

---

### 3. PAC-BAYES BOUNDS FOR SUPERVISED CLASSIFICATION

3.1. INTRODUCTION. PAC-Bayes theory was first developed in the framework of supervised classification (see [6, 7, 8, 9, 5]) and subsequently extended to other settings. We will not deal with these extensions in the present notes, but instead focus on supervised classification, a setting that plays a central role in statistical learning theory and requires specific techniques of proofs.

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  $\Theta$  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  $\hat{\theta}(W_1^n)$  depending on the observed sample  $W_1^n$  such that the excess risk

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

is small. The previous quantity is random, since  $\hat{\theta}$  depends on the random sample  $W_1^n$ . Therefore, how small it is 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

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

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.

Let us remark that the point of view adopted in this section is different from the point of view of the previous section in some important ways. With the tools of the last section, we could have considered the loss function

$$L_c[(x, y), q_\theta] = -\log_2[q_\theta(y | x)],$$

where  $q_\theta(y|x)$  is some family of conditional distributions indexed by  $\theta$ .

We could then have obtained results concerned with

$$\sum_{i=1}^n L_c \left[ (X_i, Y_i), \int q_\theta d\rho_i(\theta) \right],$$

where  $\rho_i \in \mathcal{M}_+^1(\Theta)$  is a random probability measure defined with the help of some prior probability measure  $\pi \in \mathcal{M}_+^1(\Theta)$  by its density

$$\frac{d\rho_i}{d\pi}(\theta) = \frac{\prod_{i=1}^n q_\theta(Y_i|X_i)}{\int \prod_{i=1}^n q_{\theta'}(Y_i|X_i) d\pi(\theta')}.$$

Namely, the previous section provides the following almost sure upper bound:

$$\sum_{i=1}^n L_c \left[ (X_i, Y_i), \int q_\theta d\rho_i(\theta) \right] \leq \inf_{\theta \in \Theta} \sum_{i=1}^n L_c [(X_i, Y_i), q_\theta] - \log_2 [\pi(\{\theta\})]. \quad (7)$$

Importantly, the previous section is concerned with a *cumulated* notion of risk, whereas this one will deal with an *instantaneous* notion of risk. For the cumulated risk, we could obtain almost sure results, this will not be possible any more for the instantaneous risk we are considering here. We will have instead to make an hypothesis about the probabilistic nature of the observations, assuming that they are independent, and obtain results holding with a probability close to one but not equal to one.

**Exercise 7** *To make the link with the previous section more specific, we can focus on*

$$q_\theta(y|x) = \frac{\exp\{-\beta \mathbf{1}[f_\theta(x) \neq y]\}}{[1 + (|Y| - 1) \exp(-\beta)]}.$$

*Show that*

$$\begin{aligned} [1 - \exp(-\beta)] \sum_{i=1}^n \int \mathbf{1}[f_\theta(X_i) \neq Y_i] d\rho_i(\theta) \\ \leq \inf_{\theta \in \Theta} \beta \sum_{i=1}^n \mathbf{1}[f_\theta(X_i) \neq Y_i] - \log_2 [\pi(\{\theta\})]. \end{aligned}$$

*(Hint : consider a lower bound of the left-hand side of equation (7, page 25) and an upper bound of its right-hand side, using the fact that  $\log(1+z) \leq z$ .) Conclude that, almost surely*

$$\begin{aligned} \sum_{i=1}^n \int L[(X_i, Y_i), \theta] d\rho_i(\theta) \\ \leq \inf_{\theta \in \Theta} \frac{1}{1 - \beta/2} \sum_{i=1}^n L[(X_i, Y_i), \theta] - \frac{\log_2[\pi(\{\theta\})]}{\beta(1 - \beta/2)}. \end{aligned}$$

(Hint: use the inequality  $\exp(-\beta) \leq 1 - \beta + \beta^2/2$ ,  $\beta \geq 0$ .) How would you choose the value of the parameter  $\beta$  ?

3.2. 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  $\bar{\mathbb{P}}$  will be the empirical measure

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

of an i.i.d. sample  $(W_i)_{i=1}^n$  drawn from  $\mathbb{P}^{\otimes n} \in \mathcal{M}_+^1(\mathcal{W}^n)$ . 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

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

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

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

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

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

Let us mention here some important identity.

**PROPOSITION 3.2** *For any probability measures  $\pi$  and  $\rho$  on some 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 define

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

The expectations with respect to  $\rho$  and  $\pi$  of  $h$  and the log-Laplace transform of  $h$  are linked by the identities

$$\int h d\rho - \mathcal{K}(\rho, \pi) + \mathcal{K}(\rho, \pi_{\exp(h)}) = \log \left[ \int \exp(h) d\pi \right] \quad (8)$$

$$= \int h d\pi + \int_0^1 (1 - \alpha) \mathbf{Var}[h d\pi_{\exp(\alpha h)}] d\alpha. \quad (9)$$

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

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

and 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 26). 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  $\lambda$ , giving

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

**Exercise 9** 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  $\lambda$ . 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 (8, page 27) 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

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

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 10** *The second part of Proposition 3.1 (page 26) is proved in the same way and left as an exercise.*

**3.3. PAC-BAYES BOUNDS.** We are now going to make Proposition 3.1 uniform with respect to  $\theta$ . The PAC-Bayes approach to this is to randomize  $\theta$ , so we will consider now joint distributions on  $(W_1, \dots, W_n, \theta)$ , where the distribution of  $(W_1, \dots, W_n)$  is still  $\mathbb{P}^{\otimes n}$  and the conditional distribution of  $\theta$  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  $\rho$  will be compared with a prior (meaning non-random) probability measure  $\pi \in \mathcal{M}_+^1(\Theta)$ .

**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  $\lambda \in \mathbb{R}_+$ ,*

$$\int \exp \left[ \sup_{\rho \in \mathcal{M}_+^1(\Theta)} n\lambda \left\{ \Phi_\lambda [L(\mathbb{P}, \rho)] - L(\bar{\mathbb{P}}, \rho) \right\} - \mathcal{K}(\rho, \pi) \right] d\mathbb{P}^{\otimes n} \leq 1, \quad (10)$$

*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 (10) is a consequence of equation (8, page 27), showing that

---

\*We will assume that  $\rho$  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  $\sigma$ -algebra we consider on  $\Theta$  is generated by a countable family of subsets. See [1][page 50] for more details

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

and of the fact that  $\Phi_\lambda$  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 26), because we can no more cast our deviation inequality in such a way that  $\lambda$  appears on some non-random side of the inequality. Nevertheless, we can get a more explicit bound from some specific choice of the set  $\Lambda$ .

**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  $n \geq 5$ ,

$$B_{\lfloor \log(n)^2 \rfloor - 1}(q, e, \epsilon) \leq B(q, e, \epsilon) \stackrel{\text{def}}{=} \sqrt{\frac{2\bar{v}(q)\{e + \log[\log(n)^2/\epsilon]\}}{n}} \cosh[\log(n)^{-1}] + \frac{2\{e + \log[\log(n)^2/\epsilon]\}}{n} \cosh[\log(n)^{-1}]^2, \quad (11)$$

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

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

PROOF. Let us put

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

According to equation (9, page 27) 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

$$p - q \leq \sqrt{2\delta\bar{v}(q)} \cosh(t/m) + 2\delta(1-q) \cosh(t/m)^2. \quad (12)$$

Let us consider now the case when  $\hat{\lambda} > 1$ . In this case

$$p - q \leq \sqrt{2\delta\bar{v}(p)} \hat{\lambda} = 2\delta.$$

In conclusion, applying Proposition 3.3 (page 30) we see that with probability at least  $1 - \epsilon$ , for any posterior distribution  $\rho$ ,

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

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.  $\square$

## 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] F.M.J. Willems, Y.M. Shtarkov, and T.J. Tjalkens. The context-tree weighting method: basic properties. *IEEE Trans. Inform. Theory*, 41(3):653–664, 1995.