

Chapitre III

Méthodes générales d'estimation

1 Quelques considérations préliminaires

Nous allons, dans ce chapitre, supposer que nous travaillons sur le *modèle d'échantillonnage*, c'est à dire un modèle stochastique dans lequel on observe une suite de v.a. i.i.d. X_1, \dots, X_n, \dots à valeurs dans (E, \mathcal{E}) et de loi inconnue P_X . Nous nous intéressons ici à un paramètre inconnu μ fonction de P , par exemple $\mu = \mathbb{E}[X_1]$ et à des suites d'estimateurs $\hat{\mu}_n(X_1, \dots, X_n)$ de μ . Pour décrire les propriétés asymptotiques d'une telle suite, on introduit les notions suivantes.

Définition 9 Soit X_1, \dots, X_n, \dots une suite de variables aléatoires i.i.d. de loi P_X et μ un paramètre inconnu fonction de P_X . Soit $\hat{\mu}_n(X_1, \dots, X_n)$ un estimateur de μ fonction de X_1, \dots, X_n . La suite d'estimateurs $(\hat{\mu}_n)_{n \geq 1}$ est dite convergente (consistante) si

$$\hat{\mu}_n \xrightarrow[n \rightarrow +\infty]{P} \mu.$$

Si $\mu \in \mathbb{R}$, elle est dite asymptotiquement normale s'il existe une quantité positive σ^2 , telle que

$$\sqrt{n}(\hat{\mu}_n - \mu) \rightsquigarrow \mathcal{N}(0, \sigma^2).$$

La normalité asymptotique de la suite $(\hat{\mu}_n)_{n \geq 1}$ entraîne sa convergence.

Remarque importante : Lorsque l'on considère des suites d'estimateurs asymptotiquement normaux, ils sont d'autant meilleurs (asymptotiquement) que leur variance asymptotique est plus petite. En effet, si $\sqrt{n}(\hat{\mu}_n - \mu) \rightsquigarrow \mathcal{N}(0, \sigma^2)$, si Z est $\mathcal{N}(0, 1)$ et $\mathbb{P}[|Z| \leq t_\alpha] = 1 - \alpha$, alors

$$\mathbb{P}\left[|\hat{\mu}_n - \mu| \leq n^{-1/2}t_\alpha\sigma\right] \xrightarrow[n \rightarrow +\infty]{} 1 - \alpha,$$

de sorte que la variance asymptotique σ^2 donne la précision de l'approximation de μ par $\hat{\mu}_n$ lorsque n est suffisamment grand. Si σ est indépendant de μ et connu, on en déduit immédiatement des intervalles de confiance asymptotiques pour μ :

$$\mathbb{P}\left[\hat{\mu}_n - n^{-1/2}t_\alpha\sigma \leq \mu \leq \hat{\mu}_n + n^{-1/2}t_\alpha\sigma\right] \xrightarrow[n \rightarrow +\infty]{} 1 - \alpha.$$

Si σ est inconnu ou dépend de μ , il faut l'estimer. Si l'on dispose d'une suite d'estimateurs convergents $\hat{\sigma}_n^2$ de σ^2 , le Théorème de Slutsky entraîne que $\sqrt{n}\hat{\sigma}_n^{-1}(\hat{\mu}_n - \mu) \rightsquigarrow \mathcal{N}(0, 1)$, de sorte que

$$\mathbb{P}\left[\hat{\mu}_n - n^{-1/2}t_\alpha\hat{\sigma}_n \leq \mu \leq \hat{\mu}_n + n^{-1/2}t_\alpha\hat{\sigma}_n\right] \xrightarrow[n \rightarrow +\infty]{} 1 - \alpha. \quad (\text{III.1.1})$$

On obtient évidemment des résultats analogues pour les intervalles de confiance unilatères. Lorsque $\sigma = \sigma(\mu)$ est une fonction continue et connue de μ , on peut l'estimer par $\sigma(\hat{\mu}_n)$. Attention, les intervalles de confiance (III.1.1) peuvent être vraiment très asymptotiques !

2 Loi empirique et estimateurs empiriques

Nous ferons ici sur P_X des hypothèses minimales, par exemple que P_X a des moments jusqu'à l'ordre 2, ou 4, ..., ou que sa fonction de répartition est continue, dérivable, etc.

2.1 Loi empirique, moments et espérances empiriques

2.1.1 Premier exemple, la moyenne empirique

L'exemple le plus classique d'estimateur *empirique* est la *moyenne empirique* des observations, $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$. Si l'on suppose que P_X a un moment d'ordre 2, c'est à dire que les X_i ont une variance σ^2 , on obtient immédiatement que

$$\bar{X}_n \xrightarrow[n \rightarrow +\infty]{P} \mathbb{E}[X_1] \quad \text{et} \quad \sqrt{n} (\bar{X}_n - \mathbb{E}[X_1]) \rightsquigarrow \mathcal{N}(0, \sigma^2).$$

Donc la suite des estimateurs \bar{X}_n de l'espérance des X_i est convergente et asymptotiquement normale.

2.1.2 Loi empirique

Si l'on parle de moyenne empirique pour \bar{X}_n , c'est parce que \bar{X}_n est effectivement la moyenne (c'est-à-dire l'espérance) d'une certaine loi de probabilité, dite *loi empirique*, associée aux observations.

Définition 10 *Étant donné des points x_1, \dots, x_n dans un ensemble (E, \mathcal{E}) , on appelle loi empirique des x_i la probabilité $\nu_n = n^{-1} \sum_{i=1}^n \delta_{x_i}$ sur E (où δ_x désigne la mesure de Dirac au point x), c'est à dire la probabilité définie par*

$$\nu_n(A) = n^{-1} \sum_{i=1}^n \mathbb{1}_A(x_i) \quad \text{pour tout } A \in \mathcal{E}. \quad (\text{III.2.1})$$

Les quantités associées à la loi empirique sont également dites empiriques et l'on parlera ainsi de moyenne empirique des x_i , de variance empirique des x_i , de fonction de répartition empirique des x_i , etc.

Lorsque les x_i sont des réalisations de variables aléatoires X_i i.i.d., c'est à dire que $x_i = X_i(\omega)$ pour tout i , la loi empirique de l'échantillon X_1, \dots, X_n est une fonction de $\omega \in \Omega$ qui s'écrit

$$\nu_n(\omega) = n^{-1} \sum_{i=1}^n \delta_{X_i(\omega)}.$$

Dans ce cas, ν_n est une probabilité aléatoire, c'est à dire une application (mesurable) de l'ensemble Ω dans l'ensemble des probabilités à support fini sur (E, \mathcal{E}) , qui est une fonction des X_i , de même que toutes les quantités empiriques associées, en particulier la fréquence empirique de l'évènement A : $\nu_n(A) = n^{-1} \sum_{i=1}^n \mathbb{1}_A(X_i)$.

On notera que si les observations sont toutes distinctes, ν_n est simplement la loi uniforme sur l'ensemble $\{X_1, \dots, X_n\}$. Comme les X_i sont i.i.d. de loi P_X , $n\nu_n(A) \sim \mathcal{B}(n, P_X(A))$.

2.1.3 Espérances empiriques

Si la variable aléatoire Y a la loi ν_n définie par (III.2.1), pour toute fonction h , $\mathbb{E}_{\nu_n}[h(Y)] = n^{-1} \sum_{i=1}^n h(x_i)$. Dans le cas de la loi empirique d'un échantillon X_1, \dots, X_n , ceci s'écrit $\mathbb{E}_{\nu_n}[h(Y)] = n^{-1} \sum_{i=1}^n h(X_i)$ et comme ν_n est aléatoire, les espérances par rapport à ν_n le sont aussi (ce sont des fonctions de ω). En particulier, la moyenne empirique \bar{X}_n des X_i définie plus haut est tout simplement l'espérance (moyenne) $\mathbb{E}_{\nu_n}[Y]$ de la loi empirique ν_n . La loi des grands nombres et le T.L.C. impliquent que si les observations X_i sont i.i.d. de loi P_X ,

- i) si h est P -intégrable, alors $\mathbb{E}_{\nu_n}[h(Y)] \xrightarrow[n \rightarrow +\infty]{P} \mathbb{E}_{P_X}[h(Y)]$;
- ii) si $\mathbb{E}_{P_X}[h^2(Y)] < +\infty$, $\sqrt{n} (\mathbb{E}_{\nu_n}[h(Y)] - \mathbb{E}_{P_X}[h(Y)]) \rightsquigarrow \mathcal{N}(0, \text{Var}_{P_X}(h(Y)))$.

2.1.4 Variance empirique

On a souvent besoin d'estimer la variance σ^2 de la loi de X , alors même que la moyenne est inconnue, de sorte qu'il ne suffit pas d'estimer le second moment $\mathbb{E}[X^2]$ pour y parvenir. On utilise alors la variance empirique $\text{Var}_{\nu_n}(Y)$, c'est à dire la variance $\hat{\sigma}_n^2$ de la loi empirique qui s'écrit

$$\hat{\sigma}_n^2 = n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \mathbb{E}_{\nu_n}[Y^2] - (\mathbb{E}_{\nu_n}[Y])^2 = n^{-1} \sum_{i=1}^n X_i^2 - \bar{X}_n^2. \quad (\text{III.2.2})$$

L'estimateur obtenu est biaisé d'après l'inégalité de Jensen $\mathbb{E}[\bar{X}_n^2] > (\mathbb{E}[\bar{X}_n])^2$ (sauf si X_i est presque sûrement constante). On a donc $\mathbb{E}[\hat{\sigma}_n^2] < \sigma^2$. Plus précisément, en développant le carré, on trouve que $\mathbb{E}[\hat{\sigma}_n^2] = n^{-1}(n-1)\sigma^2$. Il s'ensuit que l'estimateur modifié $\hat{s}_n^2 = (n-1)^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ est un estimateur sans biais de σ^2 . Asymptotiquement, les deux estimateurs, $\hat{\sigma}_n^2$ et \hat{s}_n^2 sont équivalents.

Pour calculer la loi limite de l'un ou de l'autre de ces estimateurs, on ne peut utiliser le T.L.C. car les composantes $(X_i - \bar{X}_n)^2$ de la somme ne sont pas indépendantes. Il faut avoir recours à une astuce. On remarque que le T.L.C. s'applique aux variables $t_n^2 = n^{-1} \sum_{i=1}^n (X_i - \mathbb{E}[X_i])^2$ et donne $\sqrt{n}(t_n^2 - \sigma^2) \rightsquigarrow \mathcal{N}(0, v^2)$ avec $v^2 = \text{Var}((X_i - \mathbb{E}[X_i])^2)$. On montre ensuite par un calcul direct de la différence $\hat{\sigma}_n^2 - t_n^2$ que $\sqrt{n}(\hat{\sigma}_n^2 - t_n^2) \rightarrow 0$. Le théorème de Slutsky permet alors de conclure que

$$\sqrt{n}(\hat{\sigma}_n^2 - \sigma^2) = \sqrt{n}(\hat{\sigma}_n^2 - t_n^2) + \sqrt{n}(t_n^2 - \sigma^2) \rightsquigarrow \mathcal{N}(0, v^2).$$

Pour traiter le cas de \hat{s}_n^2 il suffit de remarquer à nouveau que $\sqrt{n}(\hat{\sigma}_n^2 - \hat{s}_n^2) \rightarrow 0$.

2.2 Fonction de répartition et quantiles empiriques

2.2.1 Fonction de répartition empirique

On rappelle que si X est une v.a.r. de loi P_X , sa fonction de répartition F_X est donnée par

$$F_X(t) = \mathbb{P}[X \leq t] = P_X[[-\infty; t]]$$

et que F_X est une fonction croissante (pas nécessairement strictement), de 0 (quand $t \rightarrow -\infty$) à 1 (quand $t \rightarrow +\infty$), continue à droite et ayant des limites à gauche. Elle a au plus un nombre dénombrable de sauts d'amplitude nécessairement > 0 . En particulier la

loi empirique $\nu_n = n^{-1} \sum_{i=1}^n \delta_{X_i}$ a une fonction de répartition que l'on note F_n et que l'on appelle *fonction de répartition empirique des X_i* :

$$F_n(t) = \nu_n([-\infty; t]) = n^{-1} \sum_{i=1}^n \mathbb{1}_{[-\infty; t]}(X_i). \quad (\text{III.2.3})$$

Pour la calculer, le plus simple est d'ordonner en croissant l'ensemble des X_i (avec possibilité d'ex-aequo).

Définition 11 *Étant donné un n -échantillon X_1, \dots, X_n , on appelle statistique d'ordre associée, l'ensemble des valeurs prises par les X_i (en conservant les valeurs multiples, s'il y en a, avec leur multiplicité) ordonnées en croissant de la plus petite à la plus grande. On note l'échantillon ordonné $X_{(1)}, \dots, X_{(n)}$ avec $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$.*

On peut alors vérifier, en notant $X_{(0)} = -\infty$ et $X_{(n+1)} = +\infty$, que

$$F_n(t) = i/n \quad \text{si } t \in [X_{(i)}; X_{(i+1)}[, \quad 0 \leq i \leq n.$$

Une conséquence immédiate de (III.2.3) est que $nF_n(t) \sim \mathcal{B}(n, F_X(t))$, ce qui entraîne que $\lim_{n \rightarrow +\infty} F_n(t) = F_X(t)$ quel que soit t . Il n'est pas très difficile de vérifier que, puisque F_X et F_n sont croissantes, cette convergence est en fait uniforme en t , ce qui donne

Théorème 4 (Glivenko-Cantelli) *Si l'on dispose de n v.a.r. i.i.d. X_1, \dots, X_n de fonction de répartition F_X et si F_n désigne la fonction de répartition empirique,*

$$\sup_{t \in \mathbb{R}} |F_n(t) - F_X(t)| \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} 0.$$

De plus, pour t fixé, le T.L.C. nous dit que

$$\sqrt{n} (F_n(t) - F_X(t)) \rightsquigarrow \mathcal{N}(0; F_X(t)[1 - F_X(t)]),$$

et l'inégalité de Hoeffding pour les binomiales que

$$\mathbb{P} [|F_n(t) - F_X(t)| \geq x] \leq 2 \exp [-2nx^2],$$

ce qui suggère que la différence $|F_n(t) - F_X(t)|$ est d'ordre de grandeur $n^{-1/2}$, même pour n petit. On peut en fait démontrer (mais c'est très difficile) un résultat beaucoup plus fin et non-asymptotique très utile en pratique :

Théorème 5 (Massart) *Supposons que l'on dispose de n v.a.r. i.i.d. X_1, \dots, X_n de fonction de répartition F_X et notons F_n la fonction de répartition empirique. Alors, pour tout $n \geq 1$ et tout $x > 0$,*

$$\mathbb{P} \left[\sup_{t \in \mathbb{R}} |F_n(t) - F_X(t)| \geq x \right] \leq 2 \exp [-2nx^2]. \quad (\text{III.2.4})$$

2.2.2 Quantiles et quantiles empiriques

Si l'on considère un n -échantillon de loi $\mathcal{N}(\mu; 1)$ où μ désigne un paramètre inconnu, on peut caractériser μ par le fait qu'il est l'espérance de la loi et l'estimer par la moyenne empirique. Si l'on remplace ici la densité gaussienne standard par la densité de Cauchy, $f(x) = [\pi(1+x^2)]^{-1}$, qui est elle aussi symétrique, on ne peut plus caractériser son centre de symétrie comme étant l'espérance de la loi car cette espérance n'existe pas. Donc pour estimer le paramètre de symétrie μ dans le modèle $\{P_\mu, \mu \in \mathbb{R}\}$ avec P_μ de densité $[\pi(1+(x-\mu)^2)]^{-1}$, on ne peut pas utiliser la moyenne empirique. Il convient donc de caractériser μ autrement. On peut, par exemple, dire que μ est la *médiane* de la loi, c'est-à-dire le point $x_{1/2}$ tel que $F_\mu(x_{1/2}) = 1/2$, où F_μ est la fonction de répartition de P_μ . Plus généralement, on peut considérer le point x_p tel que $F_\mu(x_p) = p \in]0; 1[$. La définition de x_p est évidemment très simple lorsque la fonction F est continue et strictement croissante parce qu'elle est alors inversible : quel que soit $p \in]0; 1[$, il existe x_p unique tel que $F(x_p) = p$. La situation devient plus complexe lorsque F a des discontinuités ou n'est pas strictement croissante (**dessins**). Pour définir des analogues de x_p dans cette situation, on va utiliser la notion d'*inverse généralisée* de F .

Définition 12 *Étant donné une fonction de répartition F sur \mathbb{R} , on appelle inverse généralisée de F et l'on note F^{-1} la fonction définie sur $[0; 1]$ par*

$$F^{-1}(x) = \inf\{t \in \mathbb{R} \mid F(t) \geq x\}, \quad x \in [0; 1],$$

avec les conventions habituelles : $\inf \mathbb{R} = -\infty$ et $\inf \emptyset = +\infty$.

Les propriétés suivantes de F^{-1} sont immédiates : $F^{-1}(0) = -\infty$, F^{-1} est croissante (au sens large) et lorsque F est continue et strictement croissante, F^{-1} est la fonction réciproque de F sur $]0; 1[$ au sens usuel. De plus, la continuité à droite de F implique que $F \circ F^{-1}(x) \geq x$. Par définition,

$$F(t) \geq x \quad \text{si} \quad t \geq F^{-1}(x) \quad \text{et} \quad F(t) < x \quad \text{si} \quad t < F^{-1}(x),$$

de sorte que

$$F(t) \geq x \quad \text{équivaut à} \quad t \geq F^{-1}(x). \quad (\text{III.2.5})$$

Attention : Si F n'est pas bijective, il existe $a < b$ tels que $F(b) = F(a)$, donc $F^{-1} \circ F(b) \leq a < b$; si F n'est pas continue, il existe a tel que $F(a_-) < F(a)$, donc $F \circ F^{-1}(t) = F(a) > t$ pour tout $t \in]F(a_-); F(a)[$. Il y a donc équivalence entre le fait que $F \circ F^{-1} = F^{-1} \circ F = Id$ et le fait que F soit inversible au sens usuel.

On peut alors définir les quantiles de n'importe quelle fonction de répartition de la manière suivante :

Définition 13 *Étant donné un nombre $p \in [0; 1]$ et une fonction de répartition F sur \mathbb{R} , le p -quantile $x_p(F)$ est donné par $x_p(F) = F^{-1}(p)$. Les quantiles sont aussi appelés pourcentiles. La médiane correspond à $p = 1/2$ et les quartiles à $p = 1/4$ et $p = 3/4$.*

On a évidemment toujours $x_0(F) = -\infty$. On peut alors définir les *quantiles empiriques* $x_p(n)$ comme les quantiles de la fonction de répartition empirique, soit

$$x_p(n) := x_p(F_n) = \inf\{t \in \mathbb{R} \mid F_n(t) \geq p\} \quad \text{pour } 0 \leq p \leq 1.$$

On vérifie que

$$x_p(n) = X_{(i)} \quad \text{si} \quad \frac{i-1}{n} < p \leq \frac{i}{n} \quad \text{pour } 1 \leq i \leq n,$$

ce qui permet d'évaluer $x_p(n)$ pour tous les $p \in]0; 1]$. De manière équivalente,

$$x_p(n) = X_{(i)} \quad \text{pour} \quad np \leq i < np + 1. \quad (\text{III.2.6})$$

Dans la mesure où la fonction de répartition empirique converge vers la vraie fonction de répartition, il est facile de prédire que, si la fonction de répartition se comporte bien au voisinage d'un quantile, le p -quantile empirique va converger vers le p -quantile correspondant. Un résultat dans cette direction est le suivant.

Théorème 6 *Soit X_1, \dots, X_n un n -échantillon de f.r. inconnue F et p un réel dans $]0; 1[$. On suppose qu'il existe une solution x_p à l'équation $F(x) = p$ et que F est continue et dérivable au point x_p de dérivée $f(x_p) > 0$. Alors le quantile empirique $x_p(n)$ converge en probabilité vers x_p quand $n \rightarrow +\infty$ et*

$$\sqrt{n} [x_p(n) - x_p] \rightsquigarrow \mathcal{N} \left(0, \frac{p(1-p)}{f^2(x_p)} \right) \quad \text{quand } n \rightarrow +\infty.$$

Attention : Ce résultat n'est pas valable si F a une tangente verticale au point x_p (ce qui correspond à $f(x_p) = +\infty$) auquel cas la convergence de $x_p(n)$ vers x_p est plus rapide. Si $f(x_p) = 0$ mais F est croissante au voisinage de x_p , la convergence est plus lente. S'il existe $a < b$ tels que $F(b) = F(a) = p$, alors $x_n(p)$ ne converge même pas vers x_p .

3 Applications de la méthode delta

Supposons que l'on sache estimer correctement une certaine fonction strictement monotone $\phi(\theta)$ du paramètre, par exemple un moment ou un quantile de P_θ , par la méthode empirique vue précédemment. Si l'on dispose d'une suite d'estimateurs $\hat{\phi}_n$ convergents et asymptotiquement normaux de $\phi(\theta)$, et si ϕ^{-1} est une fonction dérivable, la méthode delta consiste à estimer θ par $\hat{\theta}_n = \phi^{-1}(\hat{\phi}_n)$ et le Théorème 1 montre que la suite $(\hat{\theta}_n)_{n \geq 1}$ est asymptotiquement normale. On obtient en particulier le corollaire suivant.

Corollaire 1 *Soit une suite de variables aléatoires $(X_n)_{n \geq 1}$ i.i.d. de loi P_θ avec $\theta \in \Theta$, où Θ est un intervalle ouvert de \mathbb{R} et soit ϕ une bijection de Θ dans $\phi(\Theta)$ telle que ϕ soit continuellement dérivable, de dérivée $\phi' \neq 0$. Si, de plus, $\sqrt{n} [\hat{\phi}_n - \phi(\theta)] \rightsquigarrow \mathcal{N}(0, \sigma^2(\theta))$, alors*

$$\sqrt{n} [\phi^{-1}(\hat{\phi}_n) - \theta] \rightsquigarrow \mathcal{N}(0, [\sigma(\theta)/\phi'(\theta)]^2).$$

Démonstration : Le Théorème 1 appliqué avec $g = \phi^{-1}$ montre que

$$\sqrt{n} [\phi^{-1}(\hat{\phi}_n(X_1, \dots, X_n)) - \theta] \rightsquigarrow (\phi^{-1})'(\phi(\theta)) \mathcal{N}(0, \sigma^2(\theta))$$

et $(\phi^{-1})'(\phi(\theta)) = 1/\phi'(\theta)$. \square

De même, si l'on dispose d'intervalles de confiance $[\hat{\phi}_n - a_n; \hat{\phi}_n + b_n]$ pour $\phi(\theta)$ de niveau $1 - \alpha$, on en déduit (cas ϕ croissante) des intervalles de confiance de même niveau pour θ puisque

$$\mathbb{P}_\theta [\phi^{-1}(\hat{\phi}_n - a_n) \leq \theta \leq \phi^{-1}(\hat{\phi}_n + b_n)] \geq 1 - \alpha.$$

3.1 La méthode des moments

Elle est ainsi appelée parce que l'estimateur initial de $\phi(\theta)$ est un moment de P_θ de la forme $\mathbb{E}_\theta [X^k]$ ou, plus généralement $\mathbb{E}_\theta[h(X)]$. On peut en particulier utiliser la moyenne empirique \bar{X}_n comme estimateur initial.

Exemple 1 : le modèle uniforme La loi uniforme sur un intervalle $[a; b]$ de \mathbb{R} , notée $\mathcal{U}([a; b])$, est la loi qui modélise un tirage “au hasard” dans $[a; b]$, c'est à dire que la probabilité d'un sous-intervalle de $[a; b]$ est proportionnelle à sa longueur. La densité correspondante s'écrit $(b - a)^{-1} \mathbb{1}_{[a; b]}(x)$ et si X a la loi $\mathcal{U}([a; b])$, alors $\mathbb{E}[X] = (a + b)/2$ et $\text{Var}(X) = (b - a)^2/12$.

De cette famille de lois à deux paramètres, on peut déduire de nombreux modèles à un paramètre tels que $\{\mathcal{U}([0; \theta]), \theta > 0\}$, $\{\mathcal{U}([-\theta; \theta]), \theta > 0\}$, $\{\mathcal{U}([\theta; \theta + 1]), \theta \in \mathbb{R}\}$ ou $\{\mathcal{U}([\theta - 1; \theta + 1]), \theta \in \mathbb{R}\}$. À titre d'illustration analysons le premier modèle pour lequel $\mathbb{E}_\theta[X] = \theta/2$ et $\text{Var}_\theta(X) = \theta^2/12$. On peut utiliser $2\bar{X}_n$ comme estimateur (sans biais) de θ et son risque quadratique est $\mathbb{E}_\theta[(2\bar{X}_n - \theta)^2] = (3n)^{-1}\theta^2$. De plus $\sqrt{n}(2\bar{X}_n - \theta) \rightsquigarrow \mathcal{N}(0; \theta^2/3)$. On a donc

$$\sqrt{3n}(2\theta^{-1}\bar{X}_n - 1) \rightsquigarrow \mathcal{N}(0; 1),$$

ce qui permet de construire des intervalles de confiance asymptotiques pour θ . En effet si Z est $\mathcal{N}(0; 1)$ et $\mathbb{P}[|Z| \leq z_\alpha] = 1 - \alpha$, alors

$$\mathbb{P}\left[\sqrt{3n}|2\theta^{-1}\bar{X}_n - 1| \leq z_\alpha\right] = \mathbb{P}\left[\frac{2\bar{X}_n}{1 + z_\alpha/\sqrt{3n}} \leq \theta \leq \frac{2\bar{X}_n}{1 - z_\alpha/\sqrt{3n}}\right] \xrightarrow{n \rightarrow +\infty} 1 - \alpha.$$

Exemple 2 : les lois gamma On rappelle que, pour tout $t > 0$, la fonction $x \mapsto x^{t-1}e^{-x}$ est intégrable sur \mathbb{R}_+ et son intégrale $\Gamma(t) = \int_0^{+\infty} x^{t-1}e^{-x} dx$ définit la fonction gamma sur $]0; +\infty[$. On a en particulier

$$\Gamma(1) = 1; \quad \Gamma(t + 1) = t\Gamma(t); \quad \Gamma(n + 1) = n! \text{ si } n \in \mathbb{N}; \quad \Gamma(1/2) = \sqrt{\pi}.$$

Il s'ensuit que, pour $t > 0, \lambda > 0$, la quantité $[\Gamma(t)]^{-1}\lambda^t x^{t-1}e^{-\lambda x} \mathbb{1}_{\mathbb{R}_+}(x)$ est une densité de probabilité par rapport à la mesure de Lebesgue. La loi associée est dite *loi gamma de paramètres t et λ* et notée $\Gamma(t, \lambda)$. Si X est une variable aléatoire de loi $\Gamma(t, \lambda)$ alors $\mathbb{E}[X] = t/\lambda$, $\text{Var}(X) = t/\lambda^2$ et $\theta X \sim \Gamma(t, \lambda/\theta)$ pour $\theta > 0$. Le cas particulier $t = 1$ correspond à la *loi exponentielle de paramètre λ* , notée $\mathcal{E}(\lambda)$, de densité $\lambda e^{-\lambda x} \mathbb{1}_{\mathbb{R}_+}(x)$. Si X et Y sont deux variables aléatoires indépendantes de lois respectives $\Gamma(t_1, \lambda)$ et $\Gamma(t_2, \lambda)$ alors $X + Y$ a la loi $\Gamma(t_1 + t_2, \lambda)$. Enfin la loi $\chi^2(n)$ du chi-carré à n degrés de liberté est la loi $\Gamma(n/2, 1/2)$.

Si X_1, \dots, X_n est un n -échantillon de loi $\mathcal{E}(\lambda)$, de paramètre inconnu, alors $\bar{X}_n \rightarrow \lambda^{-1}$ et $\sqrt{n}(\bar{X}_n - \lambda^{-1}) \rightsquigarrow \mathcal{N}(0; \lambda^{-2})$. L'application de la méthode des moments conduit à estimer λ par $\hat{\lambda}_n = \bar{X}_n^{-1} \rightarrow \lambda$ lorsque $n \rightarrow +\infty$. La méthode delta implique que $\sqrt{n}(\hat{\lambda}_n - \lambda) \rightsquigarrow \mathcal{N}(0; \lambda^2)$. On peut aisément obtenir des intervalles de confiance asymptotiques pour λ . Si Z est $\mathcal{N}(0; 1)$ et $\mathbb{P}[|Z| \leq z_\alpha] = 1 - \alpha$, alors $\mathbb{P}[\sqrt{n}|\lambda\bar{X}_n - 1| \leq z_\alpha] \simeq 1 - \alpha$ pour n grand, ce qui donne un intervalle de confiance approché de la forme $(1 \pm n^{-1/2}z_\alpha)/\bar{X}_n$.

Exemple 3 : paramètres de translation et de changement d'échelle Dans un certain nombre de problèmes statistiques, les paramètres sont liés à des transformations affines d'une loi fixe connue. L'exemple le plus classique est celui de la gaussienne. Si X a une loi $\mathcal{N}(\mu, \sigma^2)$, alors $X = \sigma Y + \mu$ où Y est une gaussienne standard et le paramètre μ correspond à une translation, le paramètre σ à une homothétie. De tels paramètres apparaissent naturellement lorsque l'on fait des changements d'unités. Un autre exemple est celui des lois exponentielles. Si $X \sim \mathcal{E}(\lambda)$ alors $X = Y/\lambda$ avec $Y \sim \mathcal{E}(1)$. Enfin, si $X \sim \mathcal{U}([a, b])$, alors $X = (b - a)Y + a$ avec $Y \sim \mathcal{U}([0, 1])$.

Ces trois exemples ne sont que des illustrations d'une situation plus générale, pour laquelle on a une loi de base connue, de fonction de répartition F sur \mathbb{R} et l'on considère la famille des lois des variables $X = aY + b$ où Y a la fonction de répartition F et les paramètres sont $a > 0$ et $b \in \mathbb{R}$. La fonction de répartition G de X est donnée par

$$G(t) = \mathbb{P}[aY + b \leq t] = \mathbb{P}[Y \leq a^{-1}(t - b)] = F(a^{-1}(t - b)).$$

Si Y a une densité f , alors X a la densité $g(x) = a^{-1}f(a^{-1}(x - b))$. Les paramètres a et b sont appelés respectivement paramètres d'échelle et de translation.

Comme nous venons de le voir, un certain nombre de modèles statistiques sont de ce type, avec des paramètres qui sont fonction de paramètres d'échelle et de translation. Pour $\mathcal{E}(\lambda)$, par exemple, λ est l'inverse du paramètre d'échelle. Dans de tels modèles, les moments de la loi sont des fonctions simples des paramètres d'échelle ou de translation. Pour un modèle de translation avec $X = Y + \mu$ et $\mathbb{E}[Y] = M_1$, constante connue, on a $\mathbb{E}[X] = M_1 + \mu$ et la méthode des moments fournit $\hat{\mu}_n = \bar{X}_n - M_1$. Il est immédiat de voir qu'il s'agit d'estimateurs convergents et, si $\mathbb{E}[Y^2] < +\infty$, asymptotiquement normaux. Le raisonnement est le même dans le cas d'un paramètre d'échelle θ avec $X = \theta Y$ et $M_1 \neq 0$. On peut estimer θ par \bar{X}_n/M_1 et vérifier que l'on a encore des estimateurs convergents et asymptotiquement normaux si $\mathbb{E}[Y^2] < +\infty$.

Pour un modèle à deux paramètres, de la forme $X = \sigma Y + \mu$, en supposant, pour simplifier, que $\mathbb{E}[Y] = 0$, $\text{Var}(Y) = 1$ et $\mathbb{E}[Y^4] < +\infty$, on pourra, comme précédemment, estimer μ par \bar{X}_n et $\sigma^2 = \text{Var}(X)$ par la variance empirique $\hat{\sigma}_n^2$ (ou éventuellement \hat{s}_n^2 comme on l'a expliqué précédemment). Dans ce cas, le Théorème de Slutsky implique, une fois de plus, que

$$\sqrt{n}(\bar{X}_n - \mu)\hat{\sigma}_n^{-1} \rightsquigarrow \mathcal{N}(0, 1),$$

à partir de quoi l'on peut construire des intervalles de confiance asymptotiques pour μ .

3.2 Utilisation des quantiles empiriques

Reprenons d'abord le problème qui a motivé notre introduction des quantiles, à savoir l'estimation du paramètre μ pour la famille de densités $[\pi(1 + (x - \mu)^2)]^{-1}$ de Cauchy. Il s'agit encore d'un problème d'estimation d'un paramètre de translation mais la méthode des moments ne s'applique pas. On sait que $\mu = x_{1/2}$ est la médiane de la loi. On peut donc l'estimer par la médiane empirique $x_{1/2}(n)$ et l'on déduit du Théorème 6 que

$$\sqrt{n}(x_{1/2}(n) - \mu) \rightsquigarrow \mathcal{N}(0, \pi^2/4).$$

On peut utiliser la même méthode pour estimer la moyenne d'une gaussienne $\mathcal{N}(\mu, 1)$ et obtenir $\sqrt{n}(x_{1/2}(n) - \mu) \rightsquigarrow \mathcal{N}(0, \pi/2)$ alors que la méthode des moments donne $\sqrt{n}(\bar{X}_n - \mu) \rightsquigarrow \mathcal{N}(0, 1)$, ce qui est un peu meilleur parce que la variance asymptotique est plus faible.

D'un autre côté, la médiane empirique a des avantages que n'a pas la moyenne comme le montre l'exemple suivant. Supposons que $\mu = 0$ et que l'on ait 100 observations, dont une, la centième, est aberrante, suite à une mauvaise manipulation, et égale à 30. Alors on a $\bar{X}_n = 10^{-2} \sum_{i=1}^{99} X_i + (3/10)$, de loi (conditionnellement au fait que $X_{100} = 30$) $\mathcal{N}(3/10, 99 \cdot 10^{-4})$, ce qui est très différent de la loi $\mathcal{N}(0, 10^{-2})$ attendue. Il est facile de voir qu'alors, avec une grande probabilité, $\bar{X}_n \in [1/10, 5/10]$ car l'écart-type est à peu près $1/10$. Au contraire, la médiane empirique ne sera que très peu modifiée par une telle erreur, car remplacer une des observations par une observation aberrante peut, au pire, décaler la médiane d'un rang dans la statistique d'ordre, ce qui aura très peu d'effet, les observations du milieu de l'échantillon étant à une distance de l'ordre de n^{-1} les unes des autres. Donc la médiane est beaucoup plus stable que la moyenne en cas de données erronées. On dit qu'elle est *robuste*.

4 La méthode du maximum de vraisemblance

4.1 Le cas particulier d'un ensemble fini de paramètres

4.1.1 Tests de rapports de vraisemblance entre deux lois de probabilités

On s'intéresse ici au problème suivant : on dispose d'une observation \mathbf{X} (éventuellement multidimensionnelle, de la forme $\mathbf{X} = (X_1, \dots, X_n)$) dont la loi est soit P , soit Q mais on ignore laquelle des deux est la vraie et l'on veut tester entre ces deux hypothèses : la loi $P_{\mathbf{X}}$ de \mathbf{X} est P ou $P_{\mathbf{X}} = Q$. Comme on ne considère que deux lois, on peut toujours trouver une mesure dominante μ (par exemple $P + Q$) et décider P ou Q selon que $(dP/d\mu)(\mathbf{X}) > (dQ/d\mu)(\mathbf{X})$ ou $(dP/d\mu)(\mathbf{X}) < (dQ/d\mu)(\mathbf{X})$. En cas d'égalité, on fait un choix arbitraire. En réécrivant les choses de manière un peu différente, on obtient un test qui choisit P ou Q selon que $\log[dP/d\mu(\mathbf{X})] - \log[dQ/d\mu(\mathbf{X})]$ est positif ou négatif, c'est à dire celle des deux lois pour laquelle ce que l'on appelle la *log-vraisemblance* de \mathbf{X} est la plus grande.

Définition 14 *Étant donné une famille \mathcal{P} de probabilités sur E , dominée par une mesure μ et une variable aléatoire \mathbf{X} , la vraisemblance et la log-vraisemblance de $P \in \mathcal{P}$ sont données respectivement par*

$$\frac{dP}{d\mu}(\mathbf{X}) \quad \text{et} \quad \log \left(\frac{dP}{d\mu}(\mathbf{X}) \right),$$

avec la convention habituelle $\log 0 = -\infty$.

Les tests entre deux lois de probabilités P et Q qui choisissent entre elles selon que la valeur du rapport des vraisemblances $[(dP/d\mu)(\mathbf{X})]/[(dQ/d\mu)(\mathbf{X})]$ est plus ou moins grande sont dits tests de rapport de vraisemblances.

On remarque que si $\mathbf{X} = (X_1, \dots, X_n)$ où les $X_i \in E$ sont i.i.d. de loi \bar{P} , donc $P = \bar{P}^{\otimes n}$ est une loi produit, et que la mesure μ s'écrit $\nu^{\otimes n}$ où la mesure ν sur E domine tout, on aura

$$\frac{dP}{d\mu}(\mathbf{X}) = \prod_{i=1}^n \frac{d\bar{P}}{d\nu}(X_i) \quad \text{et} \quad \log \left(\frac{dP}{d\mu}(\mathbf{X}) \right) = \sum_{i=1}^n \log \left(\frac{d\bar{P}}{d\nu}(X_i) \right).$$

L'intérêt d'utiliser la log-vraisemblance, c'est qu'elle se présente dans ce cas sous la forme d'une somme de variables i.i.d.

Lemme 2 Si les variables X_1, \dots, X_n sont i.i.d. de loi P , si $Q \neq P$ et si P et Q sont dominées par ν , alors

$$\mathbb{P}_P \left[\sum_{i=1}^n \log \left(\frac{dQ}{d\nu}(X_i) \right) \geq \sum_{i=1}^n \log \left(\frac{dP}{d\nu}(X_i) \right) \right] \leq \rho^n(P, Q) \xrightarrow{n \rightarrow +\infty} 0,$$

où $\rho(P, Q)$ est définie par

$$0 \leq \rho(P, Q) = \int_E \sqrt{\frac{dP}{d\nu}(x) \frac{dQ}{d\nu}(x)} d\nu(x) < 1.$$

Démonstration : On remarque d'abord, en remplaçant ν par $\mu = f \cdot \nu$ que $\rho(P, Q)$ a une valeur indépendante du choix de la mesure dominante ν et ne dépend donc que de P et Q , d'où la notation. Par ailleurs, d'après l'Inégalité de Cauchy-Schwarz et le Théorème de Fubini,

$$0 \leq \rho(P, Q) \leq 1 \quad \text{et} \quad \rho(P, Q) < 1 \quad \text{si} \quad P \neq Q; \quad \rho(P^{\otimes n}, Q^{\otimes n}) = \rho^n(P, Q).$$

Pour montrer la borne de probabilité, on utilise que, quelle que soit la v.a.r. Z ,

$$\mathbb{P}[Z \geq 0] \leq \mathbb{E}[e^{tZ}] \quad \text{pour tout } t \geq 0. \quad (\text{III.4.1})$$

En effet, si $t \geq 0$,

$$\mathbb{E}[e^{tZ}] \geq \mathbb{E}[e^{tZ} \mathbf{1}_{Z \geq 0}] \geq \mathbb{E}[\mathbf{1}_{Z \geq 0}] = \mathbb{P}[Z \geq 0].$$

En appliquant (III.4.1) avec

$$Z = \log \left(\prod_{i=1}^n \frac{dQ}{d\nu}(X_i) \right) - \log \left(\prod_{i=1}^n \frac{dP}{d\nu}(X_i) \right) \quad \text{et} \quad t = 1/2,$$

on obtient la borne $\rho(P^{\otimes n}, Q^{\otimes n}) = \rho^n(P, Q)$ cherchée. \square

On a évidemment le même résultat en échangeant les rôles de P et Q , ce qui montre que les erreurs des tests de rapports de vraisemblance entre P et Q tendent vers 0 quand $n \rightarrow +\infty$.

4.1.2 Maximum de vraisemblance sur un ensemble fini

On suppose ici que l'ensemble Θ des paramètres est un ensemble fini quelconque et que le modèle est identifiable. Dans ce cas, le modèle est dominé par une mesure positive ν , par exemple par $\sum_{\theta \in \Theta} P_\theta$, avec des densités $f_\theta = dP_\theta/d\nu$. On observe X_1, \dots, X_n i.i.d. de loi inconnue P_θ . Faire tous les tests de rapport de vraisemblances entre les points de Θ revient à comparer entre elles toutes les vraisemblances et comme on est sur un ensemble fini, il y en a une qui sera plus grande que toutes les autres, donc un $\hat{\theta}$ tel que les tests de $P_{\hat{\theta}}$ contre n'importe quel autre P_θ acceptent $P_{\hat{\theta}}$. Il est alors naturel de choisir $\hat{\theta}$ comme estimateur de θ . Cet estimateur est dit *estimateur du maximum de vraisemblance* parce qu'il correspond au point pour lequel la vraisemblance est maximale. On peut aisément analyser ses performances lorsque n est grand. En effet, si θ désigne la vraie valeur du paramètre, d'après le Lemme 2 et la finitude de Θ .

$$\begin{aligned} \mathbb{P}_\theta [\hat{\theta} \neq \theta] &\leq \mathbb{P}_\theta \left[\exists \theta' \neq \theta \in \Theta \left| \sum_{i=1}^n \log \left(\frac{dP_{\theta'}}{d\nu}(X_i) \right) \geq \sum_{i=1}^n \log \left(\frac{dP_\theta}{d\nu}(X_i) \right) \right| \right] \\ &\leq \sum_{\theta' \in \Theta, \theta' \neq \theta} \mathbb{P}_\theta \left[\sum_{i=1}^n \log \left(\frac{dP_{\theta'}}{d\nu}(X_i) \right) \geq \sum_{i=1}^n \log \left(\frac{dP_\theta}{d\nu}(X_i) \right) \right] \xrightarrow{n \rightarrow +\infty} 0. \end{aligned}$$

Donc, sur un ensemble fini, les estimateurs du maximum de vraisemblance sont toujours convergents.

4.2 La méthode générale du maximum de vraisemblance

4.2.1 Définition

Définition 15 *Étant donné un modèle statistique $\{P_\theta, \theta \in \Theta\}$ dominé, la famille des densités correspondantes $dP_\theta/d\mu$, l'observation \mathbf{X} et la fonction de vraisemblance $\theta \mapsto (dP_\theta/d\mu)(\mathbf{X})$ de Θ dans \mathbb{R}_+ , on appelle estimateur du maximum de vraisemblance de θ tout point $\hat{\theta}$ de Θ tel que $(dP_{\hat{\theta}}/d\mu)(\mathbf{X}) = \sup_{\theta \in \Theta} (dP_\theta/d\mu)(\mathbf{X})$.*

Pour un modèle d'échantillonnage avec des observations X_1, \dots, X_n de densités individuelles f_θ , la log-vraisemblance s'écrit $\sum_{i=1}^n \log[f_\theta(X_i)]$ et un estimateur du maximum de vraisemblance $\hat{\theta}_n$ est caractérisé par le fait que

$$\sum_{i=1}^n \log \left(f_{\hat{\theta}_n}(X_i) \right) = \sup_{\theta \in \Theta} \sum_{i=1}^n \log \left(f_\theta(X_i) \right). \quad (\text{III.4.2})$$

Comme on vient de le voir, si Θ est fini un estimateur du maximum de vraisemblance existe et, dans un modèle d'échantillonnage, il est asymptotiquement convergent, donc asymptotiquement unique. Dans le cas général, on peut rencontrer tous les problèmes habituels liés à l'optimisation sur un ensemble infini, à savoir :

- le maximum de vraisemblance n'existe pas ;
- le maximum de vraisemblance n'est pas unique ;
- le maximum de vraisemblance existe mais on ne sait pas le calculer.

Nous verrons ci-après des exemples de telles situations. Néanmoins, dans les bons cas, c'est à dire lorsque Θ n'est pas trop gros (Θ est un compact d'un espace euclidien par exemple) et l'application $\theta \mapsto P_\theta$ est suffisamment régulière, on peut montrer que les choses se passent bien, comme dans le cas où Θ est fini.

4.2.2 Quelques exemples de modèles d'échantillonnage

Le modèle gaussien ordinaire Si l'on dispose de n observations i.i.d. de loi $\mathcal{N}(\mu, \sigma^2)$, avec deux paramètres μ et σ^2 inconnus, la log-vraisemblance s'écrit

$$-\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2. \quad (\text{III.4.3})$$

On remarque d'abord que, indépendamment de la valeur de σ^2 , la maximisation de la log-vraisemblance par rapport à μ correspond à la minimisation de $\sum_{i=1}^n (X_i - \mu)^2$, dont la solution est donnée par $\hat{\mu}_n = \bar{X}_n$. Il suffit alors de maximiser (III.4.3) par rapport à σ^2 avec $\mu = \bar{X}_n$, ce qui conduit à la variance empirique $\hat{\sigma}_n^2$ donnée par (III.2.2).

La loi de Poisson On rappelle que la loi de Poisson de paramètre $\lambda > 0$, notée $\mathcal{P}(\lambda)$, est une loi sur \mathbb{N} et si N_n est une telle variable, $\mathbb{P}_\lambda[N = k] = e^{-\lambda} \lambda^k / (k!)$ pour $k \in \mathbb{N}$. Comme dans tout modèle statistique correspondant à des v.a. à valeurs dans un espace dénombrable, il y a une mesure dominante "naturelle" qui est la mesure de comptage et la densité correspondante est $f_\lambda(k) = \mathbb{P}_\lambda[N = k]$. L'espace des paramètres est ici $\Lambda =]0, +\infty[$ qui est ouvert. Le modèle de Poisson est obtenu comme limite de modèles

binomiaux. Si $N_n \sim \mathcal{B}(n, \lambda_n/n)$ et $\lambda_n \xrightarrow{n \rightarrow +\infty} \lambda$, alors $N_n \rightsquigarrow \mathcal{P}(\lambda)$. C'est à dire que la loi de Poisson modélise des phénomènes où l'on compte le nombre de succès dans un très grand nombre d'essais chacun ayant une probabilité de succès très faible, comme pour les phénomènes radioactifs. Les propriétés de cette loi sont classiques : $\mathbb{E}_\lambda[N] = \lambda = \text{Var}_\lambda(N)$ et si les variables N_j , $1 \leq j \leq J$, sont indépendantes de lois respectives $\mathcal{P}(\lambda_j)$, alors $\sum_{j=1}^J N_j \sim \mathcal{P}\left(\sum_{j=1}^J \lambda_j\right)$.

Si l'on dispose de n v.a. indépendantes, N_1, \dots, N_n de loi $\mathcal{P}(\lambda)$, la log-vraisemblance s'écrit

$$\sum_{i=1}^n [-\lambda + N_i \log \lambda - \log(N_i!)] = -n\lambda + (\log \lambda) \sum_{i=1}^n N_i - \sum_{i=1}^n \log(N_i!),$$

dont la maximisation est immédiate et conduit à $\hat{\lambda}_n = n^{-1} \sum_{i=1}^n N_i$ (qui est aussi l'estimateur de la méthode des moments), sauf si toutes les observations sont nulles, auquel cas $\sum_{i=1}^n N_i = 0$ qui n'est pas une valeur licite puisque $0 \notin \Lambda$. On notera que cet évènement peut très bien se produire puisque il a une probabilité positive $e^{-n\lambda}$ qui peut être grande si n et λ sont petits. Par contre, asymptotiquement, la probabilité de cet évènement tend vers zéro.

Lois uniformes sur $[0, \theta]$ Pour un n -échantillon, la vraisemblance s'écrit

$$\theta^{-n} \prod_{i=1}^n \mathbb{1}_{[0, \theta]}(X_i) = \theta^{-n} \prod_{i=1}^n \mathbb{1}_{[X_i, +\infty]}(\theta) = \theta^{-n} \mathbb{1}_{[X_{(n)}, +\infty]}(\theta),$$

où $X_{(n)}$ est la plus grande des observations. Cette fonction de θ est discontinue en $X_{(n)}$, nulle avant puis positive et décroissante et atteint donc son maximum en $X_{(n)}$ qui est l'estimateur du maximum de vraisemblance, lequel ne s'obtient donc pas toujours en annulant une dérivée. Dans ce cas, il est facile de calculer sa loi. Sa fonction de répartition s'écrit $\mathbb{P}_\theta[X_{(n)} \leq t] = (t/\theta)^n$ pour $0 \leq t \leq \theta$ et la densité correspondante est $g_\theta(x) = \theta^{-n} x^{n-1} \mathbb{1}_{[0, \theta]}(x)$, de sorte que $\mathbb{E}_\theta[X_{(n)}] = n\theta/(n+1) < \theta$. C'est un estimateur biaisé, ce qui était prévisible puisque $X_{(n)} < \theta$, \mathbb{P}_θ -p.s.. Comme $\mathbb{E}_\theta[X_{(n)}^2] = n\theta^2/(n+2)$, le risque quadratique de l'estimateur du maximum de vraisemblance s'écrit

$$\mathbb{E}_\theta[(X_{(n)} - \theta)^2] = \frac{n\theta^2}{n+2} + \theta^2 - 2\frac{n\theta^2}{n+1} = \frac{2\theta^2}{(n+1)(n+2)}.$$

On peut remarquer que si l'on utilise la méthode des moments à partir de $\mathbb{E}_\theta[X_i] = \theta/2$, on trouve l'estimateur $2\bar{X}_n$ qui est sans biais et de variance $\theta^2/(3n)$. On voit que, dans ce cas, l'estimateur du maximum de vraisemblance, bien que biaisé, est bien meilleur que celui de la méthode des moments puisque son risque quadratique est, asymptotiquement, infiniment plus petit.

Lois uniformes sur $[\theta, \theta + 1]$ Dans ce cas, la vraisemblance s'écrit

$$\prod_{i=1}^n \mathbb{1}_{[\theta, \theta+1]}(X_i) = \mathbb{1}_{[\theta, +\infty]}(X_{(1)}) \mathbb{1}_{[-\infty, \theta+1]}(X_{(n)}) = \mathbb{1}_{[X_{(n)} - 1, X_{(1)}]}(\theta).$$

Cette vraisemblance ne prend que les deux valeurs 0 et 1 et son maximum est atteint en tout point de l'intervalle (p.s. non vide) $[X_{(n)} - 1, X_{(1)}]$. Il n'y a donc pas, dans

ce cas, unicité du maximum de vraisemblance. On peut néanmoins vérifier aisément en calculant, comme on l'a fait déjà plus haut, la loi de $X_{(n)}$ et celle de $X_{(1)}$, que $X_{(1)} \xrightarrow[n \rightarrow +\infty]{} \theta$ et $X_{(n)} \xrightarrow[n \rightarrow +\infty]{} \theta + 1$, \mathbb{P}_θ -p.s., de sorte qu'asymptotiquement tous ces estimateurs sont très proches et tous convergent vers θ . Par contre, lorsque n est petit, on tombe sur une réelle indétermination.

Translatées de la loi de Cauchy La loi de Cauchy a la densité $\pi^{-1}(1+x^2)^{-1}$. Elle donne lieu à un modèle de translation correspondant à la famille des densités $\pi^{-1}(1+(x-\theta)^2)^{-1}$ avec $\theta \in \mathbb{R}$. Dans un tel modèle, la log-vraisemblance d'un n -échantillon s'écrit $-n \log \pi - \sum_{i=1}^n \log(1+(X_i-\theta)^2)$. Il est facile de voir qu'une telle fonction de θ a au moins un maximum sur \mathbb{R} , lequel s'obtient en annulant la dérivée mais qu'il est impossible de résoudre analytiquement l'équation correspondante car cela revient à chercher les racines d'un polynôme de degré $2n-1$. La méthode fondée sur la médiane empirique est clairement beaucoup plus simple.

4.2.3 Extension

On est souvent amené, dans un modèle statistique, à considérer plusieurs paramétrages. Un exemple simple est le suivant. Si l'on dispose d'observations de loi exponentielle, les densités peuvent s'écrire $\theta \exp(-\theta x) \mathbb{1}_{\mathbb{R}_+}(x)$ avec $\theta > 0$ et l'espérance est θ^{-1} . On peut tout aussi bien paramétrer la même famille de lois par son espérance en considérant les densités $\lambda^{-1} \exp(-x/\lambda) \mathbb{1}_{\mathbb{R}_+}(x)$ avec $\lambda > 0$. C'est évidemment la même famille de lois. On a simplement fait un changement de paramètre bijectif $\lambda = \theta^{-1}$. On voit immédiatement que les estimateurs du maximum de vraisemblance construits à partir d'un n -échantillon vérifient $\hat{\lambda}_n = \hat{\theta}_n^{-1}$. C'est un phénomène général. Si l'on fait, dans un modèle statistique, un changement de paramètre bijectif $\lambda = g(\theta)$, les estimateurs du maximum de vraisemblance construits à partir d'un n -échantillon sont liés par $\hat{\lambda}_n = g(\hat{\theta}_n)$.

De manière plus générale, si l'on cherche à estimer une fonction $\lambda = g(\theta)$ du paramètre, on appellera estimateur du maximum de vraisemblance de λ l'estimateur $g(\hat{\theta}_n)$, même si g n'est pas bijective. Dans tous les cas, pour étudier les propriétés asymptotiques de $g(\hat{\theta}_n)$ connaissant celles de $\hat{\theta}_n$ il suffira d'appliquer la méthode delta.

4.3 Application au modèle exponentiel à un paramètre

4.3.1 Introduction

Les modèles exponentiels constituent une classe générique de modèles statistiques qui inclut un certain nombre d'exemples classiques. Que peut-on trouver de commun entre la famille des lois $\mathcal{N}(\theta, 1)$ et celle des lois de Poisson $\mathcal{P}(\theta)$? Si l'on écrit la densité $\mathcal{N}(\theta, 1)$ par rapport à la mesure dominante $\mathcal{N}(0, 1)$ de densité $(2\pi)^{-1/2} \exp[-x^2/2]$, on trouve $\exp[-\theta^2/2] \exp[\theta x]$ et, pour la loi de Poisson et la mesure dominante $(x!)^{-1} d\mu(x)$ où μ est la mesure de comptage sur \mathbb{N} , on trouve $e^{-\theta} \exp(-x \log \theta)$. Dans les deux cas, après avoir choisi une mesure dominante convenable ne dépendant que des paramètres connus, la densité se décompose comme un produit de deux facteurs, le premier ne dépendant que des paramètres et le second étant l'exponentielle du produit d'une fonction de la variable par une fonction du paramètre. Comme on va le voir, cette décomposition particulière intervient également dans d'autres modèles classiques, d'où l'intérêt d'étudier de façon systématique cette classe de modèles.

4.3.2 Préliminaires d'analyse

Étant donné, sur l'espace mesurable E , une mesure positive ν et une fonction mesurable T , non presque sûrement constante pour ν , on peut considérer l'ensemble

$$H = \left\{ \eta \in \mathbb{R} \mid \int_E \exp[\eta T(x)] d\nu(x) < +\infty \right\}. \quad (\text{III.4.4})$$

On remarquera que si T est constante ν p.s., on arrive à des trivialités.

Si cet ensemble n'est pas vide, la convexité de la fonction exponentielle entraîne que H est un intervalle : si η_1 et η_2 appartiennent à H , il en va de même de $\lambda\eta_1 + (1-\lambda)\eta_2$ pour $\lambda \in [0; 1]$ puisque $\exp[\lambda\eta_1 T + (1-\lambda)\eta_2 T] \leq \lambda \exp[\eta_1 T] + (1-\lambda) \exp[\eta_2 T]$. On supposera dans toute la suite que cet intervalle n'est pas réduit à un seul point, donc d'intérieur $\overset{\circ}{H}$ non vide. Un résultat classique d'analyse concernant les transformées de Laplace dit la chose suivante :

Proposition 4 *La fonction $\eta \mapsto J(\eta) = \int_E \exp[\eta T(x)] d\nu(x)$ est indéfiniment différentiable sur $\overset{\circ}{H}$ de dérivées successives*

$$J^{(i)}(\eta) = \int [T(x)]^i \exp[\eta T(x)] d\nu(x).$$

4.3.3 Le modèle exponentiel naturel

Étant donné ν et T comme ci-dessus, on peut définir une famille de densités par rapport à la mesure ν , indexée par H , $\{g_\eta, \eta \in H\}$, en posant

$$g_\eta(x) = \exp[\eta T(x) - A(\eta)] \quad \text{avec } A(\eta) = \log(J(\eta)) = \log\left(\int_E \exp[\eta T(x)] d\nu(x)\right). \quad (\text{III.4.5})$$

L'ensemble des lois $R_\eta = g_\eta \cdot \nu$ forme un *modèle exponentiel naturel* (ou *canonique*) indexé par l'espace des paramètres naturels H . Par exemple, les lois exponentielles ordinaires de densités $\eta \exp[-\eta x]$ par rapport à la mesure de Lebesgue ν sur \mathbb{R}_+ forment un tel modèle avec $H =]0; +\infty[$, $T(x) = -x$ et $A(\eta) = -\log \eta$. Il est important de noter que si la famille des lois R_η est définie de manière unique, sa représentation sous la forme d'un modèle exponentiel ne l'est pas. En effet, on peut remplacer la mesure ν par la mesure $a\nu$ avec $a > 0$ et la densité g_η par $a^{-1}g_\eta$, c'est à dire A par $A + \log a$, sans changer le modèle statistique.

i) Ce modèle est identifiable car si on avait $g_\eta = g_{\eta'}$ pour $\eta \neq \eta'$, on aurait

$$(\eta - \eta')T(x) = A(\eta) - A(\eta') \quad \nu \text{ p.s.},$$

ce qui impliquerait que T est p.s. constante, ce que l'on a exclu.

ii) Les densités g_η sont strictement positives, donc toutes les lois R_η sont mutuellement absolument continues et toutes équivalentes à ν .

iii) Les propriétés de la transformée de Laplace vues ci-dessus impliquent que la fonction A définie par (III.4.5) est indéfiniment différentiable sur $\overset{\circ}{H}$, les dérivées se calculant en dérivant sous le signe \int . On obtient ainsi

$$A'(\eta) = \frac{\int_E T(x) \exp[\eta T(x)] d\nu(x)}{\int_E \exp[\eta T(x)] d\nu(x)} = e^{-A(\eta)} \int_E T(x) \exp[\eta T(x)] d\nu(x) = \mathbb{E}_\eta[T(X)]. \quad (\text{III.4.6})$$

Un calcul analogue montre que

$$A''(\eta) = \text{Var}_\eta(T(X)) \geq 0, \quad (\text{III.4.7})$$

les espérances étant prises, dans les deux cas, pour la loi R_η . La variance de $T(X)$ ne peut s'annuler (sinon $T(X)$ serait ν -p.s. constante), donc la fonction A est strictement convexe.

iv) Si l'on dispose d'un n -échantillon X_1, \dots, X_n de loi R_η du modèle exponentiel précédent, la famille des lois $R_\eta^{\otimes n}$ du vecteur $\mathbf{X} = (X_1, \dots, X_n)$ forme un nouveau modèle exponentiel par rapport à la mesure $\nu^{\otimes n}$, de densités de la forme

$$g_\eta(x_1) \cdots g_\eta(x_n) = \exp \left[\eta \sum_{i=1}^n T(x_i) - nA(\eta) \right].$$

Donc lorsque l'on dispose d'un n -échantillon, il suffit de vérifier que le modèle est exponentiel en considérant le cas particulier $n = 1$.

4.3.4 Le modèle exponentiel général

Définition 16 *Le modèle exponentiel général à un paramètre est la donnée d'une famille de densités sur E (par rapport à une mesure μ de référence) de la forme*

$$f_\theta(x) = C(\theta) \exp[Q(\theta)T(x)]h(x) \quad \text{avec } Q(\theta), T(x) \in \mathbb{R}, \quad C(\theta), h(x) \in \mathbb{R}_+, \quad (\text{III.4.8})$$

et $\theta \in \Theta$ où Θ est un intervalle (non dégénéré) de \mathbb{R} . La fonction Q est continuellement différentiable, de dérivée $Q' \neq 0$ sur Θ et la fonction T n'est pas p.s. constante pour la mesure $\nu = h \cdot \mu$.

Comme il s'agit d'une famille de densités, on a

$$C(\theta) = \left[\int_E \exp[Q(\theta)T(x)]h(x) d\mu(x) \right]^{-1}.$$

On notera que Q' a un signe constant, donc que Q est strictement monotone, que l'application $\theta \mapsto Q(\theta)$ est injective et que $Q(\Theta)$ est un intervalle d'intérieur non vide. Par ailleurs, on peut remarquer que la représentation précédente n'est pas unique puisque, quel que soit $\alpha \neq 0$, on peut toujours remplacer $Q(\theta)$ par $\alpha Q(\theta)$ et $T(x)$ par $\alpha^{-1}T(x)$.

On peut toujours associer à ce modèle exponentiel général, par le changement de paramètre $\eta = Q(\theta)$, un modèle naturel, ensemble des densités g_η , $\eta \in H$, données par (III.4.5) avec $\nu = h \cdot \mu$ et $A(\eta) = -\log C(\theta) = -\log(Q^{-1}(\eta))$. L'espace H des paramètres naturels donné par (III.4.4) contient $Q(\Theta)$; il est donc d'intérieur non vide. Remarquons qu'il se peut fort bien que H soit plus grand que $Q(\Theta)$. Donc on passe du modèle naturel au modèle général par un changement de paramètre $\theta = Q^{-1}(\eta)$ précédé, éventuellement, d'une restriction de H à $Q(\Theta)$. On notera \mathbb{E}_θ les espérances évaluées dans le modèle général et \mathbb{E}_η celles évaluées dans le modèle canonique. On a évidemment la relation

$$\int k(x) f_\theta(x) d\mu(x) = \mathbb{E}_\theta[k(X)] = \mathbb{E}_\eta[k(X)] = \int k(x) g_\eta(x) d\nu(x) \quad \text{si } \eta = Q(\theta). \quad (\text{III.4.9})$$

Les propriétés du modèle général se déduisent immédiatement de celles du modèle canonique.

- i) Le modèle est identifiable.
- ii) Les densités f_θ sont > 0 sur l'ensemble $\{x \mid h(x) > 0\}$ et toutes les lois associées $P_\theta = f_\theta \cdot \mu$ sont mutuellement absolument continues de même support que $\nu = h \cdot \mu$.
- iii) Le n -échantillon correspondant forme aussi un modèle exponentiel avec la famille des densités produit

$$f_\eta(x_1) \cdots f_\eta(x_n) = [C(\theta)]^n \exp \left[Q(\theta) \sum_{i=1}^n T(x_i) \right] h(x_1) \cdots h(x_n)$$

par rapport à la mesure $\mu^{\otimes n}$.

4.3.5 Exemples de modèles exponentiels

Le modèle binomial : Si X est une variable de loi $\mathcal{B}(n, \theta)$ avec $0 < \theta < 1$, sa densité par rapport à la mesure de comptage μ sur \mathbb{N} s'écrit

$$\binom{n}{x} \theta^x (1 - \theta)^{n-x} \mathbb{1}_{\{0, \dots, n\}}(x) = (1 - \theta)^n \exp \left[x \log \left(\frac{\theta}{1 - \theta} \right) \right] \binom{n}{x} \mathbb{1}_{\{0, \dots, n\}}(x),$$

ce qui correspond bien à une décomposition de la forme (III.4.8). Notons que ceci ne fonctionne plus si $\theta = 0$ ou 1 car la loi binomiale dégénère alors en une masse de Dirac qui n'est pas équivalente aux autres lois de la famille.

Les lois de Poisson : Comme on l'a vu, la densité par rapport à la mesure de comptage sur \mathbb{N} de la loi de Poisson $\mathcal{P}(\lambda)$, $\lambda > 0$, s'écrit

$$e^{-\lambda} \lambda^x / x! = e^{-\lambda} \exp(x \log \lambda) (x!)^{-1}.$$

Les densités gaussiennes : Elles s'écrivent

$$\frac{1}{\sigma \sqrt{2\pi}} \exp \left[-\frac{(x - \mu)^2}{2\sigma^2} \right],$$

ce qui donne les décompositions suivantes, selon que c'est μ ou σ^2 qui est le paramètre à estimer. Si σ est connu, on peut écrire la densité sous la forme

$$\frac{1}{\sigma \sqrt{2\pi}} \exp \left[-\frac{\mu^2}{2\sigma^2} \right] \exp \left[\frac{\mu x}{\sigma^2} \right] \exp \left[-\frac{x^2}{2\sigma^2} \right],$$

qui correspond à un modèle exponentiel naturel dominé par la loi $\mathcal{N}(0; \sigma^2)$ de densité $(2\pi\sigma^2)^{-1/2} \exp[-x^2/(2\sigma^2)] dx$ avec $T(x) = x/\sigma^2$. Si μ est connu, on peut choisir $Q(\sigma^2) = 1/\sigma^2$ et $T(x) = -(x - \mu)^2/2$, entre autres représentations possibles.

Les lois gamma : La densité $\Gamma(t, \lambda)$ s'écrit $[\Gamma(t)]^{-1} \lambda^t x^{t-1} e^{-\lambda x} \mathbb{1}_{\mathbb{R}_+}(x)$ et peut donc, selon que λ ou t est inconnu, se mettre sous l'une des formes suivantes :

$$[\Gamma(t)]^{-1} \lambda^t \exp(-\lambda x) x^{t-1} \mathbb{1}_{\mathbb{R}_+}(x) \quad \text{ou} \quad [\Gamma(t)]^{-1} \lambda^t \exp(t \log x) e^{-\lambda x} x^{-1} \mathbb{1}_{\mathbb{R}_+}(x).$$

Les modèles uniformes : Une loi uniforme sur un intervalle $[a; b]$ a une densité à support $[a; b]$ par définition. Les paramètres étant les extrémités de l'intervalle, le support de la loi varie avec eux, ce qui démontre que les lois uniformes ne peuvent jamais correspondre à des modèles exponentiels.

4.3.6 Le maximum de vraisemblance dans le modèle exponentiel

Commençons par étudier le modèle naturel. Si l'on dispose de n observations indépendantes X_1, \dots, X_n de densité g_η , la log-vraisemblance s'écrit alors

$$L_\eta(X_1, \dots, X_n) = \sum_{i=1}^n \eta T(X_i) - nA(\eta).$$

Comme, sur $\overset{\circ}{H}$, $A'' > 0$, L_η est une fonction strictement concave de η . Donc, s'il existe un point $\hat{\eta}_n$ intérieur à H tel que $L'_{\hat{\eta}_n}(X_1, \dots, X_n) = 0$, c'est l'unique estimateur du maximum de vraisemblance de η . Or, si l'on note $\bar{T}_n = n^{-1} \sum_{i=1}^n T(X_i)$, la solution de l'équation $L'_{\hat{\eta}_n}(X_1, \dots, X_n) = 0$ est donnée par

$$A'(\hat{\eta}_n) = \mathbb{E}_{\hat{\eta}_n}[T(X)] = \bar{T}_n \quad \text{ou} \quad \hat{\eta}_n = \phi^{-1}(\bar{T}_n) \quad \text{avec} \quad \phi(\eta) = A'(\eta). \quad (\text{III.4.10})$$

En effet, comme $\phi' = A'' > 0$, la fonction ϕ est continue et strictement croissante, donc bijective. En outre, son inverse ϕ^{-1} est différentiable et $(\phi^{-1})'[\phi(\eta)] = 1/A''(\eta)$.

Supposons que la vraie valeur η du paramètre est intérieure à H , ce qui implique que $\phi(\eta)$ est intérieur à $\phi(H)$. La loi des grands nombres et le Théorème limite central impliquent, d'après (III.4.6) et (III.4.7) que

$$\bar{T}_n \xrightarrow[n \rightarrow +\infty]{} \mathbb{E}_\eta[T(X)] = \phi(\eta) \in \phi(\overset{\circ}{H}) \quad \mathbb{P}_\eta\text{-p.s.} \quad \text{et} \quad \sqrt{n}(\bar{T}_n - \phi(\eta)) \rightsquigarrow \mathcal{N}(0, \phi'(\eta)).$$

Donc, p.s., pour n assez grand, \bar{T}_n est aussi intérieur à $\phi(H)$ et (III.4.10) a une solution unique qui converge vers $\phi^{-1}(\phi(\eta)) = \eta$. Ceci implique que, pour n grand, l'estimateur du maximum de vraisemblance $\hat{\eta}_n$ défini par (III.4.10) existe et qu'il converge vers la vraie valeur η du paramètre lorsque $n \rightarrow +\infty$. De plus, comme ϕ^{-1} est différentiable, par la méthode delta,

$$\sqrt{n}(\hat{\eta}_n - \eta) = \sqrt{n}[\phi^{-1}(\bar{T}_n) - \phi^{-1}(\phi(\eta))] \rightsquigarrow \mathcal{N}\left(0, \frac{\phi'(\eta)}{[A''(\eta)]^2}\right) = \mathcal{N}\left(0, \frac{1}{A''(\eta)}\right),$$

ou encore,

$$\sqrt{n}(\hat{\eta}_n - \eta) \rightsquigarrow \mathcal{N}(0, 1/\text{Var}_\eta(T)). \quad (\text{III.4.11})$$

Donc l'estimateur du maximum de vraisemblance est convergent et asymptotiquement normal. Remarquons qu'ici l'estimateur du maximum de vraisemblance est un estimateur de la méthode des moments appliqué à la fonction $\phi(\eta) = \mathbb{E}_\eta[T(X)]$.

Pour passer au modèle général, il suffit simplement d'appliquer la méthode delta une seconde fois. En effet, si l'on dispose de n observations indépendantes X_1, \dots, X_n de densité f_θ donnée par (III.4.8) et l'on suppose que le vrai paramètre θ est à l'intérieur de l'espace Θ des paramètres, dans le modèle naturel associé, le vrai paramètre $\eta = Q(\theta)$ est (compte-tenu des propriétés de Q), intérieur à $Q(\Theta)$, donc à H . D'après ce qui précède $\hat{\eta}_n \rightarrow \eta = Q(\theta)$, \mathbb{P}_θ p.s. Donc, pour n assez grand, $\hat{\eta}_n \in Q(\Theta)$ et l'on peut définir $\hat{\theta}_n = Q^{-1}(\hat{\eta}_n) \rightarrow Q^{-1}(\eta) = \theta$, \mathbb{P}_θ p.s. C'est l'estimateur du maximum de vraisemblance de θ puisque Q^{-1} correspond à un changement de paramètre. On peut à nouveau appliquer la méthode delta à partir de (III.4.11) en utilisant le fait que $(Q^{-1})'(Q(\theta)) = 1/Q'(\theta)$, ce qui donne, compte-tenu de (III.4.9),

$$\sqrt{n}(\hat{\theta}_n - \theta) = \sqrt{n}(Q^{-1}(\hat{\eta}_n) - Q^{-1}(\eta)) \rightsquigarrow \mathcal{N}\left(0, [\text{Var}_\theta(T)(Q'(\theta))^2]^{-1}\right).$$

On peut remarquer que $\hat{\theta}_n = Q^{-1}(\hat{\eta}_n)$ est la solution (unique) de l'équation $\mathbb{E}_{\hat{\theta}_n}[T(X)] = \bar{T}_n$. On a ainsi démontré le théorème suivant :

Théorème 7 Soit X_1, \dots, X_n un n -échantillon d'un modèle exponentiel régulier de densités $f_\theta(x) = C(\theta) \exp[Q(\theta)T(x)]h(x)$ par rapport à la loi μ , où θ décrit un intervalle Θ de \mathbb{R} . Supposons que $\theta \in \overset{\circ}{\Theta}$. Alors, \mathbb{P}_θ -p.s., pour n assez grand, l'estimateur du maximum de vraisemblance $\hat{\theta}_n$ existe et c'est l'unique solution de l'équation

$$\mathbb{E}_{\hat{\theta}_n}[T(X)] = n^{-1} \sum_{i=1}^n T(X_i). \quad (\text{III.4.12})$$

De plus,

$$\sqrt{n}(\hat{\theta}_n - \theta) \rightsquigarrow \mathcal{N}\left(0, [\text{Var}_\theta(T)(Q'(\theta))^2]^{-1}\right).$$

Remarque : Les calculs précédents nous ont permis de démontrer le Théorème 7 mais ne sont pas très commodes. En pratique on écrira simplement la log-vraisemblance que l'on maximisera en θ par un calcul direct plutôt que d'utiliser (III.4.12). On évaluera ensuite la variance asymptotique en appliquant la méthode delta.

4.3.7 Applications

Dans le modèle de Poisson $\mathcal{P}(\theta)$, $\theta > 0$, $\hat{\theta}_n = \bar{X}_n$ pourvu que $\bar{X}_n > 0$, sinon l'estimateur du maximum de vraisemblance n'existe pas. Mais il existe toujours pour n suffisamment grand puisque $\bar{X}_n \xrightarrow{n \rightarrow +\infty} \theta > 0$ avec

$$\sqrt{n}(\hat{\theta}_n - \theta) \rightsquigarrow \mathcal{N}(0, \theta) \quad \text{et} \quad \mathbb{E}_\theta \left[(\hat{\theta}_n - \theta)^2 \right] = \theta/n.$$

La condition “ n assez grand” n'est pas superflue car on peut avoir $\bar{X}_n = 0$ si n est trop petit.

Dans le modèle $\mathcal{N}(\theta, 1)$ avec $\theta \in \Theta = [0; +\infty[$, l'estimateur du maximum de vraisemblance se calcule explicitement et l'on trouve $\hat{\theta}_n = \max\{\bar{X}_n; 0\}$. Comme $\bar{X}_n \sim \mathcal{N}(\theta; n^{-1})$, asymptotiquement, si $\theta > 0$, $\hat{\theta}_n = \bar{X}_n$ qui est la solution de l'équation (III.4.12) et tout se passe comme d'habitude. Par contre, si $\theta = 0$ est sur la frontière de l'espace des paramètres, $\bar{X}_n \sim \mathcal{N}(0; n^{-1})$ et $\hat{\theta}_n = 0$ avec une probabilité 1/2, ce qui montre bien que l'on n'a ni $\hat{\theta}_n$ comme solution de l'équation (III.4.12), ni la normalité asymptotique.

4.4 Propriétés asymptotiques des estimateurs du maximum de vraisemblance

4.4.1 Convergence

Sous des hypothèses convenables, mais un peu compliquées à énoncer et que nous ne précisons pas ici, les estimateurs du maximum de vraisemblance sont convergents. Ceci dit, ces hypothèses ne sont pas toujours satisfaites, loin de là, et l'on peut fabriquer de nombreux contre-exemples. Le plus souvent, d'un point de vue pratique, on pourra calculer les estimateurs du maximum de vraisemblance et montrer, par un calcul direct, qu'ils sont convergents comme c'est le cas dans le modèle exponentiel si le vrai paramètre θ est intérieur à Θ .

4.4.2 Retour sur le modèle exponentiel

Considérons un modèle exponentiel naturel donné par la famille des densités $\{g_\eta, \eta \in H\}$ définies par (III.4.5) et notons $\ell_\eta = \log g_\eta = \eta T - A(\eta)$, les dérivations étant toujours effectuées par rapport à η . Alors, pour $\eta \in \overset{\circ}{H}$, ℓ_η , comme fonction de η , est \mathcal{C}_∞ et, d'après (III.4.6) et (III.4.7),

$$\mathbb{E}_\eta [\ell'_\eta(X)] = 0; \quad \text{Var}_\eta (\ell'_\eta(X)) = A''(\eta) = \text{Var}_\eta(T(X)) = -\mathbb{E}_\eta [\ell''_\eta(X)].$$

Il s'ensuit que l'on peut définir la fonction $I(\eta)$ par

$$I(\eta) = \int \frac{(g'_\eta(x))^2}{g_\eta(x)} d\nu(x) = \mathbb{E}_\eta [(\ell'_\eta(X))^2] = \text{Var}_\eta (\ell'_\eta(X)) = -\mathbb{E}_\eta [\ell''_\eta(X)],$$

que cette fonction est continue et que, d'après (III.4.11),

$$\sqrt{n}(\hat{\eta}_n - \eta) \rightsquigarrow \mathcal{N}(0, [I(\eta)]^{-1}).$$

Si l'on se place maintenant dans le cadre du modèle exponentiel général avec une famille de densités $\{f_\theta, \theta \in \Theta\}$, données par (III.4.8), si $Q \in \mathcal{C}_2$ et $\theta \in \overset{\circ}{\Theta}$, on peut définir de la même manière $\ell_\theta = \log f_\theta$ et

$$I(\theta) = \int \frac{(f'_\theta(x))^2}{f_\theta(x)} d\mu(x). \quad (\text{III.4.13})$$

De simples changements de variables et la méthode delta montrent que l'on a encore

$$\mathbb{E}_\theta [\ell'_\theta(X)] = 0 \quad \text{et} \quad I(\theta) = \text{Var}_\theta (\ell'_\theta(X)) = -\mathbb{E}_\theta [\ell''_\theta(X)], \quad (\text{III.4.14})$$

ainsi que la normalité asymptotique des estimateurs du maximum de vraisemblance :

$$\sqrt{n}(\hat{\theta}_n - \theta) \rightsquigarrow \mathcal{N}(0, [I(\theta)]^{-1}). \quad (\text{III.4.15})$$

La fonction $I(\theta)$ est dite *Information de Fisher* du modèle statistique et l'on peut conclure que, dans un modèle exponentiel avec $Q \in \mathcal{C}_2$, l'estimateur du maximum de vraisemblance (si $\theta \in \overset{\circ}{\Theta}$) est asymptotiquement normal de variance asymptotique l'inverse de l'Information de Fisher.

4.4.3 Normalité asymptotique

Considérons maintenant un modèle paramétrique dominé par μ et donné par une famille de densités $\{f_\theta, \theta \in \Theta \subset \mathbb{R}\}$. Sous des hypothèses convenables de dérivabilité, d'intégrabilité et de domination (on demandera en particulier que $\ell_\theta = \log f_\theta \in \mathcal{C}_2$ et que $\mathbb{E}_\theta [(\ell'_\theta(X))^2] < +\infty$) on peut encore définir l'Information de Fisher $I(\theta)$ du modèle par (III.4.13) et vérifier que (III.4.14) est encore satisfaite. Si ces hypothèses convenables (que je ne préciserai pas ici) sur la famille de densités $\{f_\theta, \theta \in \Theta\}$ sont satisfaites, le modèle est dit *régulier*.

Supposons maintenant que la vraie valeur du paramètre soit $\theta_0 \in \overset{\circ}{\Theta}$, qu'une suite d'estimateurs du maximum de vraisemblance $(\hat{\theta}_n)_{n \geq 1}$ existe et qu'elle converge vers θ_0 . La formule des accroissements finis entre θ_0 et $\hat{\theta}_n$ appliquée à la dérivée de la log-vraisemblance nous donne

$$\sum_{i=1}^n \ell'_{\hat{\theta}_n}(X_i) - \sum_{i=1}^n \ell'_{\theta_0}(X_i) = (\hat{\theta}_n - \theta_0) \sum_{i=1}^n \ell''_{\theta'_n}(X_i),$$

où (θ'_n) est une suite qui converge vers θ_0 en probabilité en même temps que $\hat{\theta}_n$. Comme la log-vraisemblance $\sum_{i=1}^n \ell_\theta(X_i)$ est dérivable, sa dérivée s'annule en $\hat{\theta}_n$, ce qui permet d'écrire

$$0 = \frac{1}{\sqrt{n}} \sum_{i=1}^n \ell'_{\hat{\theta}_n}(X_i) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \ell'_{\theta_0}(X_i) + \sqrt{n} (\hat{\theta}_n - \theta_0) \left[\frac{1}{n} \left(\sum_{i=1}^n \ell''_{\theta_0}(X_i) \right) + R_n \right],$$

avec

$$R_n = R_n(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n \left[\ell''_{\hat{\theta}_n}(X_i) - \ell''_{\theta_0}(X_i) \right].$$

Finalement,

$$\sqrt{n} (\hat{\theta}_n - \theta_0) = \left[\frac{1}{\sqrt{n}} \sum_{i=1}^n \ell'_{\theta_0}(X_i) \right] \left[-\frac{1}{n} \left(\sum_{i=1}^n \ell''_{\theta_0}(X_i) \right) - R_n \right]^{-1}.$$

Or nos hypothèses (III.4.14), le T.L.C. et la loi des grands nombres entraînent que

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \ell'_{\theta_0}(X_i) \rightsquigarrow \mathcal{N}(0, I(\theta_0)) \quad \text{et} \quad \frac{1}{n} \sum_{i=1}^n \ell''_{\theta_0}(X_i) \xrightarrow{P} -I(\theta_0).$$

Comme les hypothèses de régularité que l'on a mis sur le modèle entraînent que $R_n \xrightarrow{P} 0$, le Théorème de Slutsky implique que (III.4.15) est encore satisfaite. On obtient alors le théorème suivant.

Théorème 8 *Sous des hypothèses de régularité convenables, les estimateurs du maximum de vraisemblance $\hat{\theta}_n$, s'ils convergent vers le vrai paramètre θ_0 , sont asymptotiquement normaux de variance asymptotique l'inverse de l'Information de Fisher, c'est à dire que*

$$\sqrt{n} (\hat{\theta}_n - \theta_0) \rightsquigarrow \mathcal{N}(0, [I(\theta_0)]^{-1}).$$

4.5 Estimateurs dits “à un pas”

4.5.1 Construction

Il se trouve que les estimateurs du maximum de vraisemblance sont souvent difficiles à calculer. Or, dans un modèle régulier, tout estimateur convenablement construit (estimateur empirique construit à partir d'un moment ou d'un quantile empirique par exemple) converge à la vitesse $n^{-1/2}$, il est donc proche à la fois de la vraie valeur du paramètre et de l'estimateur du maximum de vraisemblance. Une idée naturelle dans ce contexte est de chercher à le modifier pour en faire un estimateur équivalent au maximum de vraisemblance, donc asymptotiquement aussi bon.

Le principe de la construction repose sur la *méthode de Newton* que l'on utilise pour résoudre numériquement une équation de la forme $g(x) = 0$. Partant d'une approximation x_1 de la racine x_0 , on fait comme si g était linéaire au voisinage de x_1 , c'est à dire que l'on remplace g par sa tangente $g_1(x) = g(x_1) + (x - x_1)g'(x_1)$ et que l'on résout l'équation $g_1(x) = 0$, d'où la solution $x_2 = x_1 - g(x_1)/g'(x_1)$. On itère la procédure qui produit ainsi une suite d'approximations $(x_n)_{n \geq 1}$ qui, sous des hypothèses convenables, converge vers x_0 .

Dans notre cas, et en supposant que la log-vraisemblance $L_n(\theta) = \sum_{i=1}^n \ell_\theta(X_i)$ est deux fois continuellement différentiable, on cherche une solution de l'équation $L'_n(t) = 0$.

Partant d'une solution approchée T_n , la méthode de Newton en fournit une nouvelle, $T_n - L'_n(T_n)/L''_n(T_n)$. On utilise alors le fait que $\mathbb{E}_\theta [L''_n(\theta)] = -nI(\theta)$, I est continue et T_n est proche de θ pour remplacer $L''_n(T_n)$ par $-nI(T_n)$. On obtient ainsi un nouvel estimateur

$$\tilde{\theta}_n(X_1, \dots, X_n) = T_n(X_1, \dots, X_n) + \frac{L'_n(T_n)}{nI(T_n)} = T_n(X_1, \dots, X_n) + \frac{\sum_{i=1}^n l'_{T_n}(X_i)}{nI(T_n)}. \quad (\text{III.4.16})$$

La méthode classique de Newton est itérative mais ici on ne fait que le premier pas de la méthode. En effet, notre but n'est pas de trouver la racine de l'équation $L'_n(t) = 0$ (c'est à dire l'estimateur du maximum de vraisemblance $\hat{\theta}_n$) mais d'estimer θ . Comme $\hat{\theta}_n - \theta$ est de l'ordre de $n^{-1/2}$, il suffit de remplacer $\hat{\theta}_n$ par une quantité de la forme $\hat{\theta}_n + o(n^{-1/2})$, ce qui est le cas de $\tilde{\theta}_n$.

4.5.2 Comportement asymptotique

Les propriétés asymptotiques de l'estimateur $\tilde{\theta}_n$ dépendent évidemment de celles de T_n et, pour que la méthode fonctionne, il faut partir d'un estimateur T_n suffisamment proche de θ_0 .

Théorème 9 *Considérons un modèle statistique régulier $\{f_\theta, \theta \in \Theta\}$ pour lequel il existe une suite $(T_n(X_1, \dots, X_n))_{n \geq n_0}$ d'estimateurs telle que la suite $\sqrt{n}(T_n - \theta_0)$ converge en loi lorsque θ_0 désigne le vrai paramètre avec $I(\theta_0) > 0$. Soit $\tilde{\theta}_n(X_1, \dots, X_n)$ l'estimateur défini par (III.4.16). Alors la suite $(\tilde{\theta}_n)_{n \geq 1}$ est convergente et vérifie*

$$\sqrt{n}(\tilde{\theta}_n - \theta_0) \rightsquigarrow \mathcal{N}(0, [I(\theta_0)]^{-1}).$$

On pourra en particulier utiliser, comme suites d'estimateurs préliminaires, celles qui sont asymptotiquement normales. Des exemples classiques sont les estimateurs de la méthode des moments ou ceux de la méthode des quantiles. Ceci s'applique en particulier à l'estimation du paramètre d'un modèle de translation de Cauchy pour lequel l'estimateur du maximum de vraisemblance est difficile à calculer. On partira de la médiane empirique T_n et l'on construira $\tilde{\theta}_n$ comme indiqué ci-dessus.

5 Les estimateurs de Bayes

5.1 L'approche Bayésienne

Nous avons vu que, pour évaluer la qualité d'un estimateur $\hat{\theta}$ de θ , il est commode d'utiliser son risque quadratique qui est la fonction

$$\theta \mapsto R(\theta, \hat{\theta}) = \mathbb{E}_\theta[(\theta - \hat{\theta}(\mathbf{X}))^2] = \int_{\mathcal{X}} [\theta - \hat{\theta}(x)]^2 dP_\theta(x).$$

Mais si cette fonction permet d'évaluer les performances d'un estimateur, elle permet rarement de comparer deux estimateurs $\hat{\theta}$ et $\hat{\theta}'$ entre eux. En effet, le plus souvent, leurs fonctions de risque ne sont pas comparables parce que l'on n'a ni $R(\theta, \hat{\theta}) \leq R(\theta, \hat{\theta}')$ quel que soit θ , ni l'inégalité inverse. Un exemple très simple dans le modèle binomial avec $N \sim \mathcal{B}(n, \theta)$ est donné par $\hat{\theta}_1(N) = N/n$ (maximum de vraisemblance) et $\hat{\theta}_2(N) = 1/2$ (estimateur constant). Leurs fonctions de risque respectives sont $R(\theta, \hat{\theta}_1) = n^{-1}\theta(1-\theta)$ et

$R(\theta, \hat{\theta}_2) = (\theta - 1/2)^2$. Elles ne sont pas comparables et, si n est petit, la seconde sera plus petite que la première dans un voisinage non négligeable de $1/2$. Cela pose clairement problème si l'on pense que la vraie valeur de θ est proche de $1/2$.

5.1.1 Loi a priori et risque de Bayes

Une manière de régler le problème de la comparaison des estimateurs est alors de considérer un risque “moyen”. On se donne une probabilité μ sur Θ , dite *loi a priori* et l'on intègre le risque par rapport à μ .

Définition 17 *Le risque de Bayes de l'estimateur $\hat{\theta}$ pour la loi a priori μ sur Θ est donné par*

$$R_B(\mu, \hat{\theta}) = \int_{\Theta} R(\theta, \hat{\theta}) d\mu(\theta) = \int_{\Theta} d\mu(\theta) \int_E [\theta - \hat{\theta}(x)]^2 dP_{\theta}(x). \quad (\text{III.5.1})$$

On peut soit voir ce risque comme une moyenne, soit interpréter cela en disant que le paramètre est lui-même une variable aléatoire θ de loi μ et que P_{θ} est la loi conditionnelle de \mathbf{X} sachant $\theta = \theta$:

$$P_{\theta}[A] = \mathbb{E}_Q[\mathbb{1}_A(\mathbf{X}) | \theta = \theta],$$

l'espérance étant prise par rapport à la loi jointe Q de θ et \mathbf{X} définie par

$$Q(B \times C) = \mathbb{P}[\theta \in B \text{ et } \mathbf{X} \in C] = \int_B P_{\theta}[C] d\mu(\theta).$$

Dans ce cadre, (III.5.1) s'écrit $R_B(\mu, \hat{\theta}) = \mathbb{E}_Q[(\theta - \hat{\theta}(\mathbf{X}))^2]$.

5.1.2 Loi a posteriori

Dans tous les exemples que nous considérerons, les lois auront des densités, ce qui facilite les calculs. Supposons que l'on dispose d'une mesure de référence λ sur Θ , typiquement la mesure de Lebesgue, de sorte que l'on puisse écrire $\mu = \bar{g} \cdot \lambda$ et que l'on dispose de même d'une mesure de référence ν sur E pour laquelle $P_{\theta} = f(\cdot | \theta) \cdot \nu$ ou $dP_{\theta}/d\nu(x) = f(x | \theta)$. On vérifie alors que Q a la densité $\bar{g}(\theta)f(x | \theta)$ par rapport à la mesure produit $\lambda \otimes \nu$ sur $\Theta \times E$. Dans ce cas, $f(x | \theta)$ est la densité (de la loi) conditionnelle de \mathbf{X} sachant que $\theta = \theta$, ce qui justifie la notation que nous avons adoptée.

Rappel *Étant donné deux variables aléatoires $Y \in \mathcal{Y}$ et $Z \in \mathcal{Z}$ ayant une densité jointe par rapport à une mesure produit sur $\mathcal{Y} \times \mathcal{Z}$, la densité conditionnelle $f(y | z)$ de Y sachant Z est le quotient de la densité jointe par la densité marginale de Z .*

Ici la densité jointe est $\bar{g}(\theta)f(x | \theta)$ et la densité marginale de θ est \bar{g} par définition. Mais on peut également calculer la densité conditionnelle de θ sachant que $\mathbf{X} = x$. La densité marginale \bar{f} de \mathbf{X} s'obtient par intégration en θ , d'où

$$\bar{f}(x) = \int_{\Theta} f(x | \theta) \bar{g}(\theta) d\lambda(\theta) \quad \text{et} \quad g(\theta | x) = \bar{g}(\theta) f(x | \theta) / \bar{f}(x), \quad (\text{III.5.2})$$

où $g(\theta | x)$ désigne la densité conditionnelle de θ sachant que $\mathbf{X} = x$, c'est à dire la densité par rapport à λ de la loi conditionnelle Q_x de θ sachant que $\mathbf{X} = x$.

Définition 18 *La loi conditionnelle Q_x de θ sachant que $\mathbf{X} = x$ de densité $g(\theta | x)$ par rapport à λ donnée par (III.5.2) est dite loi a posteriori de θ sachant que $\mathbf{X} = x$.*

5.1.3 Définition des estimateurs de Bayes

Si l'on prend le risque de Bayes comme référence pour mesurer la qualité d'un estimateur $\hat{\theta}$ à valeur dans Θ , il est clair qu'un estimateur optimal est un estimateur qui minimise le risque de Bayes. Le Théorème de Fubini et le fait que $\bar{g}(\theta)f(x|\theta) = g(\theta|x)\bar{f}(x)$ permettent alors de récrire le risque de Bayes (III.5.1) d'un estimateur $\hat{\theta}$ de la manière suivante :

$$\begin{aligned} R_B(\mu, \hat{\theta}) &= \int_{\Theta} \bar{g}(\theta) d\lambda(\theta) \int_E [\theta - \hat{\theta}(x)]^2 f(x|\theta) d\nu(x) \\ &= \int_{\Theta} \int_E [\theta - \hat{\theta}(x)]^2 g(\theta|x) \bar{f}(x) d\nu(x) d\lambda(\theta) \\ &= \int_E \bar{f}(x) d\nu(x) \int_{\Theta} [\theta - \hat{\theta}(x)]^2 g(\theta|x) d\lambda(\theta) \\ &= \int_E \bar{f}(x) d\nu(x) \int_{\Theta} [\theta - \hat{\theta}(x)]^2 dQ_x(\theta). \end{aligned}$$

Pour minimiser cette quantité, il suffit de minimiser, pour tout x fixé,

$$\int_{\Theta} [\theta - \hat{\theta}(x)]^2 dQ_x(\theta) = \mathbb{E}_{Q_x}[(\theta - \hat{\theta}(x))^2] = \mathbb{E}_Q[(\theta - \hat{\theta}(x))^2 | \mathbf{X} = x],$$

c'est à dire l'espérance de $[\theta - \hat{\theta}(x)]^2$ pour la loi a posteriori de θ . Or on sait bien que le minimum en t de $\mathbb{E}_{Q_x}[(\theta - t)^2]$ est obtenu pour $t = \mathbb{E}_{Q_x}[\theta]$. On obtient ainsi le résultat suivant.

Proposition 5 *L'estimateur $\tilde{\theta}(\mathbf{X})$ défini par $\tilde{\theta}(x) = \mathbb{E}_{Q_x}[\theta]$, où Q_x désigne la loi a posteriori de θ quand $\mathbf{X} = x$, minimise le risque de Bayes pour la loi a priori μ . Il est appelé estimateur de Bayes pour la loi a priori μ .*

5.1.4 Calcul des estimateurs de Bayes

Pour calculer l'estimateur de Bayes, tout revient donc à évaluer la loi a posteriori de θ . Pour ce faire, une remarque sera souvent utile. On veut calculer, pour x fixé, $g(\theta|x)$ qui est une densité par rapport à λ . Or on sait par (III.5.2) que $g(\theta|x) = \bar{g}(\theta)f(x|\theta)/\bar{f}(x)$. Si l'on peut montrer que la densité jointe $\bar{g}(\theta)f(x|\theta)$ est proportionnelle à une certaine fonction $k_x(\theta)$, c'est à dire que $\bar{g}(\theta)f(x|\theta) = \alpha(x)k_x(\theta)$ où $k_x(\theta)$, considérée comme fonction de θ , est une densité de probabilité connue par rapport à λ dont les paramètres dépendent de x , alors nécessairement, puisque $g(\theta|x)$ est aussi une densité en θ , $g(\theta|x) = k_x(\theta)$, comme on le voit facilement en intégrant en θ . On procèdera le plus souvent ainsi, en montrant que $f(x|\theta)\bar{g}(\theta)$ est proportionnelle à une densité classique avec une constante de proportionnalité qui ne dépend que de x pour trouver la loi a posteriori de θ .

Exemple : X_1, \dots, X_n sont, conditionnellement à $\theta = \tau$, des v.a.r. i.i.d. $\mathcal{N}(\tau, 1)$ et θ est $\mathcal{N}(\tau, \sigma^2)$. La densité jointe des observations s'écrit donc

$$(2\pi)^{-n/2} (2\pi\sigma^2)^{-1/2} \exp \left[-\frac{1}{2} \sum_{i=1}^n (x_i - \theta)^2 - \frac{(\theta - \tau)^2}{2\sigma^2} \right].$$

Si l'on factorise en séparant ce qui dépend de θ et ce qui n'en dépend pas, on peut récrire cela, après développement des carrés, sous la forme

$$K(x_1, \dots, x_n) \exp \left[-\frac{1}{2\sigma^2} \left(n\sigma^2\theta^2 + \theta^2 - 2\theta\sigma^2 \sum_{i=1}^n x_i - 2\theta\tau \right) \right]$$

ou encore en complétant le carré,

$$K'(x_1, \dots, x_n) \exp \left[-\frac{n\sigma^2 + 1}{2\sigma^2} \left(\theta - \frac{\sigma^2 \sum_{i=1}^n x_i + \tau}{n\sigma^2 + 1} \right)^2 \right].$$

On reconnaît là une fonction de θ proportionnelle à la densité gaussienne de moyenne $(\sigma^2 \sum_{i=1}^n x_i + \tau) / (n\sigma^2 + 1)$ et de variance $\sigma^2 / (n\sigma^2 + 1)$, ce qui implique que la loi a posteriori de θ sachant $\mathbf{X} = (X_1, \dots, X_n)$ est la loi gaussienne

$$\mathcal{N} \left([n\sigma^2 \bar{X}_n + \tau] [n\sigma^2 + 1]^{-1}, \sigma^2 [n\sigma^2 + 1]^{-1} \right).$$

L'estimateur de Bayes s'écrit donc

$$\tilde{\theta}_n = [\bar{X}_n + \tau (n\sigma^2)^{-1}] / [1 + (n\sigma^2)^{-1}].$$

Asymptotiquement (lorsque $n \rightarrow +\infty$), on retrouve la moyenne empirique, mais lorsque n est petit, les deux estimateurs peuvent être très différents selon les valeurs des deux paramètres de la loi a priori. Si celle-ci est centrée et de très grande variance, on retrouve à nouveau un estimateur proche de la moyenne empirique. On peut aussi remarquer que cet estimateur est toujours biaisé. Sa fonction de risque s'écrit

$$\begin{aligned} R(\theta, \tilde{\theta}_n) &= \mathbb{E}_\theta \left[\left(\frac{\bar{X}_n + \tau (n\sigma^2)^{-1}}{1 + (n\sigma^2)^{-1}} - \theta \right)^2 \right] \\ &= \frac{\text{Var}_\theta(\bar{X}_n)}{(1 + (n\sigma^2)^{-1})^2} + \left(\frac{\theta + \tau (n\sigma^2)^{-1}}{1 + (n\sigma^2)^{-1}} - \theta \right)^2 \\ &= \frac{1}{(1 + (n\sigma^2)^{-1})^2} \left[\frac{1}{n} + \left(\frac{\tau - \theta}{n\sigma^2} \right)^2 \right] \end{aligned}$$

et le risque de Bayes s'écrit alors, puisque $\theta - \tau \sim \mathcal{N}(0, \sigma^2)$,

$$\begin{aligned} R_B(\mu, \tilde{\theta}_n) &= \frac{1}{(1 + (n\sigma^2)^{-1})^2} \mathbb{E} \left[\frac{1}{n} + \left(\frac{\theta - \tau}{n\sigma^2} \right)^2 \right] \\ &= \frac{1}{(1 + (n\sigma^2)^{-1})^2} \frac{n\sigma^4 + \sigma^2}{n^2\sigma^4} = \frac{1}{(1 + (n\sigma^2)^{-1})^2} \frac{1}{n}. \end{aligned}$$

Comme prévu, le risque de Bayes de l'estimateur de Bayes $\tilde{\theta}_n$ est plus petit que celui de \bar{X}_n qui est égal à n^{-1} .

On pourra traiter de la même manière à titre d'exercices le cas de n observations X_1, \dots, X_n i.i.d. de loi (conditionnelle par rapport à θ) de Poisson $\mathcal{P}(\theta)$, $\theta > 0$, en mettant sur θ une loi a priori $\Gamma(t, \beta)$ ou le cas de \mathbf{X} de loi conditionnelle binomiale $\mathcal{B}(n, \theta)$ avec θ de loi a priori uniforme sur $[0; 1]$.

Quelques remarques utiles

- Les estimateurs de Bayes sont généralement biaisés.
- Pour obtenir de bons estimateurs bayésiens, il convient de prendre une loi a priori qui charge bien tout Θ , c'est à dire tous les points si Θ est dénombrable et tous les ouverts non vides si Θ est une partie d'espace euclidien. Dans le cas i.i.d. et sous des hypothèses convenables de régularité, analogues à celles que l'on utilise pour les estimateurs du maximum de vraisemblance, les estimateurs de Bayes sont convergents et asymptotiquement normaux.
- Il est plus difficile d'utiliser les estimateurs de Bayes dans un cadre cas non-paramétrique parce qu'il est compliqué de construire de "bonnes" lois a priori dans ce cadre.