
 CONCENTRATION ET INÉGALITÉS DE MARGE

OLIVIER CATONI

4 avril 2013

 1. DÉVIATIONS DES SOMMES DE VARIABLES ALÉATOIRES
 INDÉPENDANTES

Soit X_i , $1 \leq i \leq n$ un échantillon de variables indépendantes et

$$M \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n X_i$$

leur moyenne empirique. On se pose la question des déviations de M par rapport à sa moyenne

$$m \stackrel{\text{def}}{=} \mathbb{E}(M) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i).$$

Considérons les fonctions génératrices des moments

$$\begin{aligned} \psi_i(\lambda) &= \log \left\{ \mathbb{E}[\exp(\lambda X_i)] \right\}, \\ \psi(\lambda) &= \frac{1}{n} \sum_{i=1}^n \psi_i(\lambda). \end{aligned}$$

Ce sont les fonctions convexes, nulles en zéro, à valeurs dans $\mathbb{R}_+ \cup \{+\infty\}$. Considérons la fonction duale

$$\psi^*(x) = \sup_{\lambda \in \mathbb{R}_+} \lambda x - \psi(\lambda) \in \mathbb{R}_+ \cup \{+\infty\}.$$

PROPOSITION 1.1 *Les déviations de la moyennes empirique M vérifient*

$$\mathbb{P}(M \geq x) \leq \exp[-n\psi^*(x)].$$

PREUVE.

CNRS – UMR 8553, Département de Mathématiques et Applications, Ecole Normale Supérieure, 45, rue d’Ulm, F75230 Paris cedex 05, and INRIA Paris-Rocquencourt – CLAS-SIC team.

$$\begin{aligned}
\mathbb{P}(M \geq x) &= \mathbb{E}\left\{\mathbf{1}[\exp(n\lambda(M-x)) \geq 1]\right\} \\
&\leq \mathbb{E}[\exp(n\lambda(M-x))] \\
&= \exp\{n[\psi(\lambda) - \lambda x]\}, \quad \lambda \in \mathbb{R}_+.
\end{aligned}$$

Par conséquent

$$\mathbb{P}(M \geq x) \leq \inf_{\lambda \in \mathbb{R}_+} \exp\{n[\psi(\lambda) - \lambda x]\} = \exp(-n\psi^*(x)).$$

□

PROPOSITION 1.2 *Posons $\Lambda_i = \sup\{\lambda \in \mathbb{R}_+ : \psi_i(\lambda) < +\infty\}$ et $\Lambda = \min\{\Lambda_1, \dots, \Lambda_n\}$. Pour tout $\lambda \in [0, \Lambda_i[$, $\psi_i(\lambda) < +\infty$ et la fonction ψ_i est de classe \mathcal{C}^∞ sur $]0, \Lambda_i[$ lorsque $\Lambda_i > 0$. Si de plus $\mathbb{E}(|X_i|^k) < \infty$, la fonction ψ_i est de classe \mathcal{C}^k sur $[0, \Lambda_i[$.*

PREUVE. Posons $\varphi(\lambda) = \mathbb{E}[\exp(\lambda X_i)]$. Soit $\lambda \in [0, \Lambda_i[$. Par définition de Λ_i , il existe $\beta \in]\lambda, \Lambda_i[$ tel que $\psi_i(\beta) < \infty$, et donc $\varphi(\beta) < \infty$. D'après l'inégalité de Jensen

$$+\infty > \mathbb{E}[\exp(\beta X_i)] = \mathbb{E}\left\{[\exp(\lambda X_i)]^{\beta/\lambda}\right\} \geq \left\{\mathbb{E}[\exp(\lambda X_i)]\right\}^{\beta/\lambda},$$

ce qui prouve que $\varphi_i(\lambda) < \infty$, et donc que $\psi_i(\lambda) < \infty$.

Remarquons que

$$\begin{aligned}
X_i^{j-1} \exp(\beta X_i) &= X_i^{j-1} \exp(\alpha X_i) + \int_{\alpha}^{\beta} X_i^j \exp(\lambda X_i) d\lambda, \\
& \qquad \qquad \qquad 0 < \alpha < \beta < \Lambda_i, \quad j \geq 1.
\end{aligned}$$

De plus

$$\mathbb{E}\left\{\sup_{\lambda \in [\alpha, \beta]} |X_i^j \exp(\lambda X_i)|\right\} < \infty$$

En effet, considérons $\gamma \in]\beta, \Lambda_i[$ et

$$C_1 = \sup_{x \in \mathbb{R}_+} x^j \exp[-(\gamma - \beta)x],$$

$$C_2 = \sup_{x \in \mathbb{R}_+} x^j \exp(-\alpha x).$$

$$|X_i^j \exp(\lambda X_i)| \leq \begin{cases} C_1 \exp(\gamma X_i), & X_i \geq 0, \\ C_2, & X_i \leq 0. \end{cases}$$

Par conséquent

$$\mathbb{E}\left\{\sup_{\lambda \in [\alpha, \beta]} |X_i^j \exp(\lambda X_i)|\right\} \leq C_1 \mathbb{E}[\exp(\gamma X_i)] + C_2 < \infty.$$

On peut donc employer le théorème de Fubini et écrire

$$\begin{aligned} \mathbb{E}[X_i^{j-1} \exp(\beta X)] &= \mathbb{E}[X_i^{j-1} \exp(\alpha X_i)] + \mathbb{E}\left(\int_{\alpha}^{\beta} X_i^j \exp(\lambda X_i) d\lambda\right) \\ &= \mathbb{E}[X_i^{j-1} \exp(\alpha X_i)] + \int_{\alpha}^{\beta} \mathbb{E}[X_i^j \exp(\lambda X_i)] d\lambda. \end{aligned}$$

D'après le théorème de convergence dominée, $\lambda \mapsto \mathbb{E}(X_i^j \exp(\lambda X_i)) : [\alpha, \beta] \rightarrow \mathbb{R}$ est continue, donc $\beta \mapsto \mathbb{E}[X_i^{j-1} \exp(\beta X_i)]$ est de classe \mathcal{C}^1 , de dérivée $\mathbb{E}[X_i^j \exp(\beta X_i)]$, donc $\beta \mapsto \mathbb{E}[\exp(\beta X_i)]$ est de classe \mathcal{C}^∞ sur $]0, \Lambda_i[$ et il en est de même de ψ_i .

Supposons maintenant de plus que $\mathbb{E}[|X_i|^k] < \infty$. On peut dans ce cas montrer de façon analogue que

$$\mathbb{E}\left\{\sup_{\lambda \in [0, \beta]} |X_i^j \exp(\lambda X_i)|\right\} < \infty, \quad 0 < \beta < \Lambda_i,$$

en déduire que

$$\mathbb{E}[X_i^{j-1} \exp(\beta X_i)] = \mathbb{E}(X_i^{j-1}) + \int_0^{\beta} \mathbb{E}[X_i^j \exp(\lambda X_i)] d\lambda, \quad 0 \leq \beta < \Lambda_i, \quad 1 \leq j \leq k,$$

et enfin que φ_i et ψ_i sont de classe \mathcal{C}^k sur $[0, \Lambda_i[$. \square

PROPOSITION 1.3 *Supposons que $\mathbb{E}(X_i^2) < \infty$ et que $\Lambda_i > 0$. La dérivée seconde de ψ_i prend la forme d'une variance :*

$$\psi_i''(\lambda) = \frac{\mathbb{E}[X_i^2 \exp(\lambda X_i)]}{\mathbb{E}[\exp(\lambda X_i)]} - \left(\frac{\mathbb{E}[X_i \exp(\lambda X_i)]}{\mathbb{E}[\exp(\lambda X_i)]}\right)^2, \quad 0 \leq \lambda < \Lambda_i,$$

de plus

$$\psi_i(\lambda) = \lambda \mathbb{E}(X_i) + \int_0^{\lambda} (\lambda - \alpha) \psi_i''(\alpha) d\alpha, \quad 0 \leq \lambda < \Lambda_i.$$

PREUVE. D'après la proposition précédente, ψ_i est de classe \mathcal{C}^2 sur $[0, \Lambda_i[$ et

$$\psi_i'(\lambda) = \frac{\mathbb{E}[X_i \exp(\lambda X_i)]}{\mathbb{E}[\exp(\lambda X_i)]},$$

en dérivant une fois de plus, on obtient l'expression de ψ_i'' indiquée dans la proposition. Considérons la variable aléatoire Y_i de loi

$$\mathbb{P}(Y_i \in A) = \frac{\mathbb{E}[\mathbf{1}(X_i \in A) \exp(\lambda X_i)]}{\mathbb{E}[\exp(\lambda X_i)]},$$

pour tout borélien A . Elle vérifie pour toute fonction mesurable f telle que $\mathbb{E}[|f(X_i)| \exp(\lambda X_i)] < \infty$

$$\mathbb{E}[f(Y_i)] = \frac{\mathbb{E}[f(X_i) \exp(\lambda X_i)]}{\mathbb{E}[\exp(\lambda X_i)]},$$

ce qui montre que $\psi_i''(\lambda) = \mathbb{E}(Y_i^2) - \mathbb{E}(Y_i)^2$ est bien une variance. La deuxième partie de la proposition s'obtient en écrivant la formule de Taylor $\psi_i(\lambda) = \psi_i(0) + \lambda \psi_i'(0) + \int_0^\lambda (\lambda - \alpha) \psi_i''(\alpha) d\alpha$. \square

PROPOSITION 1.4 *Supposons que $\Lambda > 0$ et que $\mathbb{E}(X_i^2) < \infty$, $1 \leq i \leq n$. Posons*

$$\bar{V}(\lambda) \stackrel{\text{def}}{=} \frac{2}{\lambda^2} [\psi(\lambda) - \lambda m] = \frac{2}{\lambda^2} \int_0^\lambda (\lambda - \alpha) \psi''(\alpha) d\alpha, \quad 0 \leq \lambda < \Lambda$$

$$V(\lambda) \stackrel{\text{def}}{=} \sup_{\beta \in [0, \lambda]} \bar{V}(\beta),$$

$$v \stackrel{\text{def}}{=} V(0) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}\{[X_i - \mathbb{E}(X_i)]^2\}$$

et remarquons que V est une fonction continue et croissante à valeurs dans $\mathbb{R} + \cup\{+\infty\}$. Sous ces hypothèses

$$\begin{aligned} \mathbb{P}(M \geq m + x) &\leq \exp\left(-\frac{nx^2}{2V(x/v)}\right), \\ \mathbb{P}\left(M \geq m + \sqrt{\frac{2 \log(\epsilon^{-1})}{n}} V\left(\sqrt{\frac{2 \log(\epsilon^{-1})}{nv}}\right)\right) &\leq \epsilon. \end{aligned}$$

PREUVE. Pour tout $0 \leq \beta \leq \lambda$,

$$\psi^*(m + x) \geq \beta x - \frac{\beta^2}{2} V(\lambda),$$

si bien que

$$\mathbb{P}(M \geq m + x) \leq \exp\left[-n\left(\beta x - \frac{\beta^2}{2} V(\lambda)\right)\right].$$

On obtient la première inégalité en choisissant $\lambda = x/v$ et $\beta = x/V(\lambda) \leq \lambda$. Pour obtenir la seconde posons $\epsilon = \exp\left[-n\left(\beta x - \frac{\beta^2}{2}V(\lambda)\right)\right]$, pour obtenir dans un premier temps

$$\mathbb{P}\left(M \geq m + \frac{\beta}{2}V(\lambda) + \frac{\log(\epsilon^{-1})}{n\beta}\right) \leq \epsilon.$$

Choisissons alors $\lambda = \sqrt{\frac{2\log(\epsilon^{-1})}{nv}} \geq \beta = \sqrt{\frac{2\log(\epsilon^{-1})}{nV(\lambda)}}$ pour conclure. \square

PROPOSITION 1.5 (INÉGALITÉ DE BENNETT) *Supposons que $\mathbb{E}(X_i^2) < \infty$ et que $X_i \leq \mathbb{E}(X_i) + b$, $1 \leq i \leq n$. Introduisons la fonction*

$$h(u) = (1+u)\log(1+u) - u \geq \frac{u^2}{2(1+u/3)}, \quad u \in \mathbb{R}_+.$$

Sous ces hypothèses

$$\begin{aligned} \mathbb{P}(M \geq m+x) &\leq \exp\left[-\frac{nv}{b^2}h\left(\frac{bx}{v}\right)\right] \leq \exp\left(-\frac{nx^2}{2v + \frac{2bx}{3}}\right), \\ \mathbb{P}\left(M \geq m + \sqrt{\frac{2v\log(\epsilon^{-1})}{n}}\left(1 - \frac{b}{3v}\sqrt{\frac{2v\log(\epsilon^{-1})}{n}}\right)^{-1/2}\right) &\leq \epsilon. \end{aligned}$$

PREUVE. Remarquons tout d'abord que pour tout $\lambda \in \mathbb{R}_+$,

$$\begin{aligned} \psi^*(m+x) &\geq \lambda(x+m) - \frac{1}{n} \sum_{i=1}^n \log[\mathbb{E}(\exp(\lambda X_i))] \\ &= \lambda x - \frac{1}{n} \sum_{i=1}^n \log\left\{\mathbb{E}[\exp(\lambda(X_i - m_i))]\right\}, \end{aligned}$$

où $m_i \stackrel{\text{def}}{=} \mathbb{E}(X_i)$. On peut alors écrire

$$\begin{aligned} \mathbb{E}[\exp(\lambda(X_i - m_i))] - 1 &= \mathbb{E}[\exp(\lambda(X_i - m_i)) - 1 - \lambda(X_i - m_i)] \\ &= \mathbb{E}[\lambda^2(X_i - m_i)^2 g(\lambda(X_i - m_i))] \end{aligned}$$

où $g(y) = y^{-2}(\exp(y) - 1 - y)$. La fonction g est croissante sur \mathbb{R} . En effectuant un développement de Taylor à l'ordre deux de la fonction $z \mapsto \exp(yz)$ entre 0 et 1, on voit en effet qu'elle peut s'écrire

$$g(y) = \int_0^1 (1-z)\exp(yz) dz, \quad y \in \mathbb{R}.$$

On en déduit que

$$\mathbb{E}[\lambda^2(X_i - m_i)^2 g(\lambda(X_i - m_i))] \leq \mathbb{E}[\lambda^2(X_i - m_i)^2 g(\lambda b)], \quad 1 \leq i \leq n,$$

et donc que

$$\log \left\{ \mathbb{E}[\exp(\lambda(X_i - m_i))] \right\} \leq \lambda^2 g(\lambda b) \mathbb{E}[(X_i - m_i)^2].$$

Ainsi

$$\psi^*(m+x) \geq \lambda x - \lambda^2 v g(\lambda b) = \lambda x - \frac{v}{b^2} (\exp(\lambda b) - 1 - \lambda b).$$

Choisissons $\lambda = b^{-1} \log \left(1 + \frac{bx}{v} \right)$ pour obtenir

$$\psi^*(x) \geq \frac{v}{b^2} h \left(\frac{bx}{v} \right).$$

Montrons maintenant que $h(u) \geq \frac{u^2}{2(1+u/3)}$, $u > -1$. Calculons les dérivées de h , $h'(u) = \log(1+u)$, $h''(u) = 1/(1+u)$, puis les dérivées de $f(u) = (1+u/3)h(u) - u^2/2$. On obtient $f'(u) = h'(u)(1+u/3) + h(u)/3 - u$ qui vérifie $f'(0) = 0$ et

$$\begin{aligned} f''(u) &= h''(u)(1+u/3) + 2h'(u)/3 - 1 = \frac{1+u/3}{1+u} + \frac{2}{3} \log(1+u) - 1 \\ &= \frac{2}{3} \log(1+u) - \frac{2u}{3(1+u)} = \frac{2h(u)}{3(1+u)} \geq 0, \quad u > -1. \end{aligned}$$

La fonction f , convexe, nulle en zéro et de dérivée première nulle en zéro est donc positive.

Posons $\epsilon = \exp \left(-\frac{nx^2}{2v + \frac{2bx}{3}} \right)$. On obtient

$$\begin{aligned} x^2 &= \frac{2v \log(\epsilon^{-1})}{n} \left(1 + \frac{bx^2}{3vx} \right) \\ &\leq \frac{2v \log(\epsilon^{-1})}{n} \left(1 + \frac{bx^2}{3v} \left(\frac{2v \log(\epsilon^{-1})}{n} \right)^{-1/2} \right). \end{aligned}$$

On en déduit que

$$x^2 \leq \frac{2v \log(\epsilon^{-1})}{n} \left(1 - \frac{b}{3v} \sqrt{\frac{2v \log(\epsilon^{-1})}{n}} \right)^{-1},$$

ce qui prouve la deuxième inégalité de la proposition. \square

PROPOSITION 1.6 (INÉGALITÉ DE Hoeffding) *Supposons que $a_i \leq X_i \leq b_i$, $1 \leq i \leq n$. Dans ce cas*

$$\mathbb{P}(M \geq m + x) \leq \exp\left(-\frac{2n^2 x^2}{\sum_{i=1}^n (b_i - a_i)^2}\right),$$

$$\mathbb{P}\left(M \geq m + \sqrt{\frac{\sum_{i=1}^n (b_i - a_i)^2 \log(\epsilon^{-1})}{2n^2}}\right) \leq \epsilon.$$

PREUVE. La dérivée seconde de ψ_i est la variance d'une variable aléatoire à valeurs dans l'intervalle $[a_i, b_i]$, et ne peut donc excéder $(b_i - a_i)^2/4$. Par conséquent $\psi(\lambda) \leq \lambda m + \frac{\lambda^2}{8} \sum_{i=1}^n (b_i - a_i)^2$ d'où $\psi^*(m + x) \geq \frac{2nx^2}{\sum_{i=1}^n (b_i - a_i)^2}$.

□

2. BORNES PAC-BAYÉSIENNES POUR LES DÉVIATIONS UNIFORMES DES MOYENNES EMPIRIQUES PAR RAPPORT À LEURS ESPÉRANCES

Considérons n variables aléatoires indépendantes X_i , $1 \leq i \leq n$ à valeurs dans un espace mesurable \mathcal{X} , un espace mesurable de paramètres Θ et une fonction mesurable $f : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$ (qui peut être vue comme une famille de fonctions de \mathcal{X} dans \mathbb{R} indexée par θ). Supposons que

$$\mathbb{E}[f(X, \theta)^2] < +\infty, \quad \theta \in \Theta,$$

et posons

$$M(\theta) = \frac{1}{n} \sum_{i=1}^n f(X_i, \theta),$$

$$m(\theta) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[f(X_i, \theta)],$$

$$\psi_i(\lambda, \theta) = \log\left\{\mathbb{E} \exp[\lambda f(X_i, \theta)]\right\},$$

$$\psi(\lambda, \theta) = \frac{1}{n} \sum_{i=1}^n \psi_i(\lambda, \theta),$$

$$\Lambda = \sup\{\lambda : \psi(\lambda, \theta) < \infty, \theta \in \Theta\}$$

PROPOSITION 2.1 *Supposons que $\Lambda > 0$. Soit $\nu \in \mathcal{M}_+^1(\Theta)$ une mesure de référence sur l'espace des paramètres Θ . Pour tout $\lambda \in [0, \Lambda[$,*

$$\mathbb{E} \left[\exp \left(\sup_{\rho} \left\{ \int_{\Theta} n [\lambda M(\theta) - \psi(\lambda, \theta)] d\rho(\theta) - \mathcal{K}(\rho, \nu), \right. \right. \right. \\ \left. \left. \left. \rho \in \mathcal{M}_+^1(\Theta), \theta \mapsto \lambda M(\theta) - \psi(\lambda, \theta) \in \mathbb{L}^1(\rho), \mathcal{K}(\rho, \nu) < \infty \right\} \right) \right] \leq 1.$$

Par conséquent, avec probabilité au moins $1 - \epsilon$, pour toute probabilité $\rho \in \mathcal{M}_+^1(\Theta)$, telle que $\theta \mapsto \lambda M(\theta) - \psi(\lambda, \theta) \in \mathbb{L}^1(\rho)$ et $\mathcal{K}(\rho, \nu) < \infty$,

$$\int M(\theta) d\rho(\theta) \leq \frac{1}{\lambda} \int \psi(\lambda, \theta) d\rho(\theta) + \frac{\mathcal{K}(\rho, \nu) + \log(\epsilon^{-1})}{n\lambda}.$$

PREUVE. D'après l'inégalité de Jensen, lorsque ρ vérifie les hypothèses,

$$\begin{aligned} & \exp \left[\int_{\Theta} n [\lambda M(\theta) - \psi(\lambda, \theta)] d\rho(\theta) - \mathcal{K}(\rho, \nu) \right] \\ & \leq \int_{\Theta} \exp \left\{ n [\lambda M(\theta) - \psi(\lambda, \theta)] \right\} \mathbb{1} \left(\frac{d\rho}{d\nu}(\theta) > 0 \right) \left(\frac{d\rho}{d\nu}(\theta) \right)^{-1} d\rho(\theta) \\ & = \int_{\Theta} \exp \left\{ n [\lambda M(\theta) - \psi(\lambda, \theta)] \right\} \mathbb{1} \left(\frac{d\rho}{d\nu}(\theta) > 0 \right) d\nu(\theta) \\ & \leq \int_{\Theta} \exp \left\{ n [\lambda M(\theta) - \psi(\lambda, \theta)] \right\} d\nu(\theta). \end{aligned}$$

On peut ensuite appliquer le théorème de Fubini pour les fonctions positives.

$$\begin{aligned} & \mathbb{E} \left\{ \exp \left[\sup_{\rho \in \mathcal{M}_+^1(\Theta)} \int_{\Theta} n [\lambda M(\theta) - \psi(\lambda, \theta)] d\rho(\theta) - \mathcal{K}(\rho, \nu) \right] \right\} \\ & \leq \mathbb{E} \left[\int_{\Theta} \exp \left\{ n [\lambda M(\theta) - \psi(\lambda, \theta)] \right\} d\nu(\theta) \right] \\ & = \int_{\Theta} \mathbb{E} \left[\exp \left\{ n [\lambda M(\theta) - \psi(\lambda, \theta)] \right\} \right] d\nu(\theta) = 1. \end{aligned}$$

L'espérance indiquée dans la proposition ne porte pas forcément sur une fonction mesurable, mais cette fonction est majorée par une fonction mesurable d'espérance inférieure ou égale à 1 : c'est ce que montre la preuve et c'est le sens technique à donner à la proposition. La deuxième partie de la proposition s'obtient en appliquant l'inégalité de Markov. Là encore, l'événement en question n'est pas forcément mesurable : il faut comprendre qu'il contient un événement mesurable de probabilité au moins égale à $1 - \epsilon$. \square

$$\text{Posons } m_i(\theta) = \mathbb{E}[f(X_i, \theta)],$$

$$\begin{aligned}
v(\theta) &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left\{ [f(X_i, \theta) - m_i(\theta)]^2 \right\}, \\
\bar{V}(\lambda, \theta) &= \frac{2}{\lambda^2} [\psi(\lambda, \theta) - \lambda m(\theta)], \\
V(\lambda, \theta) &= \sup_{\beta \in [0, \lambda]} \bar{V}(\beta, \theta)
\end{aligned}$$

et supposons que $v \stackrel{\text{def}}{=} \sup_{\theta \in \Theta} v(\theta) < \infty$ et $V(\lambda) \stackrel{\text{def}}{=} \sup_{\theta \in \Theta} V(\lambda, \theta) < \infty$, $0 \leq \lambda < \Lambda'$.

PROPOSITION 2.2 *Sous les hypothèses précédentes, pour toute constante positive c ,*

$$\begin{aligned}
&\mathbb{E} \left(\sup \left\{ \int_{\Theta} [M(\theta) - m(\theta)] d\rho(\theta); \right. \right. \\
&\quad \left. \left. \rho \in \mathcal{M}_+^1(\Theta), \theta \mapsto M(\theta) - m(\theta) \in \mathbb{L}^1(\rho), \mathcal{K}(\rho, \nu) \leq c \right\} \right) \\
&\leq \inf_{\lambda \in [0, \Lambda']} \frac{\lambda \bar{V}(\lambda)}{2} + \frac{c}{\lambda n} \leq \sqrt{\frac{2c}{n} V \left(\sqrt{\frac{2c}{nv}} \right)}.
\end{aligned}$$

En particulier quand Θ est fini, en prenant $c = \log(|\Theta|)$, $\rho = \delta_{\theta}$ et $\nu(\theta) = |\Theta|^{-1}$, $\theta \in \Theta$, on obtient

$$\mathbb{E} \left\{ \sup_{\theta \in \Theta} [M(\theta) - m(\theta)] \right\} \leq \sqrt{\frac{2 \log(|\Theta|)}{n} V \left(\sqrt{\frac{2 \log(|\Theta|)}{nv}} \right)}.$$

PREUVE. D'après la preuve de la proposition précédente, l'argument de l'espérance à majorer est inférieur ou égal à

$$\frac{1}{n\lambda} \log \left\{ \int \exp \left[n [\lambda M(\theta) - \psi(\lambda, \theta)] \right] d\nu(\theta) \right\} + \frac{\lambda \bar{V}(\lambda)}{2} + \frac{c}{\lambda n},$$

et on conclut à l'aide de l'inégalité de Jensen. On obtient ainsi le premier majorant $\inf_{\lambda \in [0, \Lambda']} \frac{\lambda \bar{V}(\lambda)}{2} + \frac{c}{\lambda n}$ que l'on peut affaiblir en $\inf_{0 \leq \lambda \leq \beta} \frac{\lambda V(\beta)}{2} + \frac{c}{\lambda n}$.

On obtient le second majorant en choisissant $\beta = \sqrt{\frac{2c}{nv}}$ et $\lambda = \sqrt{\frac{2c}{nV(\beta)}} \leq \beta$.

□

PROPOSITION 2.3 *Sous les hypothèses précédentes, pour toute constante positive c , avec probabilité au moins $1 - \epsilon$,*

$$\begin{aligned} & \sup \left\{ \int_{\Theta} [M(\theta) - m(\theta)] d\rho(\theta); \right. \\ & \quad \left. \rho \in \mathcal{M}_+^1(\Theta), \theta \mapsto M(\theta) - m(\theta) \in \mathbb{L}^1(\Theta), \mathcal{K}(\rho, \nu) \leq c \right\} \\ & \leq \inf_{\lambda \in [0, \lambda']]} \frac{\lambda \bar{V}(\lambda)}{2} + \frac{c + \log(\epsilon^{-1})}{\lambda n} \leq \sqrt{\frac{2[c + \log(\epsilon^{-1})]}{n} V\left(\sqrt{\frac{2[c + \log(\epsilon^{-1})]}{nv}}\right)}. \end{aligned}$$

En particulier quand Θ est fini, avec probabilité au moins $1 - \epsilon$

$$\sup_{\theta \in \Theta} [M(\theta) - m(\theta)] \leq \sqrt{\frac{2 \log(|\Theta|/\epsilon)}{n} V\left(\sqrt{\frac{2 \log(|\Theta|/\epsilon)}{nv}}\right)}.$$

PREUVE. C'est une conséquence directe de la seconde partie de la proposition 2.1 (page 7) et de la majoration $\psi(\lambda, \theta) \leq \frac{\lambda^2 V(\lambda)}{2} + \lambda m(\theta)$. \square

PROPOSITION 2.4 *Supposons que $\Theta = \mathbb{B}_d = \{\theta \in \mathbb{R}^d; \|\theta\| \leq 1\}$ et qu'il existe deux constantes positives B et g telles que*

$$\begin{aligned} & \sup_{x \in \mathcal{X}} f(x, \theta) - \inf_{x \in \mathcal{X}} f(x, \theta) \leq B, & \theta \in \mathbb{B}_d, \\ & |f(x, \theta) - f(x, \theta')| \leq g \|\theta - \theta'\|, & x \in \mathcal{X}, \quad \theta, \theta' \in \mathbb{B}_d. \end{aligned}$$

Considérons le minimiseur du risque empirique

$$\hat{\theta} \in \arg \min_{\theta \in \mathbb{B}_d} M(\theta).$$

Avec probabilité au moins $1 - \epsilon$,

$$\begin{aligned} m(\hat{\theta}) \leq \inf_{\theta \in \mathbb{B}_d} m(\theta) + B \left\{ \sqrt{\frac{d}{2n} \log\left(1 + \frac{4g}{B} \sqrt{\frac{2n}{d}}\right) + \frac{\log(2/\epsilon)}{2n}} \right. \\ \left. + \sqrt{\frac{d}{8n}} + \sqrt{\frac{\log(2/\epsilon)}{2n}} \right\}. \end{aligned}$$

Sous ces hypothèses très simples, on voit que la qualité de l'estimation de $\inf_{\theta \in \Theta} m(\theta)$ par $\widehat{\theta}$ dépend de la dimension d de l'espace des paramètres, et plus précisément du rapport d/n entre cette dimension et la taille de l'échantillon.

PREUVE. Commençons par étendre le domaine de définition de f à \mathbb{R}^d en posant

$$f(x, \theta) = f(x, \theta/\|\theta\|), \quad \theta \in \mathbb{R}^d \setminus \mathbb{B}_d.$$

Soit $\delta > 0$ un paramètre réel positif dont nous choisirons par la suite la valeur et ν la mesure uniforme sur la boule $(1 + \delta)\mathbb{B}_d$ de rayon $1 + \delta$. Considérons pour tout $\theta \in \mathbb{B}_d$ la probabilité uniforme ρ_θ sur la boule $\theta + \delta\mathbb{B}_d$ centrée en θ de rayon θ . Le volume d'une boule de \mathbb{R}^d étant proportionnel à son rayon élevé à la puissance d , on voit immédiatement que

$$\mathcal{K}(\rho_\theta, \nu) = d \log\left(\frac{1 + \delta}{\delta}\right), \quad \theta \in \mathbb{B}_d.$$

D'après la proposition précédente et l'inégalité de Hoeffding, avec probabilité au moins $1 - \epsilon$, pour tout $\theta \in \mathbb{B}_d$,

$$\int m(\theta') d\rho_\theta(\theta') \leq \int M(\theta') d\rho_\theta(\theta') + B \sqrt{\frac{d \log(1 + \delta^{-1}) + \log(\epsilon^{-1})}{2n}}.$$

On en déduit, toujours avec probabilité $1 - \epsilon$,

$$m(\widehat{\theta}) \leq M(\widehat{\theta}) + 2g\delta + B \sqrt{\frac{d \log(1 + \delta^{-1}) + \log(\epsilon^{-1})}{2n}}.$$

Soit $\theta_* \in \mathbb{B}_d$, tel que $m(\theta_*) = \inf_{\theta \in \mathbb{B}_d} m(\theta)$ (qui existe car $\theta \mapsto m(\theta)$ est continue sur le compact \mathbb{B}_d). Avec probabilité $1 - \epsilon$

$$M(\theta_*) \leq m(\theta_*) + B \sqrt{\frac{\log(\epsilon^{-1})}{2n}}.$$

Par construction de l'estimateur $\widehat{\theta}$, $M(\widehat{\theta}) \leq M(\theta_*)$. On en déduit donc qu'avec probabilité au moins $1 - 2\epsilon$,

$$m(\widehat{\theta}) \leq m(\theta_*) + B \left\{ \sqrt{\frac{d \log(1 + \delta^{-1}) + \log(\epsilon^{-1})}{2n}} + \sqrt{\frac{\log(\epsilon^{-1})}{2n}} \right\} + 2g\delta.$$

On conclut en choisissant $\delta = \frac{B}{4g} \sqrt{\frac{d}{2n}}$ et en remplaçant ϵ par $\epsilon/2$. \square

PROPOSITION 2.5 *Supposons que $\Theta = \mathbb{R}^d$, qu'il existe une fonction mesurable $(x, \theta) \mapsto \nabla f(x, \theta) \in \mathbb{R}^d$ et des constantes positives g et H telles que*

$$|f(x, \theta) - f(x, \theta')| \leq g \|\theta - \theta'\|,$$

$$|f(x, \theta') - f(x, \theta) - \langle \nabla f(x, \theta), \theta' - \theta \rangle| \leq \frac{H}{2} \|\theta' - \theta\|^2, \quad x \in \mathcal{X}, \quad \theta, \theta' \in \mathbb{R}^d.$$

Soit $\theta_* \in \arg \min_{\theta \in \mathbb{B}_d} m(\theta)$. Introduisons la fonction

$$\chi(h) = \sup_{\theta \in \mathbb{B}_d} \frac{h}{2} \|\theta - \theta_*\|^2 - m(\theta) + m(\theta_*),$$

Dans ces conditions, le minimiseur empirique $\hat{\theta} \in \arg \min_{\theta \in \mathbb{B}_d} M(\theta)$ de m sur la boule unité vérifie avec probabilité au moins $1 - \epsilon$

$$\|\hat{\theta} - \theta_*\|^2 \leq \frac{8g^2}{nh^2} \left[\left(\frac{8H}{h} + 1 \right) d + 2 \log(\epsilon^{-1}) \right] + \frac{4\chi(h)}{h}$$

et $m(\hat{\theta}) - m(\theta_*) \leq \frac{4g^2}{nh} \left[\left(\frac{8H}{h} + 1 \right) d + 2 \log(\epsilon^{-1}) \right] + \chi(h)$.

Dans le cas où il existe $h > 0$ tel que $\chi(h) = 0$, on obtient donc une vitesse de convergence en d/n au lieu d'une vitesse en $\sqrt{d/n}$ sous les hypothèses plus faibles de la proposition précédente.

Exercice 1 *Dans le cas où $m(\theta) - m(\theta_*) \geq c \|\theta - \theta_*\|^\alpha$, $\theta \in \mathbb{B}_d$, avec $c > 0$ et $\alpha > 2$ quelle vitesse obtient-on ?*

PREUVE. Choisissons $\rho_\theta = \mathcal{N}(\theta, \beta^{-1}I)$ et $\nu = \rho_{\theta_*}$. Remarquons que $\mathcal{K}(\rho_\theta, \nu) = \frac{\beta}{2} \|\theta - \theta_*\|^2$. Nous allons appliquer la proposition 2.1 (page 7) à la fonction $(x, \theta) \mapsto f(x, \theta_*) - f(x, \theta)$. D'après l'inégalité de Hoeffding

$$\log \mathbb{E} \exp \left\{ \lambda [f(X, \theta_*) - f(X, \theta)] \right\} - \lambda [m(\theta) - m(\theta_*)] \leq \frac{\lambda^2 g^2 \|\theta - \theta_*\|^2}{2}$$

Avec probabilité au moins $1 - \epsilon$, pour tout $\theta \in \mathbb{B}_d$,

$$\begin{aligned} \int m(\theta') d\rho_\theta(\theta') - m(\theta_*) &\leq \int M(\theta') d\rho_\theta(\theta') - M(\theta_*) \\ &+ \frac{\lambda g^2}{2} \int \|\theta' - \theta_*\|^2 d\rho_\theta(\theta') + \frac{\beta \|\theta - \theta_*\|^2}{2n\lambda} + \frac{\log(\epsilon^{-1})}{n\lambda}. \end{aligned}$$

De plus

$$\begin{aligned}
\int m(\theta') \, d\rho_\theta(\theta') &= m(\theta) \\
&+ \mathbb{E} \left[\int \left[f(X, \theta') - f(X, \theta) - \langle \nabla f(X, \theta), \theta' - \theta \rangle \right] d\rho_\theta(\theta') \right] \\
&\geq m(\theta) - \frac{H}{2} \int \|\theta' - \theta\|^2 \, d\rho_\theta(\theta') = m(\theta) - \frac{Hd}{2\beta}.
\end{aligned}$$

De même $\int M(\theta') \, d\rho_\theta(\theta') \leq M(\theta) + \frac{Hd}{2\beta}$. On en déduit qu'avec probabilité au moins $1 - \epsilon$, pour tout $\theta \in \mathbb{B}_d$,

$$\begin{aligned}
m(\theta) - m(\theta_*) &\leq M(\theta) - M(\theta_*) + \frac{Hd}{\beta} + \frac{\lambda g^2 d}{2\beta} + \frac{\lambda g^2}{2} \|\theta - \theta_*\|^2 \\
&\quad + \frac{\beta \|\theta - \theta_*\|^2}{2n\lambda} + \frac{\log(\epsilon^{-1})}{n\lambda}.
\end{aligned}$$

On peut alors utiliser le fait que $m(\theta) - m(\theta_*) \geq \frac{h}{2} \|\theta - \theta_*\|^2 - \chi(h)$ et que par construction $M(\hat{\theta}) \leq M(\theta_*)$. On en conclut avec probabilité au moins $1 - \epsilon$

$$\begin{aligned}
\frac{h}{2} \|\hat{\theta} - \theta_*\|^2 &\leq \chi(h) + \frac{d}{\beta} \left(H + \frac{\lambda g^2}{2} \right) \\
&\quad + \left(\frac{\lambda g^2}{2} + \frac{\beta}{2n\lambda} \right) \|\hat{\theta} - \theta_*\|^2 + \frac{\log(\epsilon^{-1})}{n\lambda}.
\end{aligned}$$

Ainsi

$$\|\hat{\theta} - \theta_*\|^2 \left(1 - \frac{\lambda g^2}{h} - \frac{\beta}{n\lambda h} \right) \leq \frac{2\chi(h)}{h} + \frac{2d}{\beta h} \left(H + \frac{\lambda g^2}{2} \right) + \frac{2\log(\epsilon^{-1})}{hn\lambda}.$$

Choisissons alors $\lambda = \frac{h}{4g^2}$ et $\beta = \frac{n\lambda h}{4} = \frac{nh^2}{16g^2}$. On obtient

$$\frac{1}{2} \|\hat{\theta} - \theta_*\|^2 \leq \frac{2\chi(h)}{h} + \frac{32g^2 d}{nh^3} \left(H + \frac{h}{8} \right) + \frac{8g^2 \log(\epsilon^{-1})}{nh^2},$$

qui donne la première majoration de la proposition.

Pour prouver la seconde, on utilise $\|\hat{\theta} - \theta_*\|^2 \leq \frac{2}{h} [m(\hat{\theta}) - m(\theta_*) + \chi(h)]$, pour obtenir

$$m(\hat{\theta}) - m(\theta_*) \leq \frac{d}{\beta} \left(H + \frac{\lambda g^2}{2} \right) + \left(\frac{\lambda g^2}{2} + \frac{\beta}{2n\lambda} \right) \frac{2}{h} [m(\hat{\theta}) - m(\theta_*) + \chi(h)] + \frac{\log(\epsilon^{-1})}{n\lambda}.$$

On conclut de même en remplaçant λ et β par leurs valeurs. \square

RÉFÉRENCES

- [1] O. Catoni. *Statistical Learning Theory and Stochastic Optimization, Lectures on Probability Theory and Statistics, École d'Été de Probabilités de Saint-Flour XXXI – 2001*, volume 1851 of *Lecture Notes in Mathematics*. Springer, 2004. Pages 1–269.
- [2] O. Catoni. *PAC-Bayesian Supervised Classification : The Thermodynamics of Statistical Learning*, volume 56 of *IMS Lecture Notes Monograph Series*. Institute of Mathematical Statistics, 2007. Pages i-xii, 1-163.
- [3] T. Cover and J. Thomas. *Elements of Information Theory*. Wiley and Sons, New York, second edition, 2006.
- [4] Pascal Germain, Alexandre Lacasse, François Laviolette, and Mario Marchand. Pac-bayesian learning of linear classifiers. In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML '09*, pages 353–360, New York, NY, USA, 2009. ACM.
- [5] J. Langford and J. Shawe-Taylor. PAC-bayes & margins. In *Advances in Neural Information Processing Systems*, pages 423–430, 2002.
- [6] D. A. McAllester. PAC-Bayesian model averaging. In *Proceedings of the 12th annual conference on Computational Learning Theory*. Morgan Kaufmann, 1999.
- [7] D. A. McAllester. PAC-Bayesian stochastic model selection. *Mach. Learn.*, 51(1) :5–21, April 2003.
- [8] David Mcallester. Simplified pac-bayesian margin bounds. In *In COLT*, pages 203–215, 2003.
- [9] M. Seeger. PAC-Bayesian generalization error bounds for gaussian process classification. Informatics report series EDI-INF-RR-0094, Division of Informatics, University of Edinburgh, 2002.