
PAC-BAYES LEARNING BOUNDS

OLIVIER CATONI

*Lecture notes for the IFCAM Summer School
on Applied Mathematics
Indian Institute of Science,
Bangalore
July, 24-25, 2014*

As a prerequisite, we will recall Chernoff's deviation bound for sums of independent random variables and its consequences. We will then proceed to PAC-Bayes uniform bounds on sums of independent random variables depending on a parameter. These bounds apply to parameter estimators based on the minimization of an empirical risk function. In the third section, we will focus on supervised classification, and in the last section more specifically on Support Vector Machines.

1. DEVIATIONS OF SUMS OF INDEPENDENT RANDOM VARIABLES

Let X_i , $1 \leq i \leq n$ be a sample of independent real valued random variables, and let

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

be their empirical mean. We will study the deviations of M with respect to its mean (under hypotheses ensuring that it exists). When this is meaningful, we will set

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

Consider the moment generating functions

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

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

These are convex functions, taking their values in $\mathbb{R} \cup \{+\infty\}$, and sending zero to zero. Let us consider the dual function

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

PROPOSITION 1.1 *The deviations of the empirical mean M are such that*

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

PROOF. Let us remark that no hypothesis is needed, because m is not involved explicitly in this result and because $\psi^*(x) = 0$ is possible, in which case the result is trivial. The proof is based on the fact that for any $z \in \mathbb{R}_+$, $\mathbf{1}(z \geq 1) \leq z$. Indeed

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

Consequently,

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

□

PROPOSITION 1.2 *Let us put $\Lambda_i = \sup\{\lambda \in \mathbb{R}_+ : \psi_i(\lambda) < +\infty\}$ and $\Lambda = \min\{\Lambda_1, \dots, \Lambda_n\}$. For any $\lambda \in [0, \Lambda_i[$, $\psi_i(\lambda) < +\infty$ and the function ψ_i is of class \mathcal{C}^∞ on the interval $]0, \Lambda_i[$. If, moreover, $\mathbb{E}(|X_i|^k) < \infty$, the function ψ_i is of class \mathcal{C}^k on $[0, \Lambda_i[$.*

PROOF. Let us put $\varphi_i(\lambda) = \mathbb{E}[\exp(\lambda X_i)]$. Let us consider some $\lambda \in [0, \Lambda_i[$. By definition of Λ_i , there is $\beta \in]\lambda, \Lambda_i[$ such that $\psi_i(\beta) < \infty$, and therefore, $\varphi_i(\beta) < \infty$. From Jensen's inequality,

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

proving that $\varphi_i(\lambda) < \infty$, and therefore that $\psi_i(\lambda) < \infty$. Let us remark that

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

Moreover

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

Indeed, let us consider $\gamma \in]\beta, \Lambda_i[$ and

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

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

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

Consequently,

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

We can therefore use Fubini's theorem and write

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

From Lebesgue's dominated convergence theorem, $\lambda \mapsto \mathbb{E}(X_i^j \exp(\lambda X_i)) : [\alpha, \beta] \rightarrow \mathbb{R}$ is continuous, therefore $\beta \mapsto \mathbb{E}[X_i^{j-1} \exp(\beta X_i)]$ is of class \mathcal{C}^1 , and its derivative is $\mathbb{E}[X_i^j \exp(\beta X_i)]$, therefore $\beta \mapsto \mathbb{E}[\exp(\beta X_i)]$ is of class \mathcal{C}^∞ on $]0, \Lambda_i[$ and so is ψ_i .

Let us now assume moreover that $\mathbb{E}[|X_i|^k] < \infty$. In this case we can in the same way prove that for any $\beta \in]0, \Lambda_i[$,

$$\mathbb{E}\left\{\sup_{\lambda \in [0, \beta]} |X_i^j \exp(\lambda X_i)|\right\} \leq C_1 \mathbb{E}[\exp(\gamma X_i)] + \mathbb{E}(|X_i|^j) < \infty,$$

Consequently, for any $\beta \in [0, \Lambda_i[$,

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

so that φ_i and ψ_i are of class \mathcal{C}^k on the interval $[0, \Lambda_i[$. \square

PROPOSITION 1.3 *Let us assume that $\mathbb{E}(X_i^2) < \infty$ and that $\Lambda_i > 0$. The second derivative of ψ_i can be seen as a variance:*

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

moreover

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

PROOF. According to the previous proposition, ψ_i is of class \mathcal{C}^2 on $[0, \Lambda_i[$ and

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

Taking one more derivative, we get the expression of ψ_i'' given in the proposition. Let us consider the random variable Y_i distributed according to the probability measure

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

for any Borel set A . For any measurable function f satisfying $\mathbb{E}[|f(X_i)| \exp(\lambda X_i)] < \infty$, it is such that

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

This shows that $\psi_i''(\lambda) = \mathbb{E}(Y_i^2) - \mathbb{E}(Y_i)^2$ can indeed be expressed as a variance. The second part of the proposition is obtained using a Taylor expansion with integral reminder, namely

$$\psi_i(\lambda) = \psi_i(0) + \lambda \psi_i'(0) + \int_0^\lambda (\lambda - \alpha) \psi_i''(\alpha) \, d\alpha.$$

□

PROPOSITION 1.4 *Let us assume that $\Lambda > 0$ and that $\mathbb{E}(X_i^2) < \infty$, $1 \leq i \leq n$. Let us put*

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

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

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

let us remark that V is a continuous non decreasing function with values in $\mathbb{R} + \cup \{+\infty\}$. Under these hypotheses

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

PROOF. For any $0 \leq \beta \leq \lambda$,

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

so that

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

We get the first part of the inequality, choosing $\lambda = x/v$ and $\beta = x/V(\lambda) \leq \lambda$. To get the second part, let us put $\epsilon = \exp\left[-n\left(\beta x - \frac{\beta^2}{2} V(\lambda)\right)\right]$. We get first that

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

and from there, we conclude, choosing $\lambda = \sqrt{\frac{2 \log(\epsilon^{-1})}{nv}} \geq \beta = \sqrt{\frac{2 \log(\epsilon^{-1})}{nV(\lambda)}}$.

□

PROPOSITION 1.5 (BENNETT'S INEQUALITY) *Let us assume that $\mathbb{E}(X_i^2) < \infty$ and that $X_i \leq \mathbb{E}(X_i) + b$, $1 \leq i \leq n$. Let us introduce the function*

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

Under these hypotheses,

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

PROOF. Let us remark first that for any $\lambda \in \mathbb{R}_+$,

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

where $m_i \stackrel{\text{def}}{=} \mathbb{E}(X_i)$. We can then write

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

where $g(y) = y^{-2}(\exp(y) - 1 - y)$. The function g is non decreasing on \mathbb{R} . Using a Taylor expansion of order two of the function $z \mapsto \exp(yz)$ between 0 and 1, we see that it can indeed be written as

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

Consequently

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

Therefore,

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

Thus,

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

Let us choose $\lambda = b^{-1} \log\left(1 + \frac{bx}{v}\right)$, to get

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

Let us show now that $h(u) \geq \frac{u^2}{2(1+u/3)}$, $u > -1$. Let us compute the derivatives of h , $h'(u) = \log(1+u)$, $h''(u) = 1/(1+u)$, and then the derivatives

of $f(u) = (1 + u/3)h(u) - u^2/2$. We get $f'(u) = h'(u)(1 + u/3) + h(u)/3 - u$. Thus $f'(0) = 0$ and

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

The convex function f , sending zero to zero, with a null first derivative at zero, is therefore everywhere non negative.

Let us put $\epsilon = \exp\left(-\frac{nx^2}{2v + \frac{2bx}{3}}\right)$. We get

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

We deduce that

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

proving the second inequality of the proposition. \square

PROPOSITION 1.6 (Hoeffding's inequality) *Let us assume that $a_i \leq X_i \leq b_i$, $1 \leq i \leq n$. In this case,*

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

PROOF. The second derivative of ψ_i is the variance of a random variable taking its values in the interval $[a_i, b_i]$. It cannot therefore be larger than

$(b_i - a_i)^2/4$. Consequently, $\psi(\lambda) \leq \lambda m + \frac{\lambda^2}{8} \sum_{i=1}^n (b_i - a_i)^2$, and therefore

$$\psi^*(m + x) \geq \frac{2nx^2}{\sum_{i=1}^n (b_i - a_i)^2}. \quad \square$$

2. PAC-BAYES BOUNDS ON THE UNIFORM DEVIATIONS OF EMPIRICAL MEANS WITH RESPECT TO THEIR EXPECTATIONS

Let us consider n independent random variables X_i , $1 \leq i \leq n$ taking their values in a measurable space \mathcal{X} . Let us also consider a measurable parameter space Θ and a measurable function $f : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$ (which can be seen as a family of functions from \mathcal{X} to \mathbb{R} depending on the parameter θ). Let us assume that

$$\mathbb{E}[f(X_i, \theta)^2] < +\infty, \quad \theta \in \Theta, \quad 1 \leq i \leq n,$$

and let us put

$$\begin{aligned} M(\theta) &= \frac{1}{n} \sum_{i=1}^n f(X_i, \theta), \\ m(\theta) &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[f(X_i, \theta)], \\ \psi_i(\lambda, \theta) &= \log \left\{ \mathbb{E} \exp[\lambda f(X_i, \theta)] \right\}, \\ \psi(\lambda, \theta) &= \frac{1}{n} \sum_{i=1}^n \psi_i(\lambda, \theta), \\ \Lambda &= \sup \{ \lambda : \psi(\lambda, \theta) < \infty, \theta \in \Theta \} \end{aligned}$$

PROPOSITION 2.1 *Let us assume that $\Lambda > 0$. Let $\nu \in \mathcal{M}_+^1(\Theta)$ be a reference measure on the parameter space Θ . For any $\lambda \in [0, \Lambda[$,*

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

Consequently, with probability at least $1 - \epsilon$, for any probability measure $\rho \in \mathcal{M}_+^1(\Theta)$, such that $\theta \mapsto \lambda M(\theta) - \psi(\lambda, \theta) \in \mathbb{L}^1(\rho)$ and $\mathcal{K}(\rho, \nu) < \infty$,

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

PROOF. From Jensen's inequality, whenever ρ satisfies the hypotheses,

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

We can then apply Fubini's theorem for non negative functions.

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

The expectation used in the proposition is taken on a function that may not be measurable, but that is upper bounded by a measurable function with an expectation not greater than one, this is what the proof shows and this is also how the proposition should be understood. The second part of the proposition is a consequence of Markov's inequality. Here again, the involved event may not be a measurable set, it should be understood that it contains a measurable set of probability at least equal to $1 - \epsilon$. \square

Let us put $m_i(\theta) = \mathbb{E}[f(X_i, \theta)]$,

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

and let us assume that $v \stackrel{\text{def}}{=} \sup_{\theta \in \Theta} v(\theta) < \infty$ and $V(\lambda) \stackrel{\text{def}}{=} \sup_{\theta \in \Theta} V(\lambda, \theta) < \infty$, $0 \leq \lambda < \Lambda'$.

PROPOSITION 2.2 *Under the previous hypotheses, for any positive constant c ,*

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

In particular, when Θ is a finite set, taking $c = \log(|\Theta|)$, $\rho = \delta_{\theta}$ and $\nu(\theta) = |\Theta|^{-1}$, $\theta \in \Theta$, we get

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

PROOF. From the proof of the previous proposition, the argument of the expectation to be bounded is not greater than

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

and we conclude with the help of Jensen's inequality. We get in this way the first upper bound $\inf_{\lambda \in [0, \Lambda'[,} \frac{\lambda \bar{V}(\lambda)}{2} + \frac{c}{\lambda n}$ that we can weaken to get $\inf_{0 \leq \lambda \leq \beta} \frac{\lambda V(\beta)}{2} + \frac{c}{\lambda n}$. To get the second upper bound, we should choose $\beta = \sqrt{\frac{2c}{nv}}$ and $\lambda = \sqrt{\frac{2c}{nV(\beta)}} \leq \beta$. \square

PROPOSITION 2.3 *Under the previous hypotheses, for any positive constant c , with probability at least $1 - \epsilon$,*

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

In particular, when Θ is a finite set, with probability at least $1 - \epsilon$

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

PROOF. This is a direct consequence of the second part of Proposition 2.1 (page 8) and of the inequality $\psi(\lambda, \theta) \leq \frac{\lambda^2 V(\lambda)}{2} + \lambda m(\theta)$. \square

PROPOSITION 2.4 *Let us assume that $\Theta = \mathbb{B}_d = \{\theta \in \mathbb{R}^d; \|\theta\| \leq 1\}$ and that there exist two positive constants B and g such that*

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

Let us consider the value of the parameter where the empirical risk takes its minimum value

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

With probability at least $1 - \epsilon$,

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

Under those simple hypotheses, we see that the quality of the estimation of $\inf_{\theta \in \Theta} m(\theta)$ by $\hat{\theta}$ depends on the dimension d of the parameter space, and more precisely of the ratio d/n between this dimension and the sample size.

PROOF. Let us start by extending the domain of f to \mathbb{R}^d , putting

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

Let $\delta > 0$ be a positive real parameter to be set later and ν the uniform measure on the ball $(1 + \delta)\mathbb{B}_d$ of radius $1 + \delta$. Let us consider, for all $\theta \in \mathbb{B}_d$, the uniform probability measure ρ_θ on the ball $\theta + \delta\mathbb{B}_d$ centered at θ and of radius δ . As the volume of a ball in \mathbb{R}^d is proportional to its radius raised to the power d , we see that

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

From the previous proposition and Hoeffding's inequality, with probability at least $1 - \epsilon$, for any $\theta \in \mathbb{B}_d$,

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

We deduce, still with probability at least $1 - \epsilon$, that

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

Let $\theta_* \in \mathbb{B}_d$, such that $m(\theta_*) = \inf_{\theta \in \mathbb{B}_d} m(\theta)$ (it exists, since $\theta \mapsto m(\theta)$ is continuous on the compact set \mathbb{B}_d). With probability at least $1 - \epsilon$

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

According to the definition of the estimator $\hat{\theta}$, $M(\hat{\theta}) \leq M(\theta_*)$. Consequently, with probability at least $1 - 2\epsilon$,

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

To conclude, it is enough to choose $\delta = \frac{B}{4g}\sqrt{\frac{d}{2n}}$ and to replace ϵ with $\epsilon/2$.
□

PROPOSITION 2.5 *Let us assume that $\Theta = \mathbb{R}^d$, and that for some measurable function $(x, \theta) \mapsto \nabla f(x, \theta) \in \mathbb{R}^d$ and some positive constants g and H ,*

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

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

Let $\theta_ \in \arg \min_{\theta \in \mathbb{B}_d} m(\theta)$. Let us consider, for any $h > 0$, the function*

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

Under these hypotheses, the empirical minimizer, $\hat{\theta} \in \arg \min_{\theta \in \mathbb{B}_d} M(\theta)$ of m on the unit ball is such that with probability at least $1 - \epsilon$

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

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

In the case when there is $h > 0$ such that $\chi(h) = 0$, we thus get a convergence speed of order d/n instead of \sqrt{d}/n , under stronger hypotheses than in the previous proposition.

Exercise 1 *In the case when $m(\theta) - m(\theta_*) \geq c \|\theta - \theta_*\|^\alpha$, $\theta \in \mathbb{B}_d$, where $c > 0$ and $\alpha > 2$ what speed do we get?*

PROOF. Let us choose $\rho_\theta = \mathcal{N}(\theta, \beta^{-1}I)$ and $\nu = \rho_{\theta_*}$. Let us remark that $\mathcal{K}(\rho_\theta, \nu) = \frac{\beta}{2} \|\theta - \theta_*\|^2$. We are going to apply Proposition 2.1 (page 8) to the function $(x, \theta) \mapsto f(x, \theta_*) - f(x, \theta)$. From Hoeffding's inequality,

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

Consequently, with probability at least $1 - \epsilon$, for any $\theta \in \mathbb{B}_d$,

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

Moreover,

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

In the same way, $\int M(\theta') d\rho_\theta(\theta') \leq M(\theta) + \frac{Hd}{2\beta}$. We deduce that with probability at least $1 - \epsilon$, for any $\theta \in \mathbb{B}_d$,

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

We can then use the fact that $m(\theta) - m(\theta_*) \geq \frac{h}{2}\|\theta - \theta_*\|^2 - \chi(h)$ and that by construction $M(\widehat{\theta}) \leq M(\theta_*)$. We conclude that with probability at least $1 - \epsilon$

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

Thus

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

Let us then choose $\lambda = \frac{h}{4g^2}$ and $\beta = \frac{n\lambda h}{4} = \frac{nh^2}{16g^2}$. We get

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

This gives the first upper bound of the proposition.

To prove the second upper bound, let us use the fact that $\|\widehat{\theta} - \theta_*\|^2 \leq \frac{2}{h}[m(\widehat{\theta}) - m(\theta_*) + \chi(h)]$, to obtain

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

We conclude in the same way, replacing λ and β by their values. \square

3. PAC-BAYES BOUNDS FOR SUPERVISED CLASSIFICATION

In this section, we are given some i.i.d. sample $(W_i)_{i=1}^n \in \mathcal{W}^n$, where \mathcal{W} is a measurable space, and some binary measurable loss function $L : \mathcal{W} \times \Theta \rightarrow \{0, 1\}$, where Θ is a measurable parameter space. Our aim is to minimize with respect to $\theta \in \Theta$ the expected loss

$$\int L(w, \theta) d\mathbb{P}(w),$$

where \mathbb{P} is the marginal distribution of the observed sample $(W_i)_{i=1}^n$. More precisely, assuming that \mathbb{P} is unknown, we would like to find an estimator $\widehat{\theta}(W_{1:n})$ depending on the observed sample $W_{1:n} \stackrel{\text{def}}{=} (W_i)_{i=1}^n$ such that the excess risk

$$\int L(w, \widehat{\theta}) d\mathbb{P}(w) - \inf_{\theta \in \Theta} \int L(w, \theta) d\mathbb{P}(w)$$

is small. The previous quantity is random, since $\widehat{\theta}$ depends on the random sample $W_{1:n}$. Therefore its size can be understood in different ways. Here we will focus on the *deviations* of the excess risk. Accordingly, we will look for estimators providing a small risk with a probability close to one.

A typical example of such a problem is provided by supervised classification. In this setting $\mathcal{W} = \mathcal{X} \times \mathcal{Y}$, where \mathcal{Y} is a finite set, $W_i = (X_i, Y_i)$, where (X_i, Y_i) are input-output pairs, a family of measurable classification rules $\{f_\theta : \mathcal{X} \rightarrow \mathcal{Y}; \theta \in \Theta\}$ is considered and the loss function $L(w, \theta)$ is defined as the classification error

$$L[(x, y), \theta] = \mathbb{1}[f_\theta(x) \neq y].$$

Accordingly the aim is to minimize the expected classification error

$$\mathbb{P}_{X,Y}[f_\theta(X) \neq Y]$$

in view of a sample $(X_i, Y_i)_{i=1}^n$ of observations.

The point of view exposed here is a synthesis of the approaches of [9] and [2].

3.1. DEVIATION BOUNDS FOR SUMS OF BERNOULLI RANDOM VARIABLES. Given some parameter $\lambda \in \mathbb{R}$, let us consider the (normalized) log-Laplace transform of the Bernoulli distribution :

$$\Phi_\lambda(p) \stackrel{\text{def}}{=} -\frac{1}{\lambda} \log[1 - p + p \exp(-\lambda)].$$

Let us also consider the Kullback-Leibler divergence of two Bernoulli distributions

$$K(q, p) \stackrel{\text{def}}{=} q \log\left(\frac{q}{p}\right) + (1 - q) \log\left(\frac{1 - q}{1 - p}\right).$$

In the sequel $\overline{\mathbb{P}}$ will be the empirical measure

$$\overline{\mathbb{P}} = \frac{1}{n} \sum_{i=1}^n \delta_{W_i}$$

of an i.i.d. sample $(W_i)_{i=1}^n$ drawn from $\mathbb{P}^{\otimes n} \in \mathcal{M}_+^1(\mathcal{W}^n)$ (the set of probability measures on \mathcal{W}^n). We will use a short notation for integrals, putting for any $\rho, \pi \in \mathcal{M}_+^1(\Theta)$ and any integrable function $f \in \mathbb{L}_1(\mathcal{W} \times \Theta^2, \mathbb{P} \otimes \pi \otimes \rho)$

$$f(\mathbb{P}, \rho, \pi) = \int f(w, \theta, \theta') d\mathbb{P}(w) d\rho(\theta) d\pi(\theta'),$$

so that for instance $L(\mathbb{P}, \rho) = \int L(w, \theta) d\mathbb{P}(w) d\rho(\theta)$.

Let us recall first Chernoff's bound.

PROPOSITION 3.1 *For any fixed value of the parameter $\theta \in \Theta$, the identity*

$$\int \exp[-n\lambda L(\bar{\mathbb{P}}, \theta)] d\mathbb{P}^{\otimes n} = \exp\left\{-n\lambda\Phi_\lambda[L(\mathbb{P}, \theta)]\right\}$$

shows that with probability at least $1 - \epsilon$,

$$L(\mathbb{P}, \theta) \leq B_+[L(\bar{\mathbb{P}}, \theta), \log(\epsilon^{-1})/n],$$

$$\begin{aligned} \text{where } B_+(q, \delta) &= \inf_{\lambda \in \mathbb{R}_+} \Phi_\lambda^{-1}\left(q + \frac{\delta}{\lambda}\right) \\ &= \sup\left\{p \in [0, 1] : K(q, p) \leq \delta\right\}, \quad q \in [0, 1], \delta \in \mathbb{R}_+. \end{aligned}$$

Moreover

$$-\delta q \leq B_+(q, \delta) - q - \sqrt{2\delta q(1-q)} \leq 2\delta(1-q).$$

In the same way, the identity

$$\int \exp[n\lambda L(\bar{\mathbb{P}}, \theta)] d\mathbb{P}^{\otimes n} = \exp\left\{n\lambda\Phi_{-\lambda}[L(\mathbb{P}, \theta)]\right\}$$

shows that with probability at least $1 - \epsilon$

$$L(\bar{\mathbb{P}}, \theta) \leq B_-[L(\mathbb{P}, \theta), \log(\epsilon^{-1})/n],$$

$$\begin{aligned} \text{where } B_-(q, \delta) &= \inf_{\lambda \in \mathbb{R}_+} \Phi_{-\lambda}(q) + \frac{\delta}{\lambda} \\ &= \sup\left\{p \in [0, 1] : K(p, q) \leq \delta\right\}, \quad q \in [0, 1], \delta \in \mathbb{R}_+, \end{aligned}$$

and

$$-\delta q \leq B_-(q, \delta) - q - \sqrt{2\delta q(1-q)} \leq 2\delta(1-q).$$

Before proving this proposition, let us mention some important identity.

PROPOSITION 3.2 *For any probability measures π and ρ defined on the same measurable space, such that $\mathcal{K}(\rho, \pi) < \infty$, and any bounded measurable function h , let us define the transformed probability measure $\pi_{\exp(h)} \ll \pi$ by its density*

$$\frac{d\pi_{\exp(h)}}{d\pi} = \frac{\exp(h)}{Z},$$

where $Z = \int \exp(h) d\pi$. Let us moreover introduce the notation

$$\mathbf{Var}(h d\pi) = \int (h - \int h d\pi)^2 d\pi.$$

The expectations with respect to ρ and π of h and the log-Laplace transform of h are linked by the identities

$$\int h d\rho - \mathcal{K}(\rho, \pi) + \mathcal{K}(\rho, \pi_{\exp(h)}) = \log \left[\int \exp(h) d\pi \right] \quad (1)$$

$$= \int h d\pi + \int_0^1 (1 - \alpha) \mathbf{Var}[h d\pi_{\exp(\alpha h)}] d\alpha. \quad (2)$$

PROOF. The first identity is a straightforward consequence of the definitions of $\pi_{\exp(h)}$ and of the Kullback-Leibler divergence function. The second one is the Taylor expansion of order one with integral remainder of the function

$$f(\alpha) = \log \left[\int \exp(\alpha h) d\pi \right],$$

which says that $f(1) = f(0) + f'(0) + \int_0^1 (1 - \alpha) f''(\alpha) d\alpha$. \square

Exercise 1 *Prove that $f \in \mathcal{C}^\infty$. Hint : write*

$$h^k \exp(\alpha h) = h^k + \int_0^\alpha h^{k+1} \exp(\gamma h) d\gamma,$$

use Fubini's theorem to show that $\alpha \mapsto \int h^k \exp(\alpha h) d\pi$ belongs to \mathcal{C}^1 and compute its derivative.

Let us come now to the proof of Proposition 3.1 (page 16). Chernoff's inequality reads

$$\Phi_\lambda[L(\mathbb{P}, \theta)] - \frac{\log(\epsilon^{-1})}{n\lambda} \leq L(\bar{\mathbb{P}}, \theta),$$

where the inequality holds with probability at least $1 - \epsilon$. Since the left-hand side is non-random, it can be optimized in λ , giving

$$L(\mathbb{P}, \theta) \leq B_+ [L(\bar{\mathbb{P}}, \theta), \log(\epsilon^{-1})/n].$$

Exercise 2 Prove this statement in more details. For any integer $k > 1$, consider the event

$$A_k = \left\{ \sup_{\lambda \in \mathbb{R}_+} F(\lambda) - k^{-1} > L(\bar{\mathbb{P}}, \theta) \right\},$$

where $F(\lambda) = \Phi_\lambda[L(\mathbb{P}, \theta)] - \frac{\log(\epsilon^{-1})}{n\lambda}$. Show that $\mathbb{P}^{\otimes n}(A_k) \leq \epsilon$ by choosing some suitable value of λ . Remark that $A_k \subset A_{k+1}$ and conclude that $\mathbb{P}^{\otimes n}(\cup_k A_k) \leq \epsilon$.

Since

$$\lim_{\lambda \rightarrow +\infty} \Phi_\lambda^{-1}\left(q + \frac{\delta}{\lambda}\right) = \lim_{\lambda \rightarrow +\infty} \frac{1 - \exp(-\lambda q - \delta)}{1 - \exp(-\lambda)} \leq 1,$$

$$B_+(q, \delta) \leq 1.$$

Applying equation (1, page 17) to Bernoulli distributions gives

$$\lambda \Phi_\lambda(p) = \lambda q + K(q, p) - K(q, p_\lambda)$$

where

$$p_\lambda = \frac{p}{p + (1-p)\exp(\lambda)}.$$

This shows that

$$\begin{aligned} B_+(q, \delta) &= \sup \left\{ p \in [0, 1] : \Phi_\lambda(p) \leq q + \frac{\delta}{\lambda}, \lambda \in \mathbb{R}_+ \right\} \\ &= \sup \left\{ p \in [q, 1[: K(q, p) \leq \delta + K(q, p_\lambda), \lambda \in \mathbb{R}_+ \right\} \\ &= \sup \left\{ p \in [q, 1[: K(q, p) \leq \delta \right\} \\ &= \sup \left\{ p \in [0, 1] : K(q, p) \leq \delta \right\}, \end{aligned}$$

because when $q \leq p < 1$ then $\lambda = \log\left(\frac{q^{-1} - 1}{p^{-1} - 1}\right) \in \mathbb{R}_+$, $q = p_\lambda$ and therefore $K(q, p_\lambda) = 0$.

Let us remark now that $\frac{\partial^2}{\partial x^2} K(x, p) = x^{-1}(1-x)^{-1}$. Thus if $p \geq q \geq 1/2$, then

$$K(q, p) \geq \frac{(p-q)^2}{2q(1-q)},$$

so that if $K(q, p) \leq \delta$, then

$$p \leq q + \sqrt{2\delta q(1-q)}.$$

Now if $q \leq 1/2$ and $p \geq q$ then

$$K(q, p) \geq \left\{ \begin{array}{ll} \frac{(p-q)^2}{2p(1-p)}, & p \leq 1/2 \\ 2(p-q)^2, & p \geq 1/2 \end{array} \right\} \geq \frac{(p-q)^2}{2p(1-q)},$$

so that if $K(q, p) \leq \delta$, then

$$(p-q)^2 \leq 2\delta p(1-q),$$

implying that

$$p-q \leq \delta(1-q) + \sqrt{2\delta q(1-q) + \delta^2(1-q)^2} \leq \sqrt{2\delta q(1-q)} + 2\delta(1-q).$$

On the other hand,

$$K(q, p) \leq \frac{(p-q)^2}{2 \min\{q(1-q), p(1-p)\}} \leq \frac{(p-q)^2}{2q(1-p)},$$

thus when $K(q, p) = \delta$ with $p > q$, then

$$(p-q)^2 \geq 2\delta q(1-p),$$

implying that

$$p-q \geq -\delta q + \sqrt{2\delta q(1-q) + \delta^2 q^2} \geq \sqrt{2\delta q(1-q)} - \delta q.$$

Exercise 3 *The second part of Proposition 3.1 (page 16) is proved in the same way and left as an exercise.*

3.2. PAC-BAYES BOUNDS. We are now going to make Proposition 3.1 uniform with respect to θ . The PAC-Bayes approach to this [6, 7, 8, 4] is to randomize θ , so we will consider now joint distributions on $(W_{1:n}, \theta)$, where the distribution of $W_{1:n}$ is still $\mathbb{P}^{\otimes n}$ and the conditional distribution of θ given the sample is given by some transition probability kernel $\rho : \mathcal{W}^n \rightarrow \mathcal{M}_+^1(\Theta)$, called in this context a posterior distribution*. This posterior distribution ρ will be compared with a prior (meaning non-random) probability measure $\pi \in \mathcal{M}_+^1(\Theta)$.

*We will assume that ρ is a regular conditional probability kernel, meaning that for any measurable set A the map $(w_1, \dots, w_n) \mapsto \rho(w_1, \dots, w_n)(A)$ is assumed to be measurable. We will also assume that the σ -algebra we consider on Θ is generated by a countable family of subsets. See [1][page 50] for more details

PROPOSITION 3.3 *Let us introduce the notation*

$$B_\Lambda(q, \delta) = \inf_{\lambda \in \Lambda} \Phi_\lambda^{-1} \left(q + \frac{\delta}{\lambda} \right).$$

For any prior probability measure $\pi \in \mathcal{M}_+^1(\Theta)$ and any $\lambda \in \mathbb{R}_+$,

$$\int \exp \left[\sup_{\rho \in \mathcal{M}_+^1(\Theta)} n\lambda \left\{ \Phi_\lambda[L(\mathbb{P}, \rho)] - L(\bar{\mathbb{P}}, \rho) \right\} - \mathcal{K}(\rho, \pi) \right] d\mathbb{P}^{\otimes n} \leq 1, \quad (3)$$

and therefore for any finite set $\Lambda \subset \mathbb{R}_+$, with probability at least $1 - \epsilon$, for any $\rho \in \mathcal{M}_+^1(\Theta)$,

$$L(\mathbb{P}, \rho) \leq B_\Lambda \left(L(\bar{\mathbb{P}}, \rho), \frac{\mathcal{K}(\rho, \pi) + \log(|\Lambda|/\epsilon)}{n} \right),$$

PROOF. The exponential moment inequality (3) is a consequence of equation (1, page 17), showing that

$$\begin{aligned} \exp \left\{ \sup_{\rho \in \mathcal{M}_+^1(\Theta)} n\lambda \int \left\{ \Phi_\lambda[L(\mathbb{P}, \theta)] - L(\bar{\mathbb{P}}, \theta) \right\} d\rho(\theta) - \mathcal{K}(\rho, \pi) \right\} \\ \leq \int \exp \left[n\lambda \left\{ \Phi_\lambda[L(\mathbb{P}, \theta)] - L(\bar{\mathbb{P}}, \theta) \right\} \right] d\pi(\theta), \end{aligned}$$

and of the fact that Φ_λ is convex, showing that

$$\Phi_\lambda[L(\mathbb{P}, \rho)] \leq \int \Phi_\lambda[L(\mathbb{P}, \theta)] d\rho(\theta).$$

The deviation inequality follows as usual. \square

We cannot take the infimum on $\lambda \in \mathbb{R}_+$ as in Proposition 3.1 (page 16), because we can no more cast our deviation inequality in such a way that λ appears on some non-random side of the inequality. Nevertheless, we can get a more explicit bound from some specific choice of the set Λ .

PROPOSITION 3.4 *Let us define the least increasing upper bound of the variance of a Bernoulli distribution of parameter $p \in [0, 1]$ as*

$$\bar{v}(p) = \begin{cases} p(1-p), & p \leq 1/2, \\ 1/4, & \text{otherwise.} \end{cases}$$

Let us choose some positive integer parameter m and let us put

$$t = \frac{1}{4} \log \left(\frac{n}{8 \log[(m+1)/\epsilon]} \right).$$

With probability at least $1 - \epsilon$, for any $\rho \in \mathcal{M}_+^1(\Theta)$,

$$L(\mathbb{P}, \rho) \leq L(\bar{\mathbb{P}}, \rho) + B_m[L(\bar{\mathbb{P}}, \rho), \mathcal{K}(\rho, \pi), \epsilon],$$

where

$$\begin{aligned} B_m(q, e, \epsilon) &= \max \left\{ \sqrt{\frac{2\bar{v}(q)\{e + \log[(m+1)/\epsilon]\}}{n}} \cosh(t/m) \right. \\ &\quad \left. + \frac{2(1-q)\{e + \log[(m+1)/\epsilon]\}}{n} \cosh(t/m)^2, \right. \\ &\quad \left. \frac{2\{e + \log[(m+1)/\epsilon]\}}{n} \right\} \\ &\leq \sqrt{\frac{2\bar{v}(q)\{e + \log[(m+1)/\epsilon]\}}{n}} \cosh(t/m) \\ &\quad + \frac{2\{e + \log[(m+1)/\epsilon]\}}{n} \cosh(t/m)^2. \end{aligned}$$

Moreover, as soon as $n \geq 5$,

$$\begin{aligned} B_{\lfloor \log(n)^2 \rfloor - 1}(q, e, \epsilon) &\leq B(q, e, \epsilon) \stackrel{\text{def}}{=} \\ &\sqrt{\frac{2\bar{v}(q)\{e + \log[\log(n)^2/\epsilon]\}}{n}} \cosh[\log(n)^{-1}] \\ &\quad + \frac{2\{e + \log[\log(n)^2/\epsilon]\}}{n} \cosh[\log(n)^{-1}]^2, \quad (4) \end{aligned}$$

so that with probability at least $1 - \epsilon$, for any $\rho \in \mathcal{M}_+^1(\Theta)$,

$$\begin{aligned} L(\mathbb{P}, \rho) &\leq L(\bar{\mathbb{P}}, \rho) \\ &\quad + \sqrt{\frac{2\bar{v}[L(\bar{\mathbb{P}}, \rho)]\{\mathcal{K}(\rho, \pi) + \log[\log(n)^2/\epsilon]\}}{n}} \cosh[\log(n)^{-1}] \\ &\quad + \frac{2\{\mathcal{K}(\rho, \pi) + \log[\log(n)^2/\epsilon]\}}{n} \cosh[\log(n)^{-1}]^2. \end{aligned}$$

PROOF. Let us put

$$q = L(\bar{\mathbb{P}}, \rho),$$

$$\begin{aligned}\delta &= \frac{\mathcal{K}(\rho, \pi) + \log[(m+1)/\epsilon]}{n}, \\ \lambda_{\min} &= \sqrt{\frac{8 \log[(m+1)/\epsilon]}{n}}, \\ \Lambda &= \left\{ \lambda_{\min}^{1-k/m}, k = 0, \dots, m \right\}, \\ p &= B_{\Lambda}(q, \delta) = \inf_{\lambda \in \Lambda} \Phi_{\lambda}^{-1}\left(q + \frac{\delta}{\lambda}\right), \\ \hat{\lambda} &= \sqrt{\frac{2\delta}{\bar{v}(p)}}.\end{aligned}$$

According to equation (2, page 17) applied to Bernoulli distributions, for any $\lambda \in \Lambda$,

$$\Phi_{\lambda}(p) = p - \frac{1}{\lambda} \int_0^{\lambda} (\lambda - \alpha) p_{\alpha} (1 - p_{\alpha}) d\alpha \leq q + \frac{\delta}{\lambda}.$$

As moreover $p_{\alpha} \leq p$,

$$p - q \leq \inf_{\lambda \in \Lambda} \frac{\lambda \bar{v}(p)}{2} + \frac{\delta}{\lambda} = \inf_{\lambda \in \Lambda} \sqrt{2\delta \bar{v}(p)} \cosh \left[\log \left(\frac{\hat{\lambda}}{\lambda} \right) \right].$$

As $\bar{v}(p) \leq 1/4$ and $\delta \geq \frac{\log[(m+1)/\epsilon]}{n}$,

$$\sqrt{\frac{2\delta}{\bar{v}(p)}} = \hat{\lambda} \geq \lambda_{\min} = \sqrt{\frac{8 \log[(m+1)/\epsilon]}{n}}.$$

Therefore either $\lambda_{\min} \leq \hat{\lambda} \leq 1$, or $\hat{\lambda} > 1$. Let us consider these two cases separately.

If $\lambda_{\min} = \min \Lambda \leq \hat{\lambda} \leq \max \Lambda = 1$, then $\log(\hat{\lambda})$ is at distance at most t/m from some $\log(\lambda)$ where $\lambda \in \Lambda$, because $\log(\Lambda)$ is a grid with constant steps of size $2t/m$. Thus

$$p - q \leq \sqrt{2\delta \bar{v}(p)} \cosh(t/m).$$

If moreover $q \leq 1/2$, then $\bar{v}(p) \leq p(1-q)$, so that we obtain a quadratic inequality in p , whose solution is less than

$$p \leq q + \sqrt{2\delta q(1-q)} \cosh(t/m) + 2\delta(1-q) \cosh(t/m)^2.$$

If on the contrary $q \geq 1/2$, then $\bar{v}(p) = \bar{v}(q) = 1/4$ and

$$p \leq q + \sqrt{2\delta \bar{v}(q)} \cosh(t/m),$$

so that in both cases

$$p - q \leq \sqrt{2\delta\bar{v}(q)} \cosh(t/m) + 2\delta(1 - q) \cosh(t/m)^2. \quad (5)$$

Let us consider now the case when $\hat{\lambda} > 1$. In this case $\bar{v}(p) < 2\delta$, so that

$$p - q \leq \frac{\bar{v}(p)}{2} + \delta \leq 2\delta.$$

In conclusion, applying Proposition 3.3 (page 20) we see that with probability at least $1 - \epsilon$, for any posterior distribution ρ ,

$$L(\mathbb{P}, \rho) \leq p \leq q + \max\left\{2\delta, \sqrt{2\delta\bar{v}(q)} \cosh(t/m) + 2\delta(1 - q) \cosh(t/m)^2\right\},$$

which is precisely the statement to be proved.

In the special case when $m = \lfloor \log(n)^2 \rfloor - 1 \geq \log(n)^2 - 2$,

$$\frac{t}{m} \leq \frac{1}{4\lfloor \log(n)^2 - 2 \rfloor} \log\left(\frac{n}{8\log\lfloor \log(n)^2 - 1 \rfloor}\right) \leq \log(n)^{-1}$$

as soon as the last inequality holds, that is as soon as $n \geq \exp(\sqrt{2}) \simeq 4.11$ to make $\log(n)^2 - 2$ positive and

$$3\log(n)^2 - 8 + \log(n) \log\left\{8\log\lfloor \log(n)^2 - 1 \rfloor\right\} \geq 0,$$

which holds true for any $n \geq 5$, as can be checked numerically. \square

4. LINEAR CLASSIFICATION AND SUPPORT VECTOR MACHINES

We are going in this section to consider more specifically the case of linear binary classification. In this setting $\mathcal{W} = \mathcal{X} \times \mathcal{Y} = \mathbb{R}^d \times \{-1, +1\}$, $w = (x, y)$, where $x \in \mathbb{R}^d$ and $y \in \{-1, +1\}$, $\Theta = \mathbb{R}^d$, and

$$L(w, \theta) = \mathbf{1}[\langle \theta, x \rangle y \leq 0].$$

We are going to follow in this section the approach presented in [5] and [8].

Although we will stick in this presentation to the case when \mathcal{X} is a vector space of finite dimension, the results also apply to support vector machines [11, 10, 12], where the pattern space is some arbitrary space mapped to a Hilbert space \mathcal{H} by some implicit mapping $\Psi : \mathcal{X} \rightarrow \mathcal{H}$, $\Theta = \mathcal{H}$ and $L(w, \theta) = \mathbf{1}(\langle \theta, \Psi(x) \rangle y \leq 0)$. It turns out that classification algorithms do

not need to manipulate \mathcal{H} itself, but only to compute scalar products of the form $k(x_1, x_2) = \langle \Psi(x_1), \Psi(x_2) \rangle$, defining a symmetric positive kernel k on the original pattern space \mathcal{X} . The converse is also true, any positive symmetric kernel k can be represented as a scalar product in some mapped Hilbert space (this is the Moore-Aronszajn theorem). Often used kernels on \mathbb{R}^d are

$$\begin{aligned} k(x_1, x_2) &= (1 + \langle x_1, x_2 \rangle)^s, \text{ for which } \dim \mathcal{H} < \infty, \\ k(x_1, x_2) &= \exp(-\|x_1 - x_2\|^2), \text{ for which } \dim \mathcal{H} = +\infty. \end{aligned}$$

In the following, we will work in \mathbb{R}^d , which covers only the case when $\dim \mathcal{H} < \infty$, but extensions would be possible.

Let us consider, after [5, 8] as prior probability measure π the centered Gaussian measure with covariance $\beta^{-1} \mathbf{Id}$, so that

$$\frac{d\pi}{d\theta}(\theta) = \left(\frac{\beta}{2\pi}\right)^{d/2} \exp\left(-\frac{\beta\|\theta\|^2}{2}\right).$$

Let us also consider the function

$$\begin{aligned} \varphi(x) &= \frac{1}{\sqrt{2\pi}} \int_x^{+\infty} \exp(-t^2/2) dt, \quad x \in \mathbb{R} \\ &\leq \min\left\{\frac{1}{x\sqrt{2\pi}}, \frac{1}{2}\right\} \exp\left(-\frac{x^2}{2}\right), \quad x \in \mathbb{R}_+. \end{aligned}$$

Let π_θ be the measure π shifted by θ , defined by the identity

$$\int h(\theta') d\pi_\theta(\theta') = \int h(\theta + \theta') d\pi(\theta').$$

In this case

$$\mathcal{K}(\pi_\theta, \pi) = \frac{\beta}{2} \|\theta\|^2,$$

and

$$L(w, \pi_\theta) = \varphi[\sqrt{\beta}y\|x\|^{-1}\langle \theta, x \rangle].$$

Thus the randomized loss function has an explicit expression : randomization replaces the indicator function of the negative real line by a smooth approximation. As we are eventually interested in $L(w, \theta)$, we will shift things a little bit, considering along with the classification error function L some *error with margin*

$$M(w, \theta) = \mathbf{1}[y\|x\|^{-1}\langle \theta, x \rangle \leq 1].$$

Unlike $L(w, \theta)$ which is independent of the norm of θ , the margin error $M(w, \theta)$ depends on $\|\theta\|$, counting a classification error each time x is at

distance less than $\|x\|/\|\theta\|$ from the boundary $\{x' : \langle \theta, x' \rangle = 0\}$, so that the error with margin region is the complement of the open cone $\{x \in \mathbb{R}^d; y\langle \theta, x \rangle > \|x\|\}$.

Let us compute the randomized margin error

$$M(w, \pi_\theta) = \varphi \left\{ \sqrt{\beta} [y\|x\|^{-1} \langle \theta, x \rangle - 1] \right\}.$$

It satisfies the inequality

$$M(w, \pi_\theta) \geq \varphi(-\sqrt{\beta})L(w, \theta) = [1 - \varphi(\sqrt{\beta})]L(w, \theta). \quad (6)$$

Applying previous results we obtain

PROPOSITION 4.1 *With probability at least $1 - \epsilon$, for any $\theta \in \mathbb{R}^d$,*

$$L(\mathbb{P}, \theta) \leq [1 - \varphi(\sqrt{\beta})]^{-1} M(\mathbb{P}, \pi_\theta) \leq C_1(\theta),$$

where

$$C_1(\theta) = [1 - \varphi(\sqrt{\beta})]^{-1} B \left(M(\bar{\mathbb{P}}, \pi_\theta), \frac{\beta \|\theta\|^2}{2}, \epsilon \right),$$

the bound B being defined by equation (4, page 21).

We can now minimize this empirical upper-bound to define an estimator. Let us consider some estimator $\hat{\theta}$ such that

$$C_1(\hat{\theta}) \leq \inf_{\theta \in \mathbb{R}^d} C_1(\theta) + \zeta.$$

Then for any fixed parameter θ_* , $C_1(\hat{\theta}) \leq C_1(\theta_*) + \zeta$. On the other hand, with probability at least $1 - \epsilon$

$$M(\bar{\mathbb{P}}, \pi_{\theta_*}) \leq B_- \left(M(\mathbb{P}, \pi_{\theta_*}), \frac{\log(\epsilon^{-1})}{n} \right).$$

Indeed

$$\begin{aligned} & \int \exp \left\{ n\lambda [M(\bar{\mathbb{P}}, \pi_{\theta_*}) - \Phi_{-\lambda} [M(\mathbb{P}, \pi_{\theta_*})]] \right\} d\mathbb{P}^{\otimes n} \\ & \leq \int \exp \left\{ n\lambda \int \left\{ M(\bar{\mathbb{P}}, \theta) - \Phi_{-\lambda} [M(\mathbb{P}, \theta)] \right\} d\pi_{\theta_*}(\theta) \right\} d\mathbb{P}^{\otimes n} \leq 1, \end{aligned}$$

because $p \mapsto -\Phi_{-\lambda}(p)$ is convex. As a consequence

PROPOSITION 4.2 *With probability at least $1 - 2\epsilon$,*

$$L(\mathbb{P}, \widehat{\theta}) \leq \inf_{\theta_* \in \Theta} [1 - \varphi(\sqrt{\beta})]^{-1} B \left(B_- \left(M(\mathbb{P}, \pi_{\theta_*}), \frac{\log(\epsilon^{-1})}{n} \right), \frac{\beta \|\theta_*\|^2}{2}, \epsilon \right) + \zeta.$$

It is also possible to state a result in terms of empirical margins. Indeed

$$M(w, \pi_\theta) \leq M(w, \theta/2) + \varphi(\sqrt{\beta}).$$

Thus with probability at least $1 - \epsilon$, for any $\theta \in \mathbb{R}^d$,

$$L(\mathbb{P}, \theta) \leq C_2(\theta),$$

where

$$C_2(\theta) = [1 - \varphi(\sqrt{\beta})]^{-1} B \left(M(\overline{\mathbb{P}}, \theta/2) + \varphi(\sqrt{\beta}), \frac{\beta \|\theta\|^2}{2}, \epsilon \right).$$

However, C_1 and C_2 are non-convex criteria, and faster minimization algorithms are available for the usual SVM loss function, for which we are going to derive some generalization bounds now. Indeed, let us choose some positive radius R and let us put $\|x\|_R = \max\{R, \|x\|\}$, so that in the case when $\|x\| \leq R$, $\|x\|_R = R$.

$$M(w, \pi_\theta) = \varphi[\sqrt{\beta}(y\|x\|^{-1}\langle \theta, x \rangle - 1)] \leq (2 - y\|x\|_R^{-1}\langle \theta, x \rangle)_+ + \varphi(\sqrt{\beta}). \quad (7)$$

To check that this is true, consider the functions

$$\begin{aligned} f(z) &= \varphi[\sqrt{\beta}(\|x\|^{-1}z - 1)], \\ g(z) &= (2 - \|x\|_R^{-1}z)_+ + \varphi(\sqrt{\beta}), \quad z \in \mathbb{R}. \end{aligned}$$

Let us remark that they are both non increasing, that f is convex on the interval $z \in (\|x\|, \infty($ (because φ is convex on \mathbb{R}_+), and that $\sup f = \sup \varphi = 1$. Since $\|x\|_R \geq \|x\|$, for any $z \in]-\infty, \|x\|]$, $g(z) \geq 1 \geq f(z)$. Moreover, $g(2\|x\|_R) = \varphi(\sqrt{\beta}) \geq \varphi[\sqrt{\beta}(2\|x\|^{-1}\|x\|_R - 1)] = f(z)$. Since on the interval $(\|x\|, 2\|x\|_R)$, the function g is linear, the function f is convex and g is not smaller than f at the two ends, this proves that g is not smaller than f on the whole interval. Finally, on the interval $z \in (2\|x\|_R, +\infty($, the function g is constant and the function f is decreasing, so that on this interval also g is not smaller than f , and this ends the proof of (7), since the three intervals on which $g \geq f$ cover the whole real line.

Using the upper bounds (7) and (6, page 25), and Proposition 3.3 (page 20), we obtain

PROPOSITION 4.3 *With probability at least $1 - \epsilon$, for any $\theta \in \mathbb{R}^d$,*

$$\begin{aligned} L(\mathbb{P}, \theta) &\leq [1 - \varphi(\sqrt{\beta})]^{-1} B_\Lambda \left(\int (2 - y \|x\|_R^{-1} \langle \theta, x \rangle)_+ d\bar{\mathbb{P}}(x, y) + \varphi(\sqrt{\beta}), \right. \\ &\quad \left. \frac{\beta \|\theta\|^2 + 2 \log(|\Lambda|/\epsilon)}{2n} \right) \\ &= [1 - \varphi(\sqrt{\beta})]^{-1} \inf_{\lambda \in \Lambda} \Phi_\lambda^{-1} \left[C_3(\lambda, \theta) + \varphi(\sqrt{\beta}) + \frac{\log(|\Lambda|/\epsilon)}{n\lambda} \right], \end{aligned}$$

where

$$C_3(\lambda, \theta) = \int (2 - y \|x\|_R^{-1} \langle \theta, x \rangle)_+ d\bar{\mathbb{P}}(x, y) + \frac{\beta \|\theta\|^2}{2n\lambda}.$$

Let us assume now that the patterns x are in a ball, so that $\|x\| \leq R$ almost surely. In this case $\|x\|_R = R$ almost surely. Let us remark that $L(\mathbb{P}, \theta) = L(\mathbb{P}, 2R\theta)$, and let us make the previous result uniform in $\beta \in \Xi$. This leads to

PROPOSITION 4.4 *Let us assume that $\|x\| \leq R$ almost surely. With probability at least $1 - \epsilon$, for all $\theta \in \mathbb{R}^d$,*

$$\begin{aligned} L(\mathbb{P}, \theta) &\leq \inf_{\beta \in \Xi} [1 - \varphi(\sqrt{\beta})]^{-1} \inf_{\lambda \in \Lambda} \Phi_\lambda^{-1} \left[2C_4(\beta, \lambda, \theta) \right. \\ &\quad \left. + \varphi(\sqrt{\beta}) + \frac{\log(|\Xi| |\Lambda|/\epsilon)}{n\lambda} \right], \end{aligned}$$

where

$$C_4(\beta, \lambda, \theta) = \frac{1}{2} C_3(\lambda, 2R\theta) = \int (1 - y \langle \theta, x \rangle)_+ d\bar{\mathbb{P}}(x, y) + \frac{\beta R^2 \|\theta\|^2}{n\lambda},$$

and

$$\Phi_\lambda^{-1}(q) = \frac{1 - \exp(-\lambda q)}{1 - \exp(-\lambda)} \leq \frac{q}{1 - \frac{\lambda}{2}}.$$

The loss function $C_4(\lambda, \theta)$ is the most employed learning criterion for support vector machines, and is called the box constraint. It is convex in θ . There are fast algorithms to compute $\inf_\theta C_4(\lambda, \theta)$ for any fixed values of λ and β . Here we get an empirical criterion which could be used to optimize also the values of λ and β , that is to optimize the strength of the regularizing factor $\frac{\beta R^2 \|\theta\|^2}{n\lambda}$.

In this criterion, $\|\theta\|^{-1}$ can be interpreted as the margin width, that is the minimal distance of x from the separating hyperplane $\{x' : \langle \theta, x' \rangle = 0\}$ beyond which the error term $(1 - y\langle \theta, x \rangle)_+$ vanishes (for data x that are on the right side of the separating hyperplane). The speed of convergence depends on $R^2\|\theta\|^2/n$. For this reason, $R^2\|\theta\|^2$, the square of the ratio between the radius of the ball containing the data and the margin, plays the role of the dimension. The bound does not depend on d , showing that with separating hyperplanes and more generally Support Vector Machines, we can get low error rates while choosing to represent the data in a Reproducing Kernel Hilbert Space with a large, or even infinite, dimension.

We considered so far only linear hyperplanes and data centered around 0. Anyhow, this also covers affine hyperplanes and data contained in a non necessarily centered ball, through a change of coordinates. More precisely, the previous proposition has the following corollary:

COROLLARY 4.5 *Assume that almost surely $\|x - c\| \leq R$, for some $c \in \mathbb{R}^d$ and $R \in \mathbb{R}_+$. With probability at least $1 - \epsilon$, for any $\theta \in \mathbb{R}^d$, any $\gamma \in \mathbb{R}$ such that $\min_{i=1,\dots,n} \langle \theta, x_i \rangle \leq \gamma \leq \max_{i=1,\dots,n} \langle \theta, x_i \rangle$,*

$$\int \mathbb{1}[y(\langle \theta, x \rangle - \gamma) \leq 0] d\mathbb{P}(x, y) \leq \inf_{\beta \in \Xi} [1 - \varphi(\sqrt{\beta})]^{-1} \inf_{\lambda \in \Lambda} \Phi_\lambda^{-1} \left[2C_5(\beta, \lambda, \theta, \gamma) + \varphi(\sqrt{\beta}) + \frac{\log(|\Xi| |\Lambda| / \epsilon)}{n\lambda} \right],$$

where

$$C_5(\beta, \lambda, \theta, \gamma) = \int [1 - y(\langle \theta, x \rangle - \gamma)]_+ d\bar{\mathbb{P}}(x, y) + \frac{4\beta R^2 \|\theta\|^2}{n\lambda}.$$

PROOF. Let us apply the previous result to $x' = (x - c, R)$, and $\theta' = [\theta, R^{-1}(\langle \theta, c \rangle - \gamma)]$. We get that $\|x'\|^2 \leq 2R^2$ and $\|\theta'\|^2 \leq 2\|\theta\|^2$, because almost surely $-\|\theta\|R \leq \text{ess inf} \langle \theta, x - c \rangle \leq \gamma - \langle \theta, c \rangle \leq \text{ess sup} \langle \theta, x - c \rangle \leq \|\theta\|R$, so that almost surely, for the allowed values of γ , $(\langle \theta, c \rangle - \gamma)^2 \leq R^2\|\theta\|^2$. This proves that $C_4(\beta, \lambda, \theta') \leq C_5(\beta, \lambda, \theta, \gamma)$, as required to deduce the corollary from the previous proposition. \square

REFERENCES

- [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.
- [10] V. Vapnik. *Estimation of Dependences Based on Empirical Data*. Springer-Verlag, Berlin, 1982.
- [11] V. Vapnik. *Statistical learning theory*. Wiley, 1998.
- [12] V. Vapnik and A. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory Probab. Appl.*, 16:264–280, 1971.