

Rappels de statistique mathématique
Enoncé des travaux dirigés n°8

Guillaume Lacôte
 Bureau E03

✉ Guillaume.Lacote@ensae.fr

☞ <http://ensae.no-ip.com/SE222/>

Enoncé de l'exercice 1

Soit un échantillon X_1, \dots, X_n de variables aléatoires i.i.d telles que :

- X_i a une probabilité α de valoir a , et
- une probabilité $(1 - \alpha)$ de suivre une loi normale $\mathcal{N}(0, 1)$,

ce qui s'écrit encore :

$$X_1 = a \cdot \mathbb{1}_{Z_1=1} + Y_1 \cdot \mathbb{1}_{Z_1=0}$$

où $Z_1 \sim \mathcal{B}(1, \alpha)$ est indépendante de $Y_1 \sim \mathcal{N}(0, 1)$.

- ☞ Q1 Montrer que $(\delta_a + \lambda)^{\otimes n}$ est une mesure dominante pour le modèle considéré.
- ☞ Q2 Montrer que :

$$\frac{dP_{X_1, \dots, X_n}(x_1, \dots, x_n)}{d(\delta_a + \lambda)^{\otimes n}} = \prod_{i=1}^n (\alpha \mathbb{1}_a(x_i) + (1 - \alpha)(1 - \mathbb{1}_a(x_i))\phi(x_i))$$

où ϕ est la densité d'une loi normale $\mathcal{N}(0, 1)$.

Enoncé de l'exercice 2

Cet exercice s'inspire librement des travaux de Dov Samet, 2003.¹

On s'intéresse pour tout couple (X, Y) de variables aléatoires réelles à la propriété \mathcal{P} suivant

$$(\mathcal{P}) \quad \forall \mathcal{A} \text{ mesurable, } \begin{cases} \mathbb{P}(X < Y \text{ et } X \in \mathcal{A}) = \mathbb{P}(X > Y \text{ et } X \in \mathcal{A}) = \frac{1}{2}\mathbb{P}(X \in \mathcal{A}) \\ \mathbb{P}(Y < X \text{ et } Y \in \mathcal{A}) = \mathbb{P}(Y > X \text{ et } Y \in \mathcal{A}) = \frac{1}{2}\mathbb{P}(Y \in \mathcal{A}) \end{cases}$$

- ☞ Q1 Comment s'interprète la propriété \mathcal{P} ?

■ Partie 1 "*Devinez quel est le plus grand !*"

On s'intéresse au problème suivant, présenté originellement par Blackwell (1951) :

Deux nombres réels sont tirés aléatoirement, et chacun est placé dans une enveloppe. L'une d'elle est (indépendamment) tirée au hasard et vous est présentée. Vous devez deviner, au vu du nombre qu'elle contient, s'il s'agit du plus grand ou le plus petit des deux. Sauriez-vous deviner juste plus d'une fois sur deux en moyenne ?

¹Résultats présentés à la 14^{ème} conférence internationale de Théorie des jeux, Stony Brook 2003. Voir <http://www.sunysb.edu/gametheory/Conf03/twopuzzles.pdf> pour un exposé plus complet.

Pour simplifier la présentation, supposons que vous gagnez +1 si vous avez deviné juste et -1 sinon ; ce problème revient à savoir si vous pouvez vous garantir un gain espéré strictement positif.

☞ Q1 Soient X et Y les deux nombres tirés aléatoirement, et supposons que (X, Y) vérifie \mathcal{P} . Quel est votre meilleur gain espéré ? Interpréter.

☞ Q2 Cherchons néanmoins à construire une stratégie gagnante.

Pour ce faire, donnons-nous $s \in \mathbb{R}$ et soit $\sigma_s : \left(\begin{array}{l} \mathbb{R} \rightarrow \{-1, 1\} \\ x \mapsto \begin{cases} 1 & \text{si } x > s \\ -1 & \text{sinon} \end{cases} \end{array} \right)$ la *stratégie de seuil* associée à s .

(a) Soit $(x \neq y)$ une réalisation du tirage ; le nombre qui vous est présenté est soit x , soit y . Quel est votre gain espéré en jouant selon σ_s sachant que $x < s$ et $y < s$?

(b) Quel est votre gain espéré sachant que $x < s < y$?

(c) En déduire que si votre gain espéré (sur tous les tirages possibles) est strictement positif.

☞ Q3 En déduire qu'aucune paire de variables aléatoires réelles (X, Y) ne peut vérifier \mathcal{P} .

Partie 2 Le jeu des enveloppes

On s'intéresse désormais au célèbre jeu dit *des enveloppes*, introduit par Kraitchik (1953) sous la forme suivante :

On considère deux enveloppes contenant une certaine somme d'argent, mais l'une contenant le double de l'autre ; chacune a autant de chances que l'autre de contenir la plus grosse somme. L'une d'elles vous est donnée au hasard, et vous devez choisir entre la prendre, ou prendre l'autre.

A supposer qu'elle contienne x , l'autre a une chance sur deux de contenir $2x$ et une chance sur deux de contenir $\frac{x}{2}$; donc le contenu espéré de l'autre enveloppe est $\frac{5}{4}x > x$, de sorte que vous avez intérêt à changer d'enveloppe. Mais le même raisonnement s'applique aussi à la nouvelle enveloppe, de sorte que vous avez à nouveau intérêt à changer d'enveloppe, et ainsi de suite ...

Que faire dans cette situation ?

Résoudre ce paradoxe a été le prétexte à une littérature foisonnante ; voici une solution.

Notons X la somme (variable aléatoire réelle) contenue dans l'enveloppe qui vous est présentée, et Y celle contenue dans l'autre.

☞ Q1 Quel sens donner à l'expression \mathcal{H}_x : "l'autre a une chance sur deux de contenir $2x$ et une chance sur deux de contenir $\frac{x}{2}$ " ?

☞ Q2 En déduire que si cette hypothèse est vraie, alors la distribution de (X, Y) vérifie \mathcal{P} .

☞ Q3 Conclure.

Enoncé de l'exercice 3

Soient X_1, \dots, X_n n variables aléatoires i.i.d dans \mathbb{R}^k , de loi paramétrique \mathcal{P}_θ de densité

$$\frac{d\mathcal{P}_\theta(x)}{d\lambda} = f(x - \theta)$$

$\theta \in \mathbb{R}^k$ est appelé *paramètre de position*.

On cherche à étudier les estimateurs T du paramètre θ . On se restreint pour cela à la classe d'estimateurs *équivalariants*, c'est-à-dire des estimateurs T vérifiant :

$$\forall c \in \mathbb{R}^k, T(X_1 + c, \dots, X_n + c) = T(X_1, \dots, X_n) + c$$

Cette restriction est naturelle : en effet, pour tout $c \in \mathbb{R}^k$, $X_i + c \sim \mathcal{P}_{\theta+c}$, et on ne peut admettre qu'un simple changement d'échelle puisse mener à une estimation différente.

☞ Q1 Montrer que :

$$T \text{ équivalent} \Leftrightarrow \begin{array}{l} \exists T_1 : (\mathbb{R}^{k-1} \rightarrow \mathbb{R}^k) / \forall (x_1, \dots, x_n) \in \mathbb{R}^k \\ T(x_1, \dots, x_n) = T_1(x_2 - x_1, \dots, x_n - x_1) + x_1 \end{array}$$

☞ Q2 Soit $W(x, y) = w(x - y)$ une fonction de perte. Montrer que, si un estimateur T est équivalent au risque associé, pour la perte W , ne dépend pas de θ .

Un estimateur équivalent T qui minimise le risque $R_w(T, \theta)$ (c'est-à-dire au point $\theta = 0$) alors un estimateur optimal (parmi les équivalariants) pour la fonction de perte W .

☞ Q3 Soit la fonction

$$\psi(x) = \int_{\mathbb{R}^k} w(x - u) \prod_{i=1}^n f(X_i - u) du$$

et l'estimateur T^* défini par

$$\psi(T^*) = \inf_{x \in \mathbb{R}^k} \psi(x)$$

T^* est appelé *l'estimateur de Pitman* (on suppose ici que w est telle que l'équation $\psi(x^*) = \inf_{x \in \mathbb{R}^k} \psi(x)$ admette une solution et une seule).

(a) Montrer que T^* est équivalent.

(b) Montrer que, pour toute fonction φ telle que $\mathbb{E}(\|\varphi(X_1, \dots, X_n)\|) < +\infty$, on a :

$$E_\theta[\varphi(X_1, \dots, X_n) | X_2 - X_1, \dots, X_n - X_1] = \int_{\mathbb{R}^k} \varphi(X_1 + \theta - u, \dots, X_n + \theta - u) \frac{\prod_{i=1}^n f(X_i - u)}{\int_{\mathbb{R}^k} \prod_{i=1}^n f(X_i - v) dv} du$$

(c) Montrer finalement que T^* est optimal pour w .

☞ Q4 Donner l'expression de T^* lorsque $w(x - y) = \|x - y\|^2$.

☞ Q5 Dans ce cas, que vaut T^* lorsque $X_i \sim \mathcal{N}(\theta, 1)$?

Et lorsque $X_i \sim \mathcal{U}_{[-\frac{1}{2}+\theta, \frac{1}{2}+\theta]}$?

Enoncé de l'exercice 4

On considère une loi multinomiale à K modalités ($K \geq 3$) de probabilités p_1, p_2, \dots, p_K . On dispose de n observations indépendantes issues de cette loi. On notera N_k le nombre d'observations de la modalité $k \in \llbracket 1, K \rrbracket$.

On veut tester l'hypothèse : $H_0 : "p_1 + p_2 = \frac{1}{2}"$

☞ Q1 Donner les estimateurs du maximum de vraisemblance contraints $(\hat{p}_k^0)_k$ et non contraints $(\hat{p}_k)_k$ de p_1, \dots, p_K .

☞ Q2 Calculer la statistique ξ^W du test de Wald de l'hypothèse H_0 .

☞ Q3 (a) En constatant que \hat{p}_3^0 est asymptotiquement efficace sous H_0 , montrer que

$$\text{Cov}_{as}(\hat{p}_3, \hat{p}_3 - \hat{p}_3^0) = 0$$

et en déduire que

$$\mathbb{V}_{as}(\hat{p}_3 - \hat{p}_3^0) = \mathbb{V}_{as}(\hat{p}_3) - \mathbb{V}_{as}(\hat{p}_3^0)$$

(b) Calculer la loi limite de $(\hat{p}_3^0 - p_3)$, et en déduire $\mathbb{V}_{as}(\hat{p}_3^0)$.

(c) Donner un estimateur $\widehat{\mathbb{V}_{as}(\hat{p}_3 - \hat{p}_3^0)}$ convergent sous H_0 de $\mathbb{V}_{as}(\hat{p}_3 - \hat{p}_3^0)$.
En déduire l'expression de la statistique du test d'Hausman de H_0

$$\xi^H = n \frac{(\hat{p}_3 - \hat{p}_3^0)^2}{\widehat{\mathbb{V}_{as}(\hat{p}_3 - \hat{p}_3^0)}}$$

(d) Vérifier que le test d'Hausman est convergent.

(e) Montrer que les statistiques ξ^W et ξ^H sont asymptotiquement équivalentes sous H_0 .

Enoncé de l'exercice 5

La mise en œuvre de nombreuses techniques d'analyse, et notamment la quasi-totalité des techniques d'estimation et de test de modèles statistiques, nécessitent d'engendrer des réalisations x d'une variable aléatoire réelle X de loi \mathcal{L} donnée.

L'objet de cet exercice est de présenter quelques techniques permettant à un ordinateur, machine essentiellement déterministe, d'engendrer de telles variables. Il s'agit donc de construire pour toute loi \mathcal{L} une fonction récursive $f^{\mathcal{L}} : \mathbb{N} \rightarrow I$ qui énumère des réalisations d'un variable $X \sim \mathcal{L}$, i.e. telle que la distribution empirique des $(f^{\mathcal{L}}(n))_{n \in \mathbb{N}}$ converge vers la vraie distribution $F^{\mathcal{L}}$.

☞ Q1 Cherchons tout d'abord à engendrer une variable aléatoire **entière** uniforme.

Soit $N \in \mathbb{N}^*$ et considérons $\mathcal{L} = \mathcal{U}_{\llbracket 0, N-1 \rrbracket}$.

(a) Définissons pour $\phi : \llbracket 0, N-1 \rrbracket \rightarrow \llbracket 0, N-1 \rrbracket$ et $s \in \llbracket 0, N-1 \rrbracket$

$$f_{\phi, s} : \left(\begin{array}{cc} \llbracket 0, N-1 \rrbracket & \rightarrow \llbracket 0, N-1 \rrbracket \\ t & \mapsto \phi^t(s) \end{array} \right)$$

Montrer que $f_{\phi, s}$ est périodique à partir d'un certain rang ; on note $T_{\phi, s}$ sa période.

En quel sens peut-on dire que $f_{\phi, s}$ "génère" une variable uniformément distribuée $\llbracket 0, N-1 \rrbracket$?

(b) En quel sens peut-on dire que $f_{\phi, s}$ est d'autant meilleure que $T_{\phi, s}$ est grande ?

(c) Soient $a \in \llbracket 2, N-1 \rrbracket$, $c \in \llbracket 0, N-1 \rrbracket$ et $p \in \llbracket 2, N-1 \rrbracket$ et définissons $\psi_{a, c, p}$

$$\left(\begin{array}{cc} \mathbb{N} & \rightarrow \llbracket 0, N-1 \rrbracket \\ n & \mapsto (an + c) \bmod p \end{array} \right).$$

A quelles conditions $f_{\psi_{a, c, p}, s}$ est-elle un bon générateur uniforme sur $\llbracket 0, N-1 \rrbracket$?

☞ Q2 Cherchons à engendrer une variable aléatoire réelle uniforme sur $[0, 1]$.

(a) Soit $(x_n)_{n \in \mathbb{N}} \in [0, 1]^{\mathbb{N}}$, et définissons pour $f : [0, 1] \rightarrow \mathbb{R}$ continue et $n \in \mathbb{N}^*$

$$S_n^x(f) = \frac{1}{n} \sum_{k=1}^n f(x_k)$$

On dit de x qu'elle vérifie le *critère de Weyl* si pour toute f continue

$$S_n^x(f) \xrightarrow{n \rightarrow +\infty} \int_0^1 f$$

Comment s'interprète ce critère ?

(b) Soit $r \in \mathbb{R}$, et définissons $x = (\text{frac}(rk))_{k \in \mathbb{N}}$ où $\text{frac}(u) = u - E(u) \in [0, 1[$ désigne la partie fractionnaire de $u \in \mathbb{R}$.

Montrer que x vérifie le critère de Weyl ssi r est irrationnel.

☞ Q3 On cherche désormais à engendrer des réalisations d'une variable aléatoire de loi quelconque

(a) Soit \mathcal{L} la loi de densité $f_{\mathcal{L}} = \frac{1}{3}\delta_1 + \frac{1}{4}\delta_2 + \frac{5}{12}\delta_3$.

Tracer la fonction de répartition $F_{\mathcal{L}}$ de \mathcal{L} .

En déduire son inverse généralisée

$$F^{-1} : \left(\begin{array}{cc} [0, 1] & \rightarrow \mathbb{R} \\ s & \mapsto \inf\{x \in \mathbb{R} / F_{\mathcal{L}}(x) \geq s\} \end{array} \right)$$

Soit $U \sim \mathcal{U}_{[0,1]}$ et définissons $X = F_{\mathcal{L}}^{-1}(U)$; quelle est la loi de X ?

(b) Soit alors $F_{\mathcal{L}}$ la fonction de répartition supposée continue et strictement croissante d'une loi \mathcal{L} quelconque.

On définit de même $X = F_{\mathcal{L}}^{-1}(U)$ pour $U \sim \mathcal{U}_{[0,1]}$; quelle est la loi de X ?

(c) Proposer un algorithme qui génère des réalisations successives d'une variable X de loi exponentielle $\mathcal{E}(\lambda)$ de densité $f_{\mathcal{E}(\lambda)}(x) = \lambda e^{-\lambda x} \mathbf{1}_{[0,+\infty[}(x)$ pour $\lambda > 0$.

(d) Soit $(X, Y) \sim \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \sigma^2 I_2\right)$.

Soient $R^2 = X^2 + Y^2$ et $\theta = \arctan \frac{Y}{X}$.

Montrer que R^2 et θ sont indépendantes, et donner leurs lois.

En déduire un algorithme de génération d'une variable gaussienne.

(e) Même question si X suit une loi de Weibull $\mathcal{W}(\alpha, \beta)$ de densité $f_{\mathcal{W}(\alpha, \beta)}(x) = \alpha \beta x^{\beta-1} e^{-\alpha x^\beta} \mathbf{1}_{x \geq 0}$.

(f) Même question si X suit une loi de Poisson $\mathcal{P}(\lambda)$ de densité $f_{\mathcal{P}(\lambda)}(k) = \frac{\lambda^k}{k!} e^{-\lambda}$.

On pourra comparer $x_n = u_1 \times \dots \times u_n$, où chaque u_n est une réalisation de $U \sim \mathcal{U}_{[0,1]}$, à $e^{-\lambda}$.