

# Master 2 Biostatistique – UE STA305

## Travaux Dirigés

Boris Hejblum

### Exercice 1

Les variables aléatoires  $Y_i, i = 1, \dots, n$  sont indépendantes et identiquement distribuées (*iid*) suivant une loi Normale de paramètres  $\theta$  et  $\sigma^2$ . La densité de la loi de Normale est :  $f_{\theta, \sigma^2}(y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\theta)^2}{2\sigma^2}}$ . On considérera  $\sigma^2$  connu.

1. Écrire la vraisemblance et la log-vraisemblance de l'échantillon  $y_i, i = 1, \dots, n$ , en faisant apparaître  $\bar{y}_{(n)} = \frac{1}{n} \sum_{i=1}^n y_i$  sous la forme  $(\theta - \bar{y}_{(n)})^2$ . Attention : on rappelle qu'une somme de nombre au carré n'est pas égale au carré de la somme de ces nombres...
2. Écrire la dérivée première et seconde de la log-vraisemblance par rapport à  $\theta$  et l'information de Fisher pour  $\theta$ .
3. Quel est la loi *a priori* de Jeffrey pour  $\theta$ ? Est-ce qu'il définit densité propre ou impropre ?
4. En prenant cette loi *a priori*, écrire le numérateur de la loi *a posteriori* de  $\theta$ . En déduire la distribution *a posteriori* de  $\theta$ .
5. On observe un deuxième échantillon  $\{y_i\}, i = n+1, \dots, 2n$  *iid* de même loi que le premier échantillon. Quelle est la distribution *a posteriori* de  $\theta$  en prenant un *a priori* uniforme ? Faire le calcul de deux façons:
  - (a) en considérant que l'on a un échantillon *iid* de taille  $2n$
  - (b) en utilisant la distribution *a posteriori* obtenue pour le premier échantillon comme distribution *a priori* pour le second échantillon.

### Exercice 2

On considère les réalisations  $\mathbf{x} = \{x_1, \dots, x_n\}$  d'une suite de variables aléatoires *iid*  $\{X_i\}_{i=1, \dots, n}$  réelles et supérieures à 1, dont la loi  $P_\theta$  est supposée connue à un paramètre  $\theta > 0$  près. Cette loi  $P_\theta$  est une loi continue, appelée loi de Pareto de paramètres  $(\theta + 1, 1)$  dont la densité est définie, pour  $x > 1$ , par :

$$f_\theta(x) = \frac{\theta + 1}{x^{\theta+2}}$$

1. L'*a priori* utilisé pour  $\theta$  est une loi exponentielle de paramètre 1, dont la fonction de densité s'écrit :  $g(\theta) = e^{-\theta}$ . Écrire le modèle bayésien associé.
2. Montrer que la densité de la loi *a posteriori* de  $\theta|\mathbf{x}$ , notée  $p(\theta|\mathbf{x})$ , est proportionnelle à :

$$\exp(-\theta) (\theta + 1)^n \left( \prod_{i=1}^n x_i^{-\theta} \right) \quad ; \quad \theta > 0$$

3. Proposer un algorithme de Metropolis-Hastings indépendant pour estimer la loi *a posteriori* de  $\theta|X_1, \dots, X_n$ . On prendra comme loi instrumentale la loi *a priori* de  $\theta$ . Expliciter l'estimateur Bayésien de  $\theta$  construit pour le coût quadratique. Ne pas oublier de faire apparaître les calculs et la formule de la probabilité d'acceptation.
4. Quel résultat théorique garantit sa convergence ? Expliquer brièvement.

### Exercice 3

On considère les réalisations  $\mathbf{x} = \{x_1, \dots, x_n\}$  d'une suite de variables aléatoires *iid*  $\{X_i\}_{i=1, \dots, n}$  suivant une loi exponentielle de paramètre  $\lambda : \mathcal{E}(\lambda)$ , où  $\lambda > 0$  est inconnu. On prend comme loi *a priori* sur  $\lambda$  la loi Gamma  $\mathcal{G}(\alpha, \beta)$  dont la densité s'écrit :

$$g(\lambda) = \lambda^{\alpha-1} \frac{\beta^\alpha e^{-\beta\lambda}}{\Gamma(\alpha)}$$

1. Écrire est le modèle bayésien associé.
2. Quelle est la loi *a posteriori* correspondante?

### Exercice 4

*Propriétés utiles :*

- La densité de probabilité de la loi Beta de paramètres  $a > 0$  et  $b > 0$  évaluée en  $\theta$  est donnée par

$$\text{Beta}(\theta; a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}$$

- Soit  $\theta$  une variable aléatoire suivant une loi Beta de paramètres  $a$  et  $b$ . On a alors

$$\mathbb{E}[\theta] = \frac{a}{a+b}$$

On souhaite estimer la probabilité de contracter une maladie  $M$  dans l'hôpital  $A$ . On dispose pour cela de données de  $n_A$  patients indiquant s'ils ont ou non contracté la maladie. On note  $\mathbf{y}^A = (y_1^A, \dots, y_{n_A}^A)$  l'échantillon observé de la variable binaire définie par :

$$y_i^A = \begin{cases} 1 & \text{si le patient } i \text{ a contracté la maladie} \\ 0 & \text{sinon} \end{cases}$$

On note  $\theta_A \in [0, 1]$  la probabilité inconnue de contracter la maladie dans l'hôpital  $A$  et l'on suppose que les variables aléatoires  $\{Y_i^A\}_{i=1, \dots, n_A}$  sont *iid* conditionnellement à  $\theta_A$ .

1. Écrire la vraisemblance des données  $p(\mathbf{y}^A | \theta_A)$
2. On utilise une approche bayésienne, et l'on suppose que  $\theta_A$  suit *a priori* une distribution uniforme sur l'intervalle  $[0, 1]$ . Donner la forme de la densité *a posteriori*  $p(\theta_A | \mathbf{y}^A)$ . Montrer que celle-ci prend une forme paramétrique connue.
3. Cette densité *a posteriori* est-elle propre? Pourquoi ?
4. Calculer la loi marginale des observations  $f(\mathbf{y}^A)$ .
5. Donner la probabilité  $p(y_{n_A+1}^A = 1 | \mathbf{y}^A)$  qu'un nouveau patient  $n_A + 1$  contracte la maladie sachant  $\mathbf{y}^A = \{y_1^A, \dots, y_{n_A}^A\}$ .
6. On dispose maintenant des données  $y_1^B, \dots, y_{n_B}^B$  de contraction de la maladie pour  $n_B$  patients d'un second hôpital  $B$ . On note  $\theta_B$  la probabilité que le patient  $i$  de l'hôpital  $B$  ait contracté la maladie, et l'on suppose toujours l'indépendance conditionnellement à  $\theta_B$ . On souhaite tester l'hypothèse  $H_0$  selon laquelle les taux de contraction de la maladie sont les mêmes dans les hôpitaux  $A$  et  $B$ , versus l'hypothèse  $H_1$  que ces taux sont différents:

$$H_0 : \theta_B = \theta_A, \theta_A \sim U(0, 1) \text{ vs } H_1 : \theta_A \sim U([0, 1]) \perp \theta_B \sim U([0, 1])$$

où  $U([0, 1])$  dénote la distribution uniforme sur l'intervalle  $[0, 1]$ .

Écrire  $p(\mathbf{y}^A, \mathbf{y}^B | H_0)$  et  $p(\mathbf{y}^A, \mathbf{y}^B | H_1)$

7. En déduire le facteur de Bayes  $B_{10}$  de l'hypothèse  $H_1$  par rapport à l'hypothèse  $H_0$ , qui se définit comme le ratio des probabilités à posteriori :

$$B_{10} = \frac{p(\mathbf{y}^A, \mathbf{y}^B | H_1)}{p(\mathbf{y}^A, \mathbf{y}^B | H_0)}$$

## Exercice 5

Une personne effectue un test afin de détecter si elle est porteuse d'un virus A, potentiellement mortel, présent dans 0.1% de la population. Le test a les propriétés suivantes:

- Si la personne est porteuse du virus, le test sera positif dans 99% des cas

- Si la personne n'est pas porteuse du virus, le test sera négatif dans 95% des cas

Le résultat du test est positif.

1. Quelle est probabilité que la personne soit porteuse du virus?
2. Suite à un test positif, la personne a le choix de suivre ou non un traitement thérapeutique. Si elle est effectivement porteuse du virus, ce traitement rallongera sa durée de vie de 6 mois. Qu'elle soit ou non porteuse du virus, ce traitement est assez lourd, et l'on considère qu'il diminue de 1 mois la durée de vie. Déterminer le coût moyen de la décision de suivre ce traitement.

### Exercice 6

Dans cet exercice, nous nous proposons de voir comment il est possible de simuler des réalisations d'une loi puis de vérifier qu'elles sont bien issues de cette loi en ré-estimant les paramètres.

1. Proposer un algorithme basé sur la méthode par inversion, permettant de simuler la réalisation d'un échantillon de taille  $n$  d'une loi de Pareto de paramètres  $\lambda = 2$  et  $k = 5$ . La densité de la loi de Pareto est la suivante :  $f(x) = k \frac{\lambda^k}{x^{k+1}}$ .
2. Grâce à ce premier algorithme nous pouvons simuler un  $n$ -échantillon *iid*  $\mathbf{x} = x^{(1)}, \dots, x^{(n)}$  suivant une loi de Pareto de paramètres  $\lambda = 2$  et  $k = 5$ . Désormais, nous voulons vérifier que l'algorithme est valide et nous voulons ré-estimer le paramètre  $k$  ayant servi à simuler ces données. On suppose  $\lambda = 2$  connu et fixé. Pour cela, nous allons appliquer des méthodes bayésiennes avec l'*a priori* suivant pour  $k$  :  $\pi(k) = \frac{1}{200} e^{-\frac{k^2}{2 \cdot 100^2}} \mathbb{1}_{k \in ]0, \infty[}$ . Écrire le modèle bayésien associé puis calculer la loi *a posteriori* de  $k|\mathbf{x}$ .
3. Expliquer brièvement la logique de l'acceptation/rejet en fonction de la loi instrumentale de proposition et de la loi que l'on veut échantillonner. Quelle simplification apparait en prenant pour loi instrumentale la loi *a priori* du paramètre? Comment appelle-t-on ce phénomène ?
4. Proposer un algorithme de Métropolis-Hastings indépendant pour échantillonner la loi *a posteriori* de  $k|\mathbf{x}$ . On prendra comme loi instrumentale la loi *a priori* de  $k$ .
5. Expliciter l'estimateur Bayésien  $\hat{E}(k|X_1, \dots, X_n)$  de  $k$  construit pour le coût quadratique.