{n}$$

puis l'inégalité :

$$\mathbb{V}(\varphi(S_n)) \leq \mathbb{V}(\bar{Y}_n)$$

d) Comparer les risques quadratiques de \bar{Y}_n et $\varphi(S_n)$ en tant qu'estimateurs de $\exp(-\theta)$.

18. a) Compléter la fonction **Python** suivante pour qu'elle simule les estimateurs \bar{Y}_n et $\varphi(S_n)$.

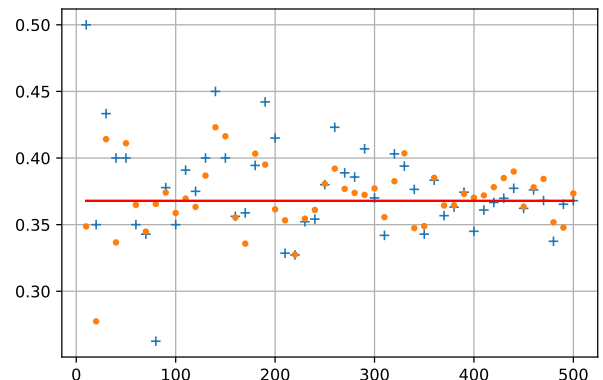
```

1 def simulEstimateurs(theta, n):
2     Xobs = _____
3     Y = np.zeros(n)
4     for k in range(n):
5         if _____:
6             Y[k] = 1
7     Ybar = np.mean(Y)
8     phi = _____
9     return [Ybar, phi]
```

b) On exécute le script **Python** ci-dessous. Commenter en lien avec la question 17.

```

1 YbarListe=[]; phiListe=[]; theta=1
2 nListe = [10*i for i in range(1, 51)]
3 for n in nListe:
4     [y,p] = simulEstimateurs(theta, n)
5     YbarListe.append(y)
6     phiListe.append(p)
7 plt.plot(nListe, YbarListe, '+')
8 plt.plot(nListe, phiListe, '.')
9 plt.plot(nListe,
10          [np.exp(-theta) for k in nListe])
11 plt.show()
```



On reprendra à la fin de la partie IV l'étude de $\varphi(S_n)$.

III. Information de Fisher

III.1. Cas discret

Dans cette section III.1, on considère I un intervalle de \mathbb{R} , θ un paramètre inconnu appartenant à I et X une variable aléatoire à valeurs dans \mathbb{N} ($X(\Omega) \subset \mathbb{N}$). On suppose qu'il existe une fonction p définie sur $I \times X(\Omega)$ telle que pour tout k élément de $X(\Omega)$:

$$\mathbb{P}([X = k]) = p(\theta, k)$$

et vérifiant, pour tout k de $X(\Omega)$, $\theta \mapsto p(\theta, k)$ est dérivable sur I .

On note de plus $h : (\theta, k) \mapsto \ln(p(\theta, k))$.

On définit enfin, sous réserve d'existence l'**information de Fisher** de X par :

$$I_X(\theta) = \sum_{k \in X(\Omega)} (\partial_1(\ln \circ p)(\theta, k))^2 p(\theta, k)$$

19. Dans cette question **19**, on considère X une variable aléatoire qui suit la loi de Bernoulli de paramètre θ (où $\theta \in]0, 1[$).

On a alors $X(\Omega) = \{0, 1\}$, $\mathbb{P}([X = 1]) = p(\theta, 1) = \theta$, $\mathbb{P}([X = 0]) = p(\theta, 0) = 1 - \theta$ et :

$$I_X(\theta) = (\partial_1(h)(\theta, 1))^2 p(\theta, 1) + (\partial_1(h)(\theta, 0))^2 p(\theta, 0)$$

Montrer :

$$I_X(\theta) = \frac{1}{\theta(1-\theta)}$$

20. Dans cette question **20**, on considère X une variable aléatoire qui suit la loi binomiale de paramètres N et θ ($N \in \mathbb{N}^*$, $\theta \in]0, 1[$).

a) Montrer :

$$I_X(\theta) = \frac{1}{(\theta(1-\theta))^2} \sum_{k=0}^N (k - N\theta)^2 \binom{N}{k} \theta^k (1-\theta)^{N-k}$$

b) En déduire :

$$I_X(\theta) = \frac{\mathbb{V}(X)}{(\theta(1-\theta))^2}$$

puis donner la valeur de $I_X(\theta)$.

21. Dans cette question **21**, on considère X une variable aléatoire qui suit la loi de Poisson de paramètre θ ($\theta \in]0, +\infty[$). Puisque $X(\Omega) = \mathbb{N}$, on a sous réserve de convergence :

$$I_X(\theta) = \sum_{k=0}^{+\infty} (\partial_1(h)(\theta, k))^2 p(\theta, k)$$

a) Montrer que la série de terme général $(\partial_1(h)(\theta, k))^2 p(\theta, k)$ converge et calculer sa somme $I_X(\theta)$.

b) Justifier :

$$I_X(\theta) = \mathbb{E} \left((\partial_1(h)(\theta, X))^2 \right)$$

III.2. Cas d'une variable gaussienne

Soit X une variable aléatoire qui suit la loi normale de moyenne θ ($\theta \in \mathbb{R}$) et de variance 1 dont la densité est notée $x \mapsto f(\theta, x)$. On définit sous réserve de convergence l'**information de Fisher** de X par :

$$I_X(\theta) = \int_{-\infty}^{+\infty} (\partial_1(\ln \circ f)(\theta, x))^2 f(\theta, x) dx$$

22. Montrer que sous réserve de convergence :

$$I_X(\theta) = \int_{-\infty}^{+\infty} (x - \theta)^2 f(\theta, x) dx$$

23. En déduire l'existence et la valeur de $I_X(\theta)$.

24. Justifier :

$$I_X(\theta) = \mathbb{E} \left((\partial_1(\ln \circ f)(\theta, X))^2 \right)$$

IV. Minoration du risque quadratique

IV.1. Inégalité de Cramer-Rao

Dans cette section IV.1, on considère I un intervalle de \mathbb{R} , θ un paramètre inconnu appartenant à I et X une variable aléatoire telle que $X(\Omega) = \llbracket 0, N \rrbracket$ ($N \in \mathbb{N}$). On suppose qu'il existe une fonction p définie sur $I \times X(\Omega)$ telle que pour tout $k \in \llbracket 0, N \rrbracket$:

$$\mathbb{P}([X = k]) = p(\theta, k)$$

et vérifiant :

- pour tout $k \in \llbracket 0, N \rrbracket$, $\theta \mapsto p(\theta, k)$ est dérivable sur I ,
- l'information de Fisher de X notée $I_X(\theta)$ définie dans la partie III est non nulle pour tout $\theta \in I$.

Le but de la section IV.1 est de démontrer l'inégalité suivante due à Cramer et Rao.

Théorème 1. (de Cramer-Rao)

Soit $f(X)$ un estimateur sans biais de $g(\theta)$ à savoir tel que $\mathbb{E}(f(X)) = g(\theta)$ où g est dérivable sur I .
On a alors :

$$\mathbb{V}(f(X)) \geq \frac{(g'(\theta))^2}{I_X(\theta)}$$

25. Montrer que pour tout θ élément de I :

$$\sum_{k=0}^N \partial_1(p)(\theta, k) = 0$$

26. En déduire que pour tout θ élément de I :

$$\mathbb{E}(\partial_1(h)(\theta, X)) = 0 \quad (E)$$

27. En dérivant partiellement par rapport à θ les deux membres de l'égalité (E), montrer que pour tout θ élément de I :

$$\mathbb{E}(\partial_{1,1}^2(h)(\theta, X)) = -\mathbb{E}\left((\partial_1(h)(\theta, X))^2\right)$$

28. Montrer que pour tout θ élément de I :

$$g'(\theta) = \sum_{k=0}^N f(k) (\partial_1(h)(\theta, k)) p(\theta, k)$$

puis :

$$g'(\theta) = \mathbb{E}((f(X) - g(\theta))(\partial_1(h)(\theta, X)))$$

29. On pose pour tout t réel :

$$L(t) = \mathbb{E}\left(\left((f(X) - g(\theta)) + t(\partial_1(h)(\theta, X))\right)^2\right)$$

- Développer le polynôme L suivant les puissances décroissantes de t .
- Calculer le discriminant Δ de L et justifier : $\Delta \leq 0$.
- En déduire l'inégalité de Cramer-Rao.

IV.2. Extension du théorème de Cramer-Rao

On reprend dans cette section IV.2 les notations et hypothèses de la partie II. On admet que, dans ce contexte, le théorème de Cramer-Rao peut se généraliser comme suit :

Théorème 2. (de Cramer-Rao)

Soit $T_n = f(X_1, \dots, X_n)$ un estimateur sans biais de $g(\theta)$ à savoir tel que $\mathbb{E}(f(X_1, \dots, X_n)) = g(\theta)$ où g est dérivable sur $]0, +\infty[$. On a alors :

$$\mathbb{V}(T_n) \geq \frac{(g'(\theta))^2}{n I_{X_1}(\theta)}$$

où $I_{X_1}(\theta)$ est l'information de Fisher d'une variable aléatoire de loi de Poisson de paramètre θ définie et calculée à la partie III.

Il s'agit dans cette section d'exploiter cette nouvelle inégalité de Cramer-Rao. On note :

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

30. Calculer $\mathbb{E}(\bar{X}_n)$ et $\mathbb{V}(\bar{X}_n)$.

31. Dédurre de la généralisation de Cramer-Rao, que \bar{X}_n a le plus petit risque quadratique parmi les estimateurs sans biais de θ .

32. Montrer que pour $g(\theta) = \exp(-\theta)$ où $\theta \in]0, +\infty[$:

$$\mathbb{V}(\varphi(S_n)) \underset{n \rightarrow +\infty}{\sim} \frac{(g'(\theta))^2}{n I_{X_1}(\theta)}$$

33. Que prouve ce résultat en terme d'optimalité de $\varphi(S_n)$ dans l'estimation de $\exp(-\theta)$?

34. À la lumière de la partie II, peut-on conclure que lorsque n est grand $\varphi(S_n)$ est le meilleur estimateur de $\exp(-\theta)$ en terme de risque quadratique ?