

TD 4 : MAXIMUM DE VRAISEMBLANCE

EXERCICE 1 (Modèles de translation) Soit X_1, \dots, X_n i.i.d. de densité commune $f_\theta(x) = f(x - \theta)$, $\theta \in \mathbb{R}$. Déterminer l'estimateur du maximum de vraisemblance dans les cas suivants :

1. f est la densité de la loi $\mathcal{N}(0, \sigma^2)$ (σ connue) ;
2. $f(x) = \frac{1}{2}e^{-|x|}$ (loi de Laplace) ;
3. $f(x) = \frac{3}{4}(1 - x^2)$ sur $[-1, 1]$.

EXERCICE 2 (Loi discrète) Soit $\theta \in]0, 1[$ un paramètre inconnu, on note X une variable aléatoire de loi définie par

$$\mathbb{P}_\theta(X = k) = (k + 1)(1 - \theta)^2\theta^k, \text{ pour tout } k \in \mathbb{N}.$$

On donne

$$\mathbb{E}_\theta[X] = \frac{2\theta}{1 - \theta} \quad \text{et} \quad \text{Var}_\theta[X] = \frac{2\theta}{(1 - \theta)^2}.$$

On souhaite estimer θ à partir d'un échantillon X_1, \dots, X_n de même loi que X .

1. Donner un estimateur $\hat{\theta}_n$ de θ par la méthode des moments.
2. L'estimateur du maximum de vraisemblance $\tilde{\theta}_n$ de θ est-il bien défini ?
3. Étudier la consistance de $\hat{\theta}_n$ et déterminer sa loi limite.

EXERCICE 3 (Loi exponentielle translatée) On observe un échantillon X_1, \dots, X_n dont la loi admet la densité

$$f_\theta(x) = \exp(-(x - \theta))\mathbb{1}_{[\theta, +\infty[}(x),$$

où θ est un paramètre réel inconnu.

1. Quels estimateurs de θ pouvez-vous proposer en utilisant les méthodes usuelles ?
2. Déterminer la loi de $n(\hat{\theta}_n - \theta)$, où $\hat{\theta}_n$ est l'estimateur du maximum de vraisemblance. En déduire un intervalle de confiance pour θ de niveau $1 - \alpha$, pour $\alpha \in]0, 1[$.
3. On souhaite tester au niveau α

$$H_0 : \theta \geq 0 \text{ contre } H_1 : \theta < 0.$$

- (a) Construire un test à partir de l'intervalle de confiance de la question 2, calculer sa puissance et donner son allure (pour n et α fixés). Quelle est sa taille α^* ?
 - (b) Proposer un autre test qui soit, lui, de taille α .
 - (c) Calculer la fonction puissance du test. La représenter en fonction de θ pour n et α fixés.
 - (d) Comment varie la puissance en fonction de α ? en fonction de n ?
4. Proposer un test de niveau α de H_0 : « La loi de X_1 est une loi exponentielle » contre H_1 : « La loi de X_1 n'est pas une loi exponentielle ». Calculer la puissance de ce test.

EXERCICE 4 (Densité en escalier) Pour tout paramètre réel θ , on définit sur \mathbb{R} la fonction

$$f_\theta(x) = \begin{cases} 1 - \theta & \text{si } -1/2 < x \leq 0, \\ 1 + \theta & \text{si } 0 < x \leq 1/2, \\ 0 & \text{sinon.} \end{cases}$$

1. Quelles conditions doit vérifier θ pour que f_θ soit une densité par rapport à la mesure de Lebesgue sur \mathbb{R} ?
2. Soit (X_1, \dots, X_n) un échantillon de densité f_θ . L'estimateur du maximum de vraisemblance $\hat{\theta}_n$ de θ est-il bien défini ?
3. Sous les conditions de la question 1, $\hat{\theta}_n$ est-il sans biais ? consistant ? asymptotiquement normal ?

EXERCICE 5 (Données censurées) Soient $\theta^* > 0$ et X_1, \dots, X_n un échantillon de variables aléatoires i.i.d. distribuées selon $\mathcal{E}(\theta^*)$. On appelle, pour tout $i \in \{1, \dots, n\}$, $Y_i = \min(X_i, 1)$.

1. Déterminer la loi de Y_1 , notée P_{θ^*} . La variable aléatoire Y_1 est-elle discrète ? Est-elle à densité ?
2. On se donne le modèle statistique $(P_\theta)_{\theta > 0}$. Donner une mesure dominante μ de ce modèle. Pour tout $\theta > 0$, déterminer une densité (notée f_θ) de P_θ par rapport à μ .
3. Calculer l'EMV $\hat{\theta}_n$ de θ^* .
4. Montrer que $\hat{\theta}_n$ est un estimateur consistant de θ^* .
5. (★) À l'aide de l'indication ci-dessous, montrer que $\hat{\theta}_n$ est asymptotiquement normal.

Indication (TCL et méthode delta multivariés) : Soit $(U_n)_{n \geq 1}$ une suite de vecteurs aléatoires à valeurs dans \mathbb{R}^2 , i.i.d et de carrés intégrables. En notant $\bar{U}_n = \frac{1}{n} \sum_{i=1}^n U_i$, on a alors

$$\sqrt{n} (\bar{U}_n - \mathbb{E}[U_1]) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}(0, \text{Var}(U_1)),$$

où $\text{Var}(U_1)$ est la matrice de variance-covariance de $U_1 = \begin{pmatrix} U_{(1)} \\ U_{(2)} \end{pmatrix}$, définie par

$$\text{Var}(U_1) = \begin{pmatrix} \text{Var}(U_{(1)}) & \text{Cov}(U_{(1)}, U_{(2)}) \\ \text{Cov}(U_{(1)}, U_{(2)}) & \text{Var}(U_{(2)}) \end{pmatrix}.$$

De plus, pour une fonction $\varphi : \mathbb{R}^2 \rightarrow \mathbb{R}$ différentiable et telle que $\nabla \varphi(\mathbb{E}[U_1]) \neq 0$, la méthode delta assure que

$$\sqrt{n} (\varphi(\bar{U}_n) - \varphi(\mathbb{E}[U_1])) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}\left(0, \nabla \varphi(\mathbb{E}[U_1])^\top \text{Var}(U_1) \nabla \varphi(\mathbb{E}[U_1])\right).$$

6. (★) Montrer que le modèle est régulier. Que pouvez-vous en déduire ?

EXERCICE 6 (★) (Loi uniforme translatée) Soient X_1, \dots, X_n des observations i.i.d. de densité f_θ , θ réel, donnée par $f_\theta(x) = \mathbb{1}_{[\theta, \theta+1]}(x)$.

1. Comment peut-on définir l'estimateur du maximum de vraisemblance de θ ?
2. Déterminer les lois exactes et asymptotiques de $n(X_{(1)} - \theta)$ et de $n(\theta + 1 - X_{(n)})$.
3. Montrer que pour toute variable aléatoire Y à valeurs dans \mathbb{R}^+ et tout $k \in \mathbb{N}^*$ (**on pourra retenir cette formule, ainsi que sa preuve**)

$$\mathbb{E}[Y^k] = k \int_0^{+\infty} t^{k-1} \mathbb{P}(Y > t) dt.$$

En déduire les deux premiers moments de $n(X_{(1)} - \theta)$.

4. Étudier la consistance des estimateurs du maximum de vraisemblance de θ .
5. Majorer $\sup_{n \geq 1} n^2 R(\hat{\theta}_n, \theta)$ pour tout estimateur du maximum de vraisemblance $\hat{\theta}_n$ de θ .
6. Proposer un intervalle de confiance pour θ au niveau 95%. Donner l'application numérique pour $\alpha = 0.05$ et $\hat{\theta} = 3$, avec $n = 10$ et $n = 100$.

EXERCICE 7 (★) (Mélange de lois uniformes) Soit X_1, \dots, X_n i.i.d. dont la loi admet la densité

$$f_{\theta, \lambda}(x) = \lambda \mathbb{1}_{[0,1]}(x) + \frac{1-\lambda}{\theta} \mathbb{1}_{[0,\theta]}(x),$$

où $\lambda \in [0, 1[$ et $\theta \geq 1$ sont deux paramètres inconnus.

1. Montrer que, si on exclut une valeur de θ à préciser, le modèle est identifiable. Sauf indication contraire, on supposera cette condition vérifiée dans la suite.
2. Montrer que la densité $f_{\theta, \lambda}$ peut s'écrire, pour tout $x \in \mathbb{R}$,

$$f_{\theta, \lambda}(x) = \left(\lambda + \frac{1-\lambda}{\theta} \right)^{\mathbb{1}_{[0,1]}(x)} \left(\frac{1-\lambda}{\theta} \right)^{\mathbb{1}_{[1,\theta]}(x)} \mathbb{1}_{[0,\theta]}(x).$$

3. Déterminer l'EMV de λ lorsque θ est connu, l'EMV de θ lorsque λ est connu, puis l'EMV de (θ, λ) lorsque les 2 paramètres sont inconnus.
4. Étudier la consistance et la loi limite de l'EMV de θ , que λ soit ou non connu.
5. Étudier la consistance de l'EMV de λ lorsque θ est connu, puis lorsque θ est inconnu. Que se passe-t-il si $\theta = 1$?
6. Construire un test de $H_0 : \theta = 1$ contre $H_1 : \theta > 1$ au niveau α , avec $\alpha \in]0, 1[$. Pour $\lambda \in [0, 1[$ et $\theta > 1$ fixés, étudier la limite de la puissance lorsque n tend vers l'infini.

EXERCICE 8 (★) (Méthode de capture-recapture) On souhaite estimer le nombre de poissons, $N \in \mathbb{N}^*$, vivant dans un bassin. Pour ce faire, nous mettons en place le protocole suivant :

- on prélève au chalut un groupe de poissons et l'on note $k \in \llbracket 1, N \rrbracket$ leur nombre ;
- ces poissons sont marqués puis relâchés dans le bassin ;
- pour un nombre de prises $n \in \mathbb{N}^*$ défini et aussi grand qu'on veut, on pêche au hasard puis on relâche successivement n poissons ;
- de retour à quai, on reporte le nombre $x \in \llbracket 0, n \rrbracket$ de poissons marqués parmi les n que l'on a pêchés puis relâchés.

On fait ensuite appeler à un statisticien afin de proposer une estimation de N à partir de cette expérience.

1. Quelle est ici la quantité mesurée (et enregistrée)? En déduire une modélisation probabiliste de l'expérience, puis un modèle statistique paramétrique indicé par un intervalle Θ .
2. L'estimateur du maximum de vraisemblance est-il défini? Proposer un estimateur \hat{N} de N .
3. Montrer que \hat{N} est convergent et asymptotiquement normal lorsque n tend vers l'infini.
4. Au regard de la question précédente, quelle valeur de k est pertinente pour l'expérience?

EXERCICE 9 (★) (Gauss et la loi normale¹) En statistique, la loi normale a été introduite par Gauss pour un problème d'estimation de paramètre en astronomie. Dans le langage moderne, voici le problème :

1. Cet exercice est inspiré de l'article *Plaidoyer pour la loi normale*, de Aimé Fuchs.

soit X_1, \dots, X_n i.i.d. de densité commune $f_\theta(x) = f(x - \theta)$, $\theta \in \mathbb{R}$. On a donc affaire à un modèle de translation dans lequel f correspond à la densité de la loi des erreurs et sur laquelle on fait les hypothèses “naturelles” suivantes : f est paire, de classe C^1 et strictement positive sur \mathbb{R} . La question à laquelle Gauss souhaitait répondre est la suivante : quelles sont les densités f pour lesquelles l’estimateur des moindres carrés correspond à celui du maximum de vraisemblance ? Comme nous allons le voir, seules les lois normales centrées satisfont cette propriété.

1. Montrer que \bar{X}_n est l’estimateur des moindres carrés.
2. On note $g = \ln f$. Montrer que si \bar{X}_n est l’EMV alors pour tout n -uplet (x_1, \dots, x_n) de réels, la fonction g doit vérifier $\sum_{i=1}^n g'(x_i - \bar{x}_n) = 0$.
3. En prenant $x_2 = \dots = x_n = x_1 - nt$, en déduire que, $\forall t \in \mathbb{R}^*$, $\forall k \in \mathbb{N}^*$, on a $\frac{g'(kt)}{kt} = \frac{g'(t)}{t}$.
4. En déduire qu’il existe une constante c telle que, pour tout $t \in \mathbb{R}^*$, $\frac{g'(t)}{t} = c$.
5. Conclure.

EXERCICE 10 (Maximum de vraisemblance et donnée aberrante) Le but de cet exercice est de montrer sur un exemple très simple que l’EMV n’est pas robuste à la présence de données aberrantes. Afin de simplifier les calculs, on fera les approximations suivantes : $\Phi^{-1}(0.975) \approx 2$, $\Phi(-3) \approx 10^{-3}$, $\Phi(-7) \approx 0$ et $99/10^4 \approx 10^{-2}$.

1. Soit $\theta \in \mathbb{R}$ un paramètre inconnu et X_1, \dots, X_{100} i.i.d. selon une loi $\mathcal{N}(\theta, 1)$. Si on note $\hat{I} = [\bar{X}_{100} - 0.2 ; \bar{X}_{100} + 0.2]$, montrer que $\mathbb{P}(\theta \in \hat{I}) \approx 0.95$. On suppose désormais $\theta = 0$, ce qui donne donc $\mathbb{P}(0 \in \hat{I}) \approx 0.95$, i.e. la vraie valeur du paramètre se situe 95% du temps dans l’intervalle (aléatoire) \hat{I} .
2. On suppose maintenant la présence d’une seule donnée aberrante (dysfonctionnement du dispositif d’enregistrement des données, etc.) et on veut voir l’incidence de celle-ci sur le résultat précédent. Plus précisément, on suppose que X_1, \dots, X_{99} sont i.i.d. selon une loi $\mathcal{N}(0, 1)$ et que $X_{100} = 50$. Donner la loi de \bar{X}_{100} sachant $X_{100} = 50$.
3. On note toujours $\hat{I} = [\bar{X}_{100} - 0.2 ; \bar{X}_{100} + 0.2]$. En déduire que $\mathbb{P}(0 \in \hat{I} \mid X_{100} = 50) \approx 10^{-3}$ et conclure.

EXERCICE 11 (*) (Non-robustesse de l’EMV) Sur l’ensemble des densités de probabilité sur \mathbb{R} (ensemble des classes d’équivalence de fonctions de L_1 positives et d’intégrale 1), on définit le carré de la distance de Hellinger h entre f et g par

$$h^2(f, g) = \frac{1}{2} \int_{\mathbb{R}} \left(\sqrt{f(x)} - \sqrt{g(x)} \right)^2 dx.$$

1. Montrer que

$$h^2(f, g) = 1 - \int_{\mathbb{R}} \sqrt{f(x)g(x)} dx.$$

Vérifier que h est une distance sur l’ensemble des densités sur \mathbb{R} (i.e. positivité avec nullité si et seulement si égalité, symétrie et inégalité triangulaire), avec de plus $0 \leq h(f, g) \leq 1$.

2. Soit $\theta > 0$ et f_θ la densité de la loi uniforme sur $[0, \theta]$. Montrer que $h^2(f_\theta, f_{\theta'}) = 1 - \sqrt{\frac{\theta}{\theta'}}$ si $\theta \leq \theta'$.
3. Soit X_1, \dots, X_n i.i.d. de loi uniforme sur $[0, \theta]$ pour un $\theta > 0$ inconnu. Donner l’estimateur du maximum de vraisemblance $\hat{\theta}_n$, calculer $\mathbb{E}_\theta \left[\sqrt{\hat{\theta}_n} \right]$ et en déduire que $\mathbb{E}_\theta \left[h^2(f_\theta, f_{\hat{\theta}_n}) \right] = \frac{1}{2n+1}$.

4. Soit la densité

$$f_n^*(x) = 10 \left(1 - \frac{1}{n}\right) \mathbf{1}_{0 \leq x \leq 1/10} + \frac{10}{n} \mathbf{1}_{9/10 \leq x \leq 1}.$$

(a) Montrer que si, pour tout n , Y_n a pour densité f_n^* , alors la suite (Y_n) converge en loi et identifier la limite.

(b) Calculer $h^2(f_n^*, f_{1/10})$.

5. Soit $n > 1$. Les variables X_1, \dots, X_n sont désormais i.i.d. de densité f_n^* , mais on les croit toujours i.i.d. suivant une densité uniforme $(f_\theta)_{\theta > 0}$, avec θ inconnu. En particulier, l'estimateur $\hat{\theta}_n$ est le même que celui défini ci-dessus.

(a) Montrer qu'avec probabilité $1 - (1 - 1/n)^n$, on a : $9/10 \leq \hat{\theta}_n \leq 1$.

(b) Montrer que, sur l'événement $\{9/10 \leq \hat{\theta}_n \leq 1\}$, on a :

$$h^2(f_n^*, f_{\hat{\theta}_n}) \geq 1 - \frac{1}{3} \sqrt{1 - \frac{1}{n}} - \frac{1}{3\sqrt{n}}.$$

(c) En notant $\mathbb{E}_n^*[\cdot]$ l'espérance par rapport au modèle où les données X_1, \dots, X_n sont i.i.d. de densité f_n^* , en déduire que $\mathbb{E}_n^*[h^2(f_n^*, f_{\hat{\theta}_n})] \geq u_n$, avec $\lim_{n \rightarrow \infty} u_n = 2/3 \times (1 - 1/e)$.

(d) Conclure quant à la robustesse de l'estimateur du maximum de vraisemblance à une mauvaise spécification du modèle.