

LSTAT2040: Examen

Durée 3 heures

Les feuilles blanches sont vos feuilles de réponses. Inscrivez vos nom et prénom sur chacune d'entre elles et numérotez-les dans l'ordre de lecture. Seules ces feuilles seront corrigées. Les feuilles de couleur vous serviront de brouillon.

À la fin de l'examen, vous devez rendre toutes les feuilles que vous avez reçues, y compris les feuilles de brouillon et les feuilles vierges que vous n'avez pas utilisées.

Vous pouvez consulter le syllabus du cours.

Tout échange d'informations, sous quelque forme que ce soit, est interdit et sera considéré comme une tricherie.

Exercice 1 (16 pts)

X_1, \dots, X_n , et Y_1, \dots, Y_n sont des variables indépendantes de lois exponentielles (Exp). La densité de X_i est $f_i(x) = \lambda_i \theta \exp(-\lambda_i \theta x) I(x > 0)$, et la densité de Y_i est $g_i(y) = \lambda_i \exp(-\lambda_i y) I(y > 0)$. $\lambda_1, \dots, \lambda_n$ et θ sont des paramètres positifs.

Dans tout ce qui suit, *sauf indication contraire*, nous supposons que $\lambda_1, \dots, \lambda_n$ et θ sont tous inconnus et que l'échantillon que nous observons est (X_i, Y_i) , $i = 1, \dots, n$.

Pour information, la fonction de répartition (cdf) d'une $Exp(\lambda)$ est $(1 - \exp(-\lambda x)) I(x > 0)$. Sa moyenne est de $1/\lambda$ et sa variance de $1/\lambda^2$.

1 (1 pt)

Le modèle serait-il identifiable si nous observons seulement les X_i mais pas les Y_i ? Justifiez.

Solution

La densité de X_i , donnée par, $\lambda_i \theta \exp(-\lambda_i \theta x)$, n'est pas identifiable puisque nous pouvons, par exemple, remplacer λ_i par $2\lambda_i$ et θ par $\theta/2$ et obtenir la même fonction.

□

2 (1 pt)

Supposons que les λ_i soient tous connus. Trouver l'estimateur du maximum de vraisemblance (EMV) de θ basé sur les X_i .

Solution

Dans le cas énoncé, la log-vraisemblance et sa dérivée sont données par

$$\begin{aligned}\ell &= \sum_i \{\log(\lambda_i \theta) - \lambda_i \theta x_i\} \\ \partial_\theta \ell &= \frac{n}{\theta} - \sum_i \lambda_i X_i\end{aligned}$$

Ce qui conduit à l'estimateur suivant $\frac{1}{n^{-1} \sum_i \lambda_i X_i}$.

□

3 (2 pts)

Montrez que $\lambda_i X_i \sim \text{Exp}(\theta)$. Utilisez ce résultat pour trouver la distribution asymptotique de l'estimateur décrit en (2).

Solution

Il est facile de voir que $P(\lambda_i X_i \leq x) = 1 - \exp(-\theta x)$ et donc $\lambda_i X_i \sim \text{Exp}(\theta)$. Par TCL, nous pouvons affirmer que

$$\sqrt{n}(n^{-1} \sum_i \lambda_i X_i - 1/\theta) \xrightarrow{d} N(0, 1/\theta^2)$$

Il suffit d'appliquer la méthode Delta, avec la transformation $t \mapsto 1/t$, pour conclure que

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, \theta^2)$$

□

4 (2 pts)

Montrez que $\hat{\theta}$, l'EMV de θ , basée sur (X_i, Y_i) , découle de l'équation suivante

$$\frac{n}{\hat{\theta}} = 2 \sum_{i=1}^n \frac{R_i}{1 + \hat{\theta} R_i},$$

où $R_i = X_i/Y_i$. Quel est l'EMV de λ_i ?

Solution

$$L = \prod_i \theta \lambda_i^2 \exp(-\lambda_i y_i) \exp(-\lambda_i \theta x_i)$$

$$\ell = n \log(\theta) + 2 \sum_i \log(\lambda_i) - \sum_i \lambda_i y_i - \theta \sum_i \lambda_i x_i$$

$$\partial_\theta \ell = \frac{n}{\theta} - \sum_i \lambda_i x_i \text{ et } \partial_{\lambda_i} \ell = \frac{2}{\lambda_i} - y_i - \theta x_i$$

$\partial_{\lambda_i} \ell = 0 \iff \lambda_i = \frac{2}{y_i + \theta x_i}$ et $\partial_\theta \ell = 0 \iff \frac{n}{\theta} = \sum_i \lambda_i x_i$. Donc $\hat{\theta}$, l'EMV de θ , est bien défini par l'équation énoncée, et il s'ensuit que l'EMV de λ_i est $\hat{\lambda}_i = \frac{2}{y_i + \hat{\theta} x_i}$.

□

5 (2 pts)

Montrez que la cdf de R_i est donnée par $F_R(r) = 1 - (1 + \theta r)^{-1}$, pour $r > 0$.

Montrez que l'EMV de θ basé sur R_1, \dots, R_n est la même que celui défini ci-dessus.

Solution

$$\begin{aligned} P(R_i \leq r) &= P(X_i \leq r Y_i) = \int_0^\infty P(X_i \leq r y) g_i(y) dy \\ &= \lambda_i \int_0^\infty (1 - \exp(-\lambda_i \theta r y)) \exp(-\lambda_i y) dy \\ &= 1 - \lambda_i \int_0^\infty \exp(-(\lambda_i \theta r + \lambda_i) y) dy \\ &= 1 + \frac{\lambda_i}{\lambda_i \theta r + \lambda_i} \exp(-(\lambda_i \theta r + \lambda_i) y) \Big|_{y=0}^{y=\infty} \\ &= 1 - \frac{1}{\theta r + 1}, \quad \forall r > 0. \end{aligned}$$

La densité de R_i est donc donnée par

$$f_R(r) = \frac{dP(R_i \leq r)}{dr} = \theta (1 + \theta r)^{-2}, \quad \forall r > 0.$$

La Log-vraisemblance basée sur les R_i n'est donc rien que

$$\ell = \sum_i \{\log(\theta) - 2 \log(\theta r_i + 1)\}$$

Et

$$\partial_\theta \ell = \frac{n}{\theta} - 2 \sum_i \frac{r_i}{\theta r_i + 1}$$

D'où la conclusion concernant $\hat{\theta}$.

□

6 (2 pts)

Montrez que l'information de Fisher contenue dans R_1, \dots, R_n à propos de θ est donnée par

$$\frac{n}{3\theta^2}.$$

Indice : Pensez à utiliser une intégration par changement de variable de type $u = \theta r + 1$. Aussi, sachez que $\int_1^\infty \frac{(1-u)^2}{u^4} du = 1/3$.

Solution

Nous avons que

$$\begin{aligned} S &\equiv S_n(\theta) = \frac{n}{\theta} - 2 \sum_i \frac{r_i}{\theta r_i + 1} \\ H &= \partial_\theta S = -\frac{n}{\theta^2} + 2 \sum_i \frac{r_i^2}{(\theta r_i + 1)^2} \\ I_n &= -E(H) = \frac{n}{\theta^2} - 2nE\left(\frac{R^2}{(\theta R + 1)^2}\right) \end{aligned}$$

Et puisque

$$\begin{aligned} E\left(\frac{R^2}{(\theta R + 1)^2}\right) &= \int_0^\infty \frac{r^2}{(\theta r + 1)^2} \frac{\theta}{(r\theta + 1)^2} dr \\ &= \frac{1}{\theta} \int_0^\infty \frac{(\theta r)^2}{(\theta r + 1)^4} dr \\ &= \frac{1}{\theta^2} \int_1^\infty \frac{(1-u)^2}{u^4} du \\ &= \frac{1}{3\theta^2}, \end{aligned}$$

$$I_n = \frac{n}{\theta^2} - \frac{2n}{3\theta^2} = \frac{n}{3\theta^2}.$$

□

7 (2 pts)

Quelle est la distribution asymptotique de $\hat{\theta}$?

Proposez, si possible, une fonction k pour laquelle $\sqrt{n}(k(\hat{\theta}) - k(\theta)) \xrightarrow{d} N(0, 1)$. Si vous estimez que c'est impossible, veuillez expliquer pourquoi.

Solution

La théorie du maximum de vraisemblance nous dit que

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, 3\theta^2).$$

càd que $\hat{\theta} \sim_a N(\theta, 3\theta^2/n)$.

La méthode Delta nous permet de déduire que

$$\sqrt{n}(k(\hat{\theta}) - k(\theta)) \xrightarrow{d} N\left(0, 3\theta^2(k'(\theta))^2\right).$$

Il suffit de prendre $k(\theta) = \frac{1}{\sqrt{3}} \log(\theta)$ pour obtenir le résultat souhaité.

□

8 (1 pt)

Nous voulons utiliser l'algorithme de Newton-Raphson (NR) pour calculer $\hat{\theta}$. Proposer une valeur (estimateur) de départ appropriée. Justifiez votre réponse.

Décrivez le déroulement de l'algorithme NR et précisez les calculs successifs à réaliser pour parvenir à la solution.

Solution

Étant donné que $F_R(1) = 1 - (1 + \theta)^{-1}$, i.e. $\theta = F_R(1)/(1 - F_R(1))$, comme estimateur (de départ/naïf) pour θ , nous pouvons utiliser

$$\tilde{\theta} = \frac{\sum_i I(R_i \leq 1)}{\sum_i I(R_i > 1)}$$

L'algorithme NR consiste, simplement, à calculer

$$\theta_{k+1} = \theta_k + \frac{S_n(\theta_k)}{J_n(\theta_k)},$$

avec

$$S_n(\theta) = \frac{n}{\theta} - 2 \sum_i \frac{r_i}{\theta r_i + 1}$$

$$J_n(\theta) = \frac{n}{\theta^2} - 2 \sum_i \frac{r_i^2}{(\theta r_i + 1)^2}.$$

L'algorithme est arrêté lorsqu'aucun changement n'est constaté ($|\theta_{k+1} - \theta_k| \approx 0$).

□

9 (1 pt)

Nous avons relevé les 20 observations suivantes.

x	0.7	11.3	2.1	30.7	4.6	20.2	0.3	0.9	0.7	2.3	1.1	1.9	0.5	0.8	1.2	15.2	0.2	0.7	0.4	2.3
y	3.8	4.6	2.1	5.6	10.3	2.8	1.9	1.4	0.4	0.9	2.8	3.2	8.5	14.5	14.4	8.8	7.6	1.3	2.2	4.0

Ces observations ont été utilisées pour calculer $\hat{\theta}$ via l'algorithme NR. Le tableau suivant présente les itérations successives, ainsi que quelques détails relatifs aux calculs effectués (dans ces formules $r_i = x_i/y_i$).

k	θ_k	$\sum_i r_i / (\theta_k r_i + 1)$	$\sum_i r_i^2 / (\theta_k r_i + 1)^2$
0	2.333333	4.502554	1.307859
1		5.121220	1.743069
2		4.983727	1.639954
3	2.010147	4.974742	1.633345
4	2.010169	4.974705	1.633318
5	2.010169	4.974705	1.633318

Utilisez ces informations pour compléter ce tableau en calculant θ_1 et θ_2 , i.e. les valeurs de θ à l'itération 1 et 2, respectivement, de l'algorithme NR.

Solution

Selon les formules ci-dessus, la valeur θ_1 est donnée par

$$2.333333 + (20/2.333333 - 2 \times 4.502554) / (20/2.333333^2 - 2 \times 1.307859)$$

càd 1.923333.

De même la valeur θ_2 est

$$1.923333 + (20/1.923333 - 2 \times 5.12122) / (20/1.923333^2 - 2 \times 1.743069)$$

càd 2.004656.

□

8 (2 pts)

Proposer deux estimations de l'écart-type asymptotique de $\hat{\theta}_{20}$, l'EMV de θ basée sur l'échantillon observé.

Solution

L'écart-type asymptotique de $\hat{\theta}_{20}$ est donnée par $\sqrt{3\theta^2/n}$, il peut donc être estimé, dans notre cas, par

$$\sqrt{\frac{3}{20}} \times 2.010169 = 0.779$$

Une autre façon d'estimer cet écart-type est d'utiliser la "observed FI", i.e. $\sqrt{1/J_n(\theta)}$,

$$\sqrt{1/(20/2.010169^2 - 2 \times 1.633318)} = 0.771$$

□

Exercice 2 (4 pts)

Commençons par définir la loi Gamma et par donner quelques informations à son sujet.

X suit une loi Gamma de paramètres k et $\theta > 0$, i.e. $X \sim \text{Gamma}(k, \theta)$, si sa densité est donnée par

$$f_X(x) = \frac{x^{k-1} \exp(-x/\theta)}{\Gamma(k)\theta^k} \mathbf{I}(0 \leq x < \infty), \quad x \in \mathbb{R},$$

où Γ est la fonction gamma $\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} \exp(-x) dx$, $\alpha \in \mathbb{R}$. Cette fonction possède un certain nombre de propriétés, notamment le fait que $\Gamma(\alpha + 1) = \alpha\Gamma(\alpha)$, $\forall \alpha > 0$ et $\Gamma(n + 1) = n!$, $\forall n \in \mathbb{N}$. Aussi, sachez que $E(X^r) = \frac{\Gamma(k+r)}{\Gamma(k)}\theta^r$, et que $\text{Gamma}(k = n/2, \theta = 2) \equiv \chi_n^2$.

Soit Y_1, \dots, Y_n un échantillon iid de taille n ($n \geq 2$) d'une $N(0, \sigma^2)$.

1 (1 pt)

Soit $U = \sum_{i=1}^n Y_i^2$. Montrez que $U \sim \text{Gamma}(n/2, 2\sigma^2)$? Expliquez vos calculs de manière détaillée.

Solution

On a que $U = \sigma^2 Z^2$, avec $Z^2 = \sum_i (Y_i/\sigma)^2 \sim \chi_n^2 \equiv \text{Gamma}(k = n/2, \theta = 2)$. Puisque $P(U \leq u) = P(Z^2 \leq u/\sigma^2)$,

$$f_U(u) = (u/\sigma^2)' f_{Z^2}(u/\sigma^2) = \frac{1}{\sigma^2} \frac{(u/\sigma^2)^{k-1} \exp(-\frac{u}{\theta\sigma^2})}{\Gamma(k)\theta^k} = \frac{u^{k-1} \exp(-\frac{u}{\theta\sigma^2})}{\Gamma(k)(\theta\sigma^2)^k}$$

Donc $U \sim \text{Gamma}(n/2, 2\sigma^2)$.

□

2 (1 pt)

Développez un *estimateur sans biais* pour $\theta = \sigma^r$ (r est un entier positif connu), sous la forme $\hat{\theta} = aU^b$, où a et b doivent être déterminés explicitement.

Solution

On a que

$$E(\hat{\theta}) = aE(U^b) = a \frac{\Gamma(k+b)}{\Gamma(k)} (2\sigma^2)^b$$

Pour que $E(\hat{\theta}) = \sigma^r$, il suffit de choisir a et b tels que $b = r/2$ et $a = 2^{-b} \frac{\Gamma(k)}{\Gamma(k+b)}$, avec $k = n/2$.

□

3 (2 pts)

Dérivez la CRLB (Cramer-Rao Lower bound) pour $\theta = \sigma^r$. Pour $r = 2$, montrez que $\hat{\theta}$ est efficace pour θ .

Solution

Soit $g(t) = t^{r/2}$. Puisque nous savons que $I_n(\sigma^2) = \frac{n}{2\sigma^4}$ (voir cours), le plus simple est d'utiliser la relation

$$I_n(\sigma^r) = \frac{I_n(\sigma^2)}{(g'(\sigma^2))^2} = \frac{\frac{n}{2\sigma^4}}{\left(\frac{r}{2}(\sigma^2)^{(r/2-1)}\right)^2} = \frac{2n}{r^2\sigma^{2r}}.$$

Donc $CRLB(\sigma^r) = \frac{r^2\sigma^{2r}}{2n}$.

Pour $r = 2$, nous avons que $CRLB(\sigma^2) = \frac{2^2\sigma^4}{2n} = \frac{2\sigma^4}{n}$ et, dans ce cas, $\hat{\theta} = 2^{-1} \frac{\Gamma(n/2)}{\Gamma(n/2+1)} U = 2^{-1} \frac{1}{n/2} U = n^{-1} \sum Y_i^2$. $Var(\hat{\theta}) = Var(Y^2)/n$, avec $Var(Y^2) = E(Y^4) - (E(Y^2))^2 = 3\sigma^4 - \sigma^4 = 2\sigma^4$. Donc $Var(\hat{\theta}) = CRLB(\sigma^2)$.

□