

# A PAC-Bayes Sample-compression Approach to Kernel Methods : Supplementary material

Pascal Germain, Alexandre Lacoste, François Laviolette, Mario Marchand and Sara Shanian.

Département d'informatique et de génie logiciel, Université Laval  
Québec, Canada, G1V 0A6

Last Revision: July 18, 2011

This supplementary refers to the ICML 2011 paper *A PAC-Bayes Sample-compression Approach to Kernel Methods* of Germain et al. (2011).

The latest revision of this document and related source code are available at:  
<http://graal.ift.ulaval.ca/publications.php>

## Contents

<b>1 Proof of Claim 1 in the proof of Theorem 5</b>	<b>2</b>
<b>2 Proof of Theorem 1: The PAC-Bayes bound with <math>KL(Q  P)</math></b>	<b>4</b>
<b>3 Details related to the PBSC algorithms</b>	<b>8</b>
3.1 PBSC-A: The aligned case . . . . .	8
3.2 PBSC-N: The non-aligned case . . . . .	9
<b>4 Empirical Results</b>	<b>12</b>
<b>5 Another PAC-Bayes bound without <math>KL(Q  P)</math> (not stated in main the paper)</b>	<b>13</b>
<b>References</b>	<b>15</b>

# 1 Proof of Claim 1 in the proof of Theorem 5

**Claim 1 :**  $\mathbf{E}_{S \sim D^m} X_{\bar{P}} \leq e^{Ald} \cdot 2\sqrt{m}$ .

The claim uses the following lemma.

**Lemma 0.** (*Maurer (2004)*) Let  $n \geq 8$ , and suppose that  $X = (X_1, \dots, X_n)$  is a vector of iid random variables,  $0 \leq X_i \leq 1$ ,  $\mathbf{E}[X_i] = \nu$  and let  $M(X) = \frac{1}{n} \sum_{j=1}^n X_j$  be the arithmetic mean of the random variables. Then

$$\sqrt{n} \leq \mathbf{E} e^{n \text{kl}(M(X) \parallel \nu)} \leq 2\sqrt{n}.$$

*Proof of Claim 1.* Let us first recall that

$$X_{\bar{P}} \stackrel{\text{def}}{=} \mathbf{E}_{\bar{h} \sim \bar{P}} e^{(m - |\mathbf{i}_{\bar{h}}|) \cdot 2(\bar{R}_S(\bar{h}) - \bar{R}_D(\bar{h}))^2}.$$

Let us also introduce a new notation. Let  $\tilde{R}_S(\bar{h})$  be the *abstract empirical risk* computed on the examples of  $S$  that are not in the compression sequence of  $\bar{h}$ . More formally,

$$\tilde{R}_S(\overline{h_1 \dots h_k}) \stackrel{\text{def}}{=} \frac{1}{m - |\mathbf{i}_{\overline{h_1 \dots h_k}}|} \sum_{j=1}^m I\left(\neg \bigvee_{i=1}^k (h_i(x_j) \neq y_j)\right) I\left((x_j, y_j) \notin \mathbf{i}_{\overline{h_1 \dots h_k}}\right),$$

where  $\bigvee$  denotes the exclusive or. Based on the definition of Equation (9) in the main paper, one can easily show that  $\tilde{R}_S(\bar{h}) = \bar{R}_U(\bar{h})$  with  $U = S \setminus S_{\mathbf{i}_{\bar{h}}}$ , where  $\mathbf{i}_{\bar{h}}$  points to the examples belonging to the union of all compression sequences of the sc-classifiers in  $\bar{h}$ . Moreover, note that, contrarily to  $\bar{R}_S(\bar{h})$  which may contain some bias,  $\tilde{R}_S(\bar{h})$  is an arithmetic mean of truly iid  $(m - |\mathbf{i}_{\bar{h}}|)$  random variables. On another hand, these two random variables have values that are very close to each other. Indeed, since

$$0 \leq m \cdot \bar{R}_S(\bar{h}) - (m - |\mathbf{i}_{\bar{h}}|) \cdot \tilde{R}_S(\bar{h}) \leq |\mathbf{i}_{\bar{h}}|,$$

we have

$$-|\mathbf{i}_{\bar{h}}| \leq -|\mathbf{i}_{\bar{h}}| \cdot \tilde{R}_S(\bar{h}) \leq m \cdot \bar{R}_S(\bar{h}) - m \cdot \tilde{R}_S(\bar{h}) \leq |\mathbf{i}_{\bar{h}}| - |\mathbf{i}_{\bar{h}}| \cdot \tilde{R}_S(\bar{h}) \leq |\mathbf{i}_{\bar{h}}|.$$

Therefore,

$$(20) \quad \left| \bar{R}_S(\bar{h}) - \tilde{R}_S(\bar{h}) \right| \leq \frac{|\mathbf{i}_{\bar{h}}|}{m}.$$

Given a compression sequence  $S_{\mathbf{i}}$ , let  $\mathbf{i}^c$  be the vector of indices of  $\mathcal{I}$  that are not in the vector  $\mathbf{i}$ . We now have

$$\begin{aligned} \mathbf{E}_{S \sim D^m} X_{\bar{P}} &\stackrel{\text{def}}{=} \mathbf{E}_{S \sim D^m} \mathbf{E}_{\bar{h} \sim \bar{P}} e^{(m - |\mathbf{i}_{\bar{h}}|) \cdot 2(\bar{R}_S(\bar{h}) - \bar{R}_D(\bar{h}))^2} \\ &= \mathbf{E}_{\mathbf{i} \sim \bar{P}_{\mathcal{I}}} \mathbf{E}_{S_{\mathbf{i}} \sim D^{|\mathbf{i}|}} \mathbf{E}_{\mu \sim \bar{P}_{S_{\mathbf{i}}}} \mathbf{E}_{S_{\mathbf{i}^c} \sim D^{m - |\mathbf{i}|}} e^{(m - |\mathbf{i}|) \cdot 2(\bar{R}_S(\bar{h}_{\mathbf{i}}^{\mu}) - \bar{R}_D(\bar{h}_{\mathbf{i}}^{\mu}))^2}, \end{aligned}$$

Hence, to prove Claim 1, it suffices to show that for all  $\mathbf{i} \in \mathcal{I}$ ,  $S_{\mathbf{i}} \in (\mathcal{X} \times \mathcal{Y})^{|\mathbf{i}|}$ , and  $\mu \in \mathcal{M}_{S_{\mathbf{i}}}$ , we have

$$(21) \quad \mathbf{E}_{S_{\mathbf{i}^c} \sim D^{m - |\mathbf{i}|}} e^{(m - |\mathbf{i}|) \cdot 2(\bar{R}_S(\bar{h}_{\mathbf{i}}^{\mu}) - \bar{R}_D(\bar{h}_{\mathbf{i}}^{\mu}))^2} \leq e^{Ald} \cdot 2\sqrt{m}.$$

Here is the sketch of the proof of Equation (21). Justification for Line (22) to (25) follows below.

$$\begin{aligned}
& \mathbf{E}_{S_{\mathbf{i}c} \sim D^{m-|\mathbf{i}|}} e^{(m-|\mathbf{i}|) \cdot 2 (\overline{R}_S(\overline{h}_i^\mu) - \overline{R}_D(\overline{h}_i^\mu))^2} \\
&= \mathbf{E}_{S_{\mathbf{i}c} \sim D^{(m-|\mathbf{i}|)}} e^{(m-|\mathbf{i}|) \cdot 2 (\overline{R}_S(\overline{h}_i^\mu) - \widetilde{R}_S(\overline{h}_i^\mu) + \widetilde{R}_S(\overline{h}_i^\mu) - \overline{R}_D(\overline{h}_i^\mu))^2} \\
&\leq \mathbf{E}_{S_{\mathbf{i}c} \sim D^{(m-|\mathbf{i}|)}} e^{(m-|\mathbf{i}|) \cdot 2 \left( [\overline{R}_S(\overline{h}_i^\mu) - \widetilde{R}_S(\overline{h}_i^\mu)]^2 + 2 |\overline{R}_S(\overline{h}_i^\mu) - \widetilde{R}_S(\overline{h}_i^\mu)| \cdot |\widetilde{R}_S(\overline{h}_i^\mu) - \overline{R}_D(\overline{h}_i^\mu)| + [\widetilde{R}_S(\overline{h}_i^\mu) - \overline{R}_D(\overline{h}_i^\mu)]^2 \right)} \\
(22) \quad &\leq \mathbf{E}_{S_{\mathbf{i}c} \sim D^{(m-|\mathbf{i}|)}} e^{(m-|\mathbf{i}|) \cdot 2 \left( \left[ \frac{|\mathbf{i}|}{m} \right]^2 + 2 \frac{|\mathbf{i}|}{m} \cdot 1 + [\widetilde{R}_S(\overline{h}_i^\mu) - \overline{R}_D(\overline{h}_i^\mu)]^2 \right)} \\
(23) \quad &\leq \mathbf{E}_{S_{\mathbf{i}c} \sim D^{(m-|\mathbf{i}|)}} e^{4ld + (m-|\mathbf{i}|) \cdot 2 [\widetilde{R}_S(\overline{h}_i^\mu) - \overline{R}_D(\overline{h}_i^\mu)]^2} \\
(24) \quad &\leq \mathbf{E}_{S_{\mathbf{i}c} \sim D^{(m-|\mathbf{i}|)}} e^{4ld + (m-|\mathbf{i}|) \cdot \text{kl}(\widetilde{R}_S(\overline{h}_i^\mu) \| \overline{R}_D(\overline{h}_i^\mu))} \\
&= e^{4ld} \cdot \mathbf{E}_{S_{\mathbf{i}c} \sim D^{(m-|\mathbf{i}|)}} e^{(m-|\mathbf{i}|) \cdot \text{kl}(\widetilde{R}_S(\overline{h}_i^\mu) \| \overline{R}_D(\overline{h}_i^\mu))} \\
(25) \quad &\leq e^{4ld} \cdot 2\sqrt{m-|\mathbf{i}|} \\
&\leq e^{4ld} \cdot 2\sqrt{m}
\end{aligned}$$

Line (22) follows from Equation (20) and the fact that the exponential function is increasing. For Line (23), since  $|\mathbf{i}| \leq ld$ , we simply have to show that  $(m-|\mathbf{i}|) \cdot 2 \left( \frac{|\mathbf{i}|}{m^2} + \frac{2}{m} \right) \leq 4$ , which follows from direct calculations. Line (24) follows directly from the property :  $2(q-p)^2 \leq \text{kl}(q \| p)$ , which can also be proved via straightforward calculations. Finally, for Line (25), first observe that  $\widetilde{R}_S(\overline{h}_i^\mu)$  is an arithmetic mean of  $(m-|\mathbf{i}|)$  iid random variables. Thus Line (25) is simply an application of Lemma 0 with  $M(X)$  replaced by  $\widetilde{R}_S(\overline{h}_i^\mu)$ ,  $n$  replaced by  $m-|\mathbf{i}|$ , and  $\nu$  replaced by  $\overline{R}_D(\overline{h}_i^\mu)$ . Thus, Claim 1 is proved.  $\square$

## 2 Proof of Theorem 1: The PAC-Bayes bound with $KL(Q\|P)$

**Theorem 1.** For any  $D$ , any family  $(\mathcal{H}^S)_{S \in D^m}$  of sets of sc-classifiers of size at most  $l$ , any prior  $\mathcal{P}$ , any  $\delta \in (0, 1]$ , any positive real number  $C_1$ , and any margin loss function  $\zeta$  such that  $l \cdot \deg(\zeta) < m$ , we have

$$\Pr_{S \sim D^m} \left( \forall Q \text{ on } \mathcal{H}^S: \zeta_D^Q \leq \zeta(1)[C' - 1] + C' \cdot \left( \zeta_S^Q + \frac{2}{m \cdot C_1} [\zeta'(1) \cdot KL(Q\|P) + \zeta(1) \cdot \ln \frac{1}{\delta}] \right) \right) \geq 1 - \delta$$

where  $KL(\cdot\|\cdot)$  is the Kullback-Leibler divergence, and where  $C' = \frac{C_1 \cdot \frac{m}{m-l \cdot \deg \zeta}}{1 - e^{-C_1 \cdot \frac{m-l \cdot \deg \zeta}{m}}}$ .

*Proof.* Let  $S$  be any training sequence,  $d \stackrel{\text{def}}{=} \deg \zeta$ . Let  $\mathcal{F}$  be a convex function to be defined later, and  $\mathcal{D}(q, p) \stackrel{\text{def}}{=} \mathcal{F}(p) - C_1 \cdot q$ .

For each  $k \in \{0, \dots, d\}$  and any  $k$ -tuple  $(h_1, \dots, h_k)$ , let us define  $\bar{h} = \overline{h_1 \dots h_k}$  as an ‘‘abstract’’ sc-classifier whose ‘‘abstract’’ true and empirical risks are respectively defined as

$$\begin{aligned} \bar{R}_D(\overline{h_1 \dots h_k}) &\stackrel{\text{def}}{=} \mathbf{E}_{(\mathbf{x}, y) \sim D} I(\neg \bigvee_{i=1}^k (h_i(x) \neq y)), \\ \bar{R}_S(\overline{h_1 \dots h_k}) &\stackrel{\text{def}}{=} \mathbf{E}_{(x, y) \sim S} I(\neg \bigvee_{i=1}^k (h_i(x) \neq y)) = \frac{1}{m} \sum_{j=1}^m I(\neg \bigvee_{i=1}^k (h_i(x_j) \neq y_j)), \end{aligned}$$

where  $\bigvee$  denotes the exclusive or.

Considering the above definitions, we also have

$$\bar{R}_D(\overline{h_1 \dots h_k}) = \mathbf{E}_{(\mathbf{x}, y) \sim D} I(\neg \bigvee_{i=1}^k (h_i(x) \neq y)) = \mathbf{E}_{(\mathbf{x}, y) \sim D} \frac{1}{2} \left[ 1 + \prod_{i=1}^k -y h_i(x) \right].$$

Since the compression sequence size of each  $h_i$ 's is at most  $l$ , we have  $|\mathbf{i}_{\bar{h}}| \leq l \cdot k$  for any  $\bar{h} = \overline{h_1 \dots h_k}$ . Moreover, we have  $\bar{R}_D(\overline{h_1 \dots h_0}) = 1$  when  $k = 0$  because the exclusive or over an empty sequence outputs *false*. For each  $S$ , let  $\overline{\mathcal{H}}^S$  be the set of all such sc-classifiers, and for each distribution  $Q$  on  $\mathcal{H}^S$ , denote by  $\bar{Q}$  the following distribution on  $\overline{\mathcal{H}}^S$ :

$$\bar{Q}(\overline{h_1 \dots h_k}) \stackrel{\text{def}}{=} \frac{a_k}{\zeta(1)} Q(h_1) \cdot \dots \cdot Q(h_k) = \frac{a_k}{\zeta(1)} \prod_{i=1}^k Q(h_i).$$

Since  $\zeta(1) = \sum_{k=0}^d a_k$ ,  $\bar{Q}$  is a probability distribution. We also denote by  $\bar{P}$  the following distribution on  $\overline{\mathcal{H}}^S$ :

$$\bar{P}(\overline{h_1 \dots h_k}) \stackrel{\text{def}}{=} \frac{a_k}{\zeta(1)} \cdot P(h_1) \cdot \dots \cdot P(h_k) = \frac{a_k}{\zeta(1)} \prod_{i=1}^k P(h_i).$$

Moreover, for  $U = D$  and  $U = S$ , we have

$$\begin{aligned}
R_U(G_{\bar{Q}}) &= \mathbf{E}_{\bar{h} \sim \bar{Q}} \bar{R}_U(\bar{h}) \\
&= \sum_{k=0}^{\deg(\zeta)} \frac{a_k}{\zeta(1)} \mathbf{E}_{h_1 \sim Q} \cdots \mathbf{E}_{h_k \sim Q} \mathbf{E}_{(\mathbf{x}, y) \sim D} \frac{1}{2} \left[ 1 + \prod_{i=1}^k -yh_i(x) \right] \\
&= \sum_{k=0}^{\deg(\zeta)} \frac{a_k}{\zeta(1)} \mathbf{E}_{(\mathbf{x}, y) \sim D} \frac{1}{2} \left[ 1 + \prod_{i=1}^k \mathbf{E}_{h_i \sim Q} -yh_i(x) \right] \\
&= \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \mathbf{E}_{(\mathbf{x}, y) \sim D} \sum_{k=0}^{\deg(\zeta)} a_k (\mathbf{E}_{h \sim Q} -yh(x))^k \right] \\
&= \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \mathbf{E}_{(\mathbf{x}, y) \sim D} \sum_{k=0}^{\deg(\zeta)} a_k (-M_Q(x, y))^k \right] \\
(26) \quad &= \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \zeta_U^Q \right].
\end{aligned}$$

Now, let us consider the following Laplace transform

$$X_{\bar{P}} \stackrel{\text{def}}{=} \mathbf{E}_{\bar{h} \sim \bar{P}} e^{(m - |\mathbf{i}_{\bar{h}}|) \mathcal{D}(\bar{R}_S(\bar{h}), \bar{R}_D(\bar{h}))}.$$

It then follows from Markov's inequality that

$$\Pr_{S \sim D^m} \left( X_{\bar{P}} \leq \frac{1}{\delta} \mathbf{E}_{S \sim D^m} X_{\bar{P}} \right) \geq 1 - \delta.$$

By taking the logarithm on each side of the innermost inequality and by transforming the expectation over  $\bar{P}$  into an expectation over  $\bar{Q}$ , we obtain

$$(27) \quad \Pr_{S \sim D^m} \left( \forall Q: \ln \left[ \mathbf{E}_{\bar{h} \sim \bar{Q}} \frac{\bar{P}(\bar{h})}{Q(\bar{h})} e^{(m - |\mathbf{i}_{\bar{h}}|) \mathcal{D}(\bar{R}_S(\bar{h}), \bar{R}_D(\bar{h}))} \right] \leq \ln \left[ \frac{1}{\delta} \mathbf{E}_{S \sim D^m} \mathbf{E}_{\bar{h} \sim \bar{P}} e^{(m - |\mathbf{i}_{\bar{h}}|) \mathcal{D}(\bar{R}_S(\bar{h}), \bar{R}_D(\bar{h}))} \right] \right) \geq 1 - \delta.$$

Since  $\zeta'(1) = \sum_{k=1}^d k \cdot a_k$ , by straightforward calculation, one can show that

$$\mathbf{E}_{\bar{h} \sim \bar{Q}} \ln \frac{\bar{P}(\bar{h})}{Q(\bar{h})} = -\frac{\zeta'(1)}{\zeta(1)} \cdot KL(Q \| P).$$

This, together with Jensen's inequality applied to the concave  $\ln(x)$ , gives

$$(28) \quad \ln \left[ \mathbf{E}_{\bar{h} \sim \bar{Q}} \frac{\bar{P}(\bar{h})}{Q(\bar{h})} e^{(m - |\mathbf{i}_{\bar{h}}|) \mathcal{D}(\bar{R}_S(\bar{h}), \bar{R}_D(\bar{h}))} \right] \geq -\frac{\zeta'(1)}{\zeta(1)} \cdot KL(Q \| P) + \mathbf{E}_{\bar{h} \sim \bar{Q}} (m - |\mathbf{i}_{\bar{h}}|) \mathcal{D}(\bar{R}_S(\bar{h}), \bar{R}_D(\bar{h})).$$

Again from the Jensen's inequality applied to the convex function  $\mathcal{F}$ , together with Equation (26), and the fact that  $m - l \cdot d \leq (m - |\mathbf{i}_{\bar{h}}|) \leq m$ , we obtain

$$\begin{aligned}
(29) \quad \mathbf{E}_{\bar{h} \sim \bar{Q}} (m - |\mathbf{i}_{\bar{h}}|) \mathcal{D}(\bar{R}_S(\bar{h}), \bar{R}_D(\bar{h})) &\geq (m - ld) \mathbf{E}_{\bar{h} \sim \bar{Q}} \mathcal{F}(\bar{R}_D(\bar{h})) - m \mathbf{E}_{\bar{h} \sim \bar{Q}} C_1 \cdot \bar{R}_S(\bar{h}) \\
&\geq (m - ld) \mathcal{F} \left( \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \zeta_D^Q \right] \right) - m C_1 \cdot \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \zeta_S^Q \right].
\end{aligned}$$

Let us analyze the value of  $\mathbf{E}_{S \sim D^m} X_{\bar{P}}$ , appearing in the right-hand side of the innermost inequality of Equation (27). First, let us define  $\mathbf{i}^c$  as the vector of indices of  $\mathcal{I}$  that are not in the vector  $\mathbf{i}$ . Thus,  $|\mathbf{i}^c| = m - |\mathbf{i}|$ . Now, note that

$$\mathbf{E}_{S \sim D^m} \mathbf{E}_{\bar{h} \sim \bar{P}} e^{(m-|\mathbf{i}^c|) \mathcal{D}(\bar{R}_S(\bar{h}), \bar{R}_D(\bar{h}))} = \mathbf{E}_{\mathbf{i} \sim \bar{P}_{\mathcal{I}}} \mathbf{E}_{S_{\mathbf{i}} \sim D^{|\mathbf{i}|}} \mathbf{E}_{\mu \sim \bar{P}_{S_{\mathbf{i}}}} \mathbf{E}_{S_{\mathbf{i}^c} \sim D^{m-|\mathbf{i}|}} e^{|\mathbf{i}^c| \mathcal{D}(\bar{R}_S(\bar{h}_{\mathbf{i}}^\mu), \bar{R}_D(\bar{h}_{\mathbf{i}}^\mu))}.$$

Now, for each  $\bar{h}_{\mathbf{i}}^\mu \in \bar{\mathcal{H}}^S$ , define  $a_{S_{\mathbf{i}}}^\mu \stackrel{\text{def}}{=} \sum_{(x,y) \in S_{\mathbf{i}}} I(\bar{h}_{\mathbf{i}}^\mu(x) \neq y)$ . Observe that  $m \cdot \bar{R}_S(\bar{h}_{\mathbf{i}}^\mu) - a_{S_{\mathbf{i}}}^\mu$  is the number of errors made by  $\bar{h}_{\mathbf{i}}^\mu$  on  $S_{\mathbf{i}^c}$ . Since the later is iid and disjoint from the compression sequence of  $\bar{h}_{\mathbf{i}}^\mu$ , we have that  $m \cdot \bar{R}_S(\bar{h}_{\mathbf{i}}^\mu) - a_{S_{\mathbf{i}}}^\mu$  is a random variable following a binomial law of parameters  $(|\mathbf{i}^c|, \bar{R}_D(\bar{h}_{\mathbf{i}}^\mu))$ . Since,  $(m - ld) \cdot \frac{k}{m} \leq |\mathbf{i}^c| \cdot \frac{k + a_{S_{\mathbf{i}}}^\mu}{m}$  for any  $k \in \{0, \dots, |\mathbf{i}^c|\}$ , we have

$$\begin{aligned} \mathbf{E}_{S_{\mathbf{i}^c} \sim D^{m-|\mathbf{i}|}} e^{|\mathbf{i}^c| \mathcal{D}(\bar{R}_S(\bar{h}_{\mathbf{i}}^\mu), \bar{R}_D(\bar{h}_{\mathbf{i}}^\mu))} &= \mathbf{E}_{S_{\mathbf{i}^c} \sim D^{m-|\mathbf{i}|}} e^{|\mathbf{i}^c| \mathcal{F}(\bar{R}_D(\bar{h}_{\mathbf{i}}^\mu)) - C_1 |\mathbf{i}^c| \bar{R}_S(\bar{h}_{\mathbf{i}}^\mu)} \\ &= e^{|\mathbf{i}^c| \mathcal{F}(\bar{R}_D(\bar{h}_{\mathbf{i}}^\mu))} \cdot \sum_{k=0}^{|\mathbf{i}^c|} \Pr_{S \sim D^m} (m \cdot \bar{R}_S(\bar{h}_{\mathbf{i}}^\mu) - a_{S_{\mathbf{i}}}^\mu = k) e^{-C_1 |\mathbf{i}^c| \cdot \frac{k + a_{S_{\mathbf{i}}}^\mu}{m}} \\ (30) \quad &\leq e^{|\mathbf{i}^c| \mathcal{F}(\bar{R}_D(\bar{h}_{\mathbf{i}}^\mu))} \cdot \sum_{k=0}^{|\mathbf{i}^c|} \Pr_{S \sim D^m} (m \cdot \bar{R}_S(\bar{h}_{\mathbf{i}}^\mu) - a_{S_{\mathbf{i}}}^\mu = k) e^{-C_1 \cdot (m-ld) \cdot \frac{k}{m}} \\ &= e^{|\mathbf{i}^c| \mathcal{F}(\bar{R}_D(\bar{h}_{\mathbf{i}}^\mu))} \cdot \sum_{k=0}^{|\mathbf{i}^c|} \binom{|\mathbf{i}^c|}{k} (\bar{R}_D(\bar{h}_{\mathbf{i}}^\mu))^k (1 - \bar{R}_D(\bar{h}_{\mathbf{i}}^\mu))^{|\mathbf{i}^c| - k} e^{-C_1 \cdot \frac{m-ld}{m} \cdot k} \\ &= e^{|\mathbf{i}^c| \mathcal{F}(\bar{R}_D(\bar{h}_{\mathbf{i}}^\mu))} \left( 1 - \bar{R}_D(\bar{h}_{\mathbf{i}}^\mu) [1 - e^{-C_1 \cdot \frac{m-ld}{m}}] \right)^{|\mathbf{i}^c|}. \end{aligned}$$

The last equation being obtained from the Newton binomial  $\sum_{i=0}^k \binom{m}{i} x^i y^{m-i} = (x+y)^m$ .

Let us now define  $\mathcal{F}$  such that  $1 = e^{|\mathbf{i}^c| \mathcal{F}(\bar{R}_D(\bar{h}_{\mathbf{i}}^\mu))} \left( 1 - \bar{R}_D(\bar{h}_{\mathbf{i}}^\mu) [1 - e^{-C_1 \cdot \frac{m-ld}{m}}] \right)^{|\mathbf{i}^c|}$ . Equivalently, let

$$(31) \quad \mathcal{F}(\bar{R}_D(\bar{h}_{\mathbf{i}}^\mu)) \stackrel{\text{def}}{=} -\ln \left( 1 - \bar{R}_D(\bar{h}_{\mathbf{i}}^\mu) \left[ 1 - e^{-C_1 \cdot \frac{m-ld}{m}} \right] \right).$$

With this choice, we have  $\mathbf{E}_{S \sim D^m} \mathbf{E}_{\bar{h} \sim \bar{P}} e^{|\mathbf{i}^c| \mathcal{D}(\bar{R}_S(\bar{h}_{\mathbf{i}}^\mu), \bar{R}_D(\bar{h}_{\mathbf{i}}^\mu))} = 1$ .

To finish the proof, let us combine Equations (28), (29), (30) and (31), in order to rewrite the innermost inequality of Equation (27) as follows

$$\begin{aligned} (m-ld) \cdot \mathcal{F} \left( \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \zeta_D^Q \right] \right) - mC_1 \cdot \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \zeta_S^Q \right] - \frac{\zeta'(1)}{\zeta(1)} \cdot \text{KL}(Q\|P) &\leq \ln \frac{1}{\delta} \\ (m-ld) \left\{ -\ln \left( 1 - \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \zeta_D^Q \right] \left[ 1 - e^{-C_1 \cdot \frac{m-ld}{m}} \right] \right) \right\} &\leq mC_1 \cdot \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \zeta_S^Q \right] + \frac{\zeta'(1)}{\zeta(1)} \cdot \text{KL}(Q\|P) + \ln \frac{1}{\delta} \\ \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \zeta_D^Q \right] \left[ 1 - e^{-C_1 \cdot \frac{m-ld}{m}} \right] &\leq 1 - \exp \left\{ - \left( \frac{1}{m-ld} \right) \left( mC_1 \cdot \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \zeta_S^Q \right] + \frac{\zeta'(1)}{\zeta(1)} \cdot \text{KL}(Q\|P) + \ln \frac{1}{\delta} \right) \right\} \\ \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \zeta_D^Q \right] \left[ 1 - e^{-C_1 \cdot \frac{m-ld}{m}} \right] &\leq \left( \frac{1}{m-ld} \right) \left( mC_1 \cdot \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \zeta_S^Q \right] + \frac{\zeta'(1)}{\zeta(1)} \cdot \text{KL}(Q\|P) + \ln \frac{1}{\delta} \right), \end{aligned}$$

where the last transformation is an application of the inequality  $1 - e^{-x} \leq x$ . We are now able to isolate  $\zeta_D^Q$  to obtain

$$\begin{aligned}
\zeta_D^Q &\leq \left( \frac{2 \cdot \zeta(1)}{1 - e^{-C_1 \frac{m-l d}{m}}} \right) \left( \frac{1}{m-l d} \right) \left( m C_1 \cdot \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \zeta_S^Q \right] + \frac{\zeta'(1)}{\zeta(1)} \cdot \text{KL}(Q \| P) + \ln \frac{1}{\delta} \right) - \zeta(1) \\
&= \left( \frac{C_1 \frac{m}{m-l d}}{1 - e^{-C_1 \frac{m-l d}{m}}} \right) \left( \zeta(1) + \zeta_S^Q + \frac{2}{m C_1} [\zeta'(1) \cdot \text{KL}(Q \| P) + \zeta(1) \cdot \ln \frac{1}{\delta}] \right) - \zeta(1) \\
&= \zeta(1) [C' - 1] + C' \cdot \left( \zeta_S^Q + \frac{2}{m C_1} [\zeta'(1) \cdot \text{KL}(Q \| P) + \zeta(1) \cdot \ln \frac{1}{\delta}] \right),
\end{aligned}$$

where

$$C' = \frac{C_1 \cdot \frac{m}{m-l \cdot \deg \zeta}}{1 - e^{-C_1 \cdot \frac{m-l \cdot \deg \zeta}{m}}}.$$

□

### 3 Details related to the PBSC algorithms

In this section, we present the theoretical development leading the two PBSC learning algorithms and the optimization procedures. Both algorithms build a majority vote of sc-classifier of compression sequence size of at most one as defined in Section 2.2 of the main paper. The two algorithms that we present minimize a bound on the quadratic loss. Given a fixed parameter  $q$ , the loss an example having margin  $(-\alpha)$  is given by:

$$\zeta(\alpha) = \left(1 + \frac{1}{q}\alpha\right)^2.$$

The first algorithm, named *PBSC-A*, works with a strongly aligned posterior  $Q$ . The second algorithm, named *PBSC-N*, works with a non-aligned posterior  $Q$ .

#### 3.1 PBSC-A: The aligned case

As seen in Section 2.1 of the paper, the strongly aligned posterior  $Q$  is totally defined by a vector  $\mathbf{w} \stackrel{\text{def}}{=} (w_0, w_1, \dots, w_m)$ . For  $Q$  to remain a valid distribution, each component of the vector  $\mathbf{w}$  must remain in the interval  $\left[-\frac{1}{m+1}, +\frac{1}{m+1}\right]$ .

Theorem 5 suggests to minimize the bound on  $\zeta_D^Q$  given by the following expression:

$$\zeta_D^Q \leq \zeta_S^Q + \frac{\zeta(1)}{\sqrt{\frac{1}{2}(m - l \deg \zeta)}} \sqrt{4l \deg \zeta + \ln \frac{2\sqrt{m}}{\delta}}.$$

To do so, we only need to minimize the empirical risk  $\zeta_S^Q$  because the last term of the right hand side is constant. The empirical risk,  $\zeta_S^Q$ , is given by

$$\zeta_S^Q = \frac{1}{mq^2} \sum_{j=1}^m \left( q - y_j \left[ w_0 + \sum_{i=1}^m w_i k(x_i, x_j) \right] \right)^2.$$

Lets define a matrix  $\mathbf{G}$  of size  $m + 1 \times m$  as

$$G_{i,j} = \begin{cases} 1 & \text{for } i = 0, \\ k(x_i, x_j) & \text{for } 1 \leq i, j \leq m. \end{cases}$$

With this notation, the optimization problem of *PBSC-A* can be written as

$$\begin{aligned} \text{Minimize: } f_A(\mathbf{w}) &= \sum_{j=1}^m \left( q - y_j \sum_{i=0}^m w_i G_{i,j} \right)^2 \\ \text{subject to: } & |w_i| \leq \frac{1}{m+1} \text{ for } i = 0, 1, \dots, m. \end{aligned}$$

We propose to solve this optimization problem by minimizing  $f_A$  coordinate-wise, similarly as it is done for AdaBoost (Schapire et al. (1998)), with the difference that we will have to ensure that  $Q$  remains an aligned distribution at each step of the algorithm. Starting from the uniform distribution  $P$  (i.e.,  $\mathbf{w} = \mathbf{0}$ ), the learning algorithm iteratively chooses (at random)  $k \in \{0, \dots, m\}$ , and updates  $w_k \leftarrow w_k + \delta$  (without updating the other weights) according to some optimally chosen value of  $\delta$ . Let  $\mathbf{w}_\delta$  be the new weight vector obtained with such an update. After an update, the objective function becomes:

$$f_A(\mathbf{w}_\delta) = \sum_{j=1}^m \left[ q - y_j \left( \sum_{i=0}^m w_i G_{i,j} + \delta G_{k,j} \right) \right]^2.$$



The optimal value for  $\delta$  is obtained when  $\frac{df_A(\mathbf{w}_\delta)}{d\delta} = 0$ , provided that  $w_k + \delta \in [\frac{-1}{m+1}, \frac{1}{m+1}]$ . The derivative of  $f_A$  with respect to the  $\delta$  is given by

$$\begin{aligned}
\frac{\partial f_A(\mathbf{w}_\delta)}{\partial \delta} &= \sum_{j=1}^m 2 \left( q - y_j \sum_{i=0}^m w_i G_{i,j} - y_j \delta G_{k,j} \right) (-y_j G_{k,j}) \\
&= 2 \sum_{j=1}^m \left[ \delta G_{k,j}^2 + y_j G_{k,j} \left( y_j \sum_{i=0}^m w_i G_{i,j} - q \right) \right] \\
&= 2 \left[ \delta \sum_{j=1}^m G_{k,j}^2 + \sum_{j=1}^m G_{k,j} \left( \sum_{i=0}^m w_i G_{i,j} - q y_j \right) \right] \\
(32) \quad &= 2 \left[ \delta \sum_{j=1}^m G_{k,j}^2 + \sum_{j=1}^m G_{k,j} D_{\mathbf{w}}(j) \right],
\end{aligned}$$

where  $D_{\mathbf{w}}(j) \stackrel{\text{def}}{=} \sum_{i=0}^m w_i G_{i,j} - q y_j$ .

Equation (32) implies that the optimal value for  $\delta$  is given by

$$(33) \quad \delta = - \frac{\sum_{j=1}^m G_{k,j} D_{\mathbf{w}}(j)}{\sum_{j=1}^m G_{k,j}^2}.$$

Algorithm 1 presents the complete optimization procedure that we have used.

---

**Algorithm 1** : PBSC-A optimization procedure

---

- 1: **Initialize:**  $w_i = 0 \quad \forall i \in \{0, \dots, m\}$  and  $D_{\mathbf{w}}(j) = -q y_j \quad \forall j \in \{1, \dots, m\}$ .
  - 2: **repeat**
  - 3:   Choose at random  $k \in \{0, \dots, m\}$ .
  - 4:   Compute  $\delta$  given by Equation (33).
  - 5:   If  $[w_k + \delta > \frac{1}{m+1}]$  then  $\delta \leftarrow \frac{1}{m+1} - w_k$ .
  - 6:   If  $[w_k + \delta < \frac{-1}{m+1}]$  then  $\delta \leftarrow \frac{-1}{m+1} - w_k$ .
  - 7:    $w_k \leftarrow w_k + \delta$ .
  - 8:   Update  $D_{\mathbf{w}}(j) \leftarrow D_{\mathbf{w}}(j) + \delta G_{k,j} \quad \forall j \in \{1, \dots, m\}$ .
  - 9: **until** Convergence
- 

### 3.2 PBSC-N: The non-aligned case

We consider the non-aligned scenario where the posterior  $Q$  is defined by a vector  $\mathbf{v} \stackrel{\text{def}}{=} (v_+, v_1, \dots, v_{2m}, v_-)$ :

$$\begin{aligned}
Q(h_{S_{\langle \rangle}^{(\varepsilon,+)}}) &= v_+, & Q(h_{S_{\langle i \rangle}^{(\sigma,+)}}) &= v_i \frac{1}{|\mathcal{M}^1|} I(\sigma \in \mathcal{M}^1), \\
Q(h_{S_{\langle \rangle}^{(\varepsilon,-)}}) &= v_-, & Q(h_{S_{\langle i \rangle}^{(\sigma,-)}}) &= v_{m+i} \frac{1}{|\mathcal{M}^1|} I(\sigma \in \mathcal{M}^1),
\end{aligned}$$

under the constraints  $v \geq 0$  for all  $v \in \mathbf{v}$  and  $\sum_{v \in \mathbf{v}} v = 1$ .

Lets compute the Kullback-Leibler divergence  $\text{KL}(Q\|P)$  between this posterior  $Q$  and the uniform prior  $P$ . We find that

$$\begin{aligned}
\text{KL}(Q\|P) &= \mathbf{E}_{h \sim Q} \ln \frac{Q(h)}{P(h)} \\
&= \sum_{i=1}^m \sum_{s \in \{-,+\}} \int_{\mathcal{M}^1} d\sigma Q(h_{S_{(i)}}^{(\sigma,s)}) \ln \left[ \frac{Q(h_{S_{(i)}}^{(\sigma,s)})}{P(h_{S_{(i)}}^{(\sigma,s)})} \right] + \sum_{s \in \{-,+\}} Q(h_{S_{(i)}}^{(\varepsilon,s)}) \ln \left[ \frac{Q(h_{S_{(i)}}^{(\varepsilon,s)})}{P(h_{S_{(i)}}^{(\varepsilon,s)})} \right] \\
&= \sum_{i=1}^{2m} \int_{\mathcal{M}^1} d\sigma \frac{v_i}{|\mathcal{M}^1|} \ln \left[ \frac{\frac{v_i}{|\mathcal{M}^1|}}{\frac{1}{2|\mathcal{M}^1|(m+1)}} \right] + v_+ \ln \frac{v_+}{\frac{1}{2(m+1)}} + v_- \ln \frac{v_-}{\frac{1}{2(m+1)}} \\
&= \sum_{v \in \mathbf{v}} v \ln \left[ \frac{v}{\frac{1}{2(m+1)}} \right] \\
&= \ln(2m+2) + \sum_{v \in \mathbf{v}} v \ln[v].
\end{aligned}$$

Moreover, we show in the main paper that the empirical risk  $\zeta_S^Q$  is given by

$$\zeta_S^Q = \frac{1}{mq^2} \sum_{j=1}^m \left( q - y_j \left[ v_+ - v_- + \sum_{i=1}^m (v_i - v_{i+m}) k(x_i, x_j) \right] \right)^2.$$

Minimizing the bound of Theorem 1 amounts at finding  $\mathbf{v}$  that minimizes

$$C \cdot \zeta_S^Q + \text{KL}(Q\|P),$$

for some constant  $C > 0$ . However, in Theorem 1, we have  $C = \frac{m \cdot C_1}{2\zeta^*(1)}$ .

Let  $v_0 \stackrel{\text{def}}{=} v_+$  and  $v_{2m+1} \stackrel{\text{def}}{=} v_-$ . Let us define a matrix  $\mathbf{G}$  of size  $2m+2 \times m$  as

$$G_{i,j} = \begin{cases} 1 & \text{if } i = 0, \\ k(\mathbf{x}_i, \mathbf{x}_j) & \text{if } 1 \leq i, j \leq m, \\ -k(\mathbf{x}_{i-m}, \mathbf{x}_j) & \text{if } m+1 \leq i \leq 2m \text{ (and } 1 \leq j \leq m), \\ -1 & \text{if } i = 2m+1. \end{cases}$$

With this notation, the optimization problem for  $PBSC-N$  can be written as

$$\begin{aligned}
\text{Minimize: } f_N(\mathbf{v}) &= \frac{C}{mq^2} \sum_{j=1}^m \left( q - y_j \sum_{i=0}^{2m+1} v_i G_{i,j} \right)^2 + \sum_{i=0}^{2m+1} v_i \ln v_i \\
\text{subject to: } &v_i \geq 0 \quad \text{for } i = 0, 1, \dots, 2m+1, \\
&\sum_{i=0}^{2m+1} v_i = 1.
\end{aligned}$$

We propose to solve this optimization problem by minimizing  $f_N$  with a coordinate-pair descent algorithm that works iteratively by exchanging weights between two components of  $\mathbf{v}$ . Starting from the uniform distribution  $P$  (i.e.,  $v_i = \frac{1}{2m+2}$  for  $i = 0, 1, \dots, 2m+1$ ), the learning algorithm iteratively chooses (at random)  $k, l \in \{0, \dots, 2m+1\}$  (with  $k \neq l$ ), and updates  $v_k \leftarrow v_k + \delta$  and  $v_l \leftarrow v_l - \delta$  (without updating the other weights) according to some optimally chosen value of  $\delta$ . Let  $\mathbf{v}_\delta$  be the new weight vector obtained

with such an update. After an update, the objective function becomes

$$f_N(\mathbf{v}_\delta) = \frac{C}{mq^2} \sum_{j=1}^m \left[ q - y_j \left( \sum_{i=0}^{2m+1} v_i G_{i,j} + \delta G_{k,j} - \delta G_{l,j} \right) \right]^2 + \sum_{i=0}^{2m+1} I(i \notin \{k, l\}) \cdot v_i \ln v_i + (v_k + \delta) \ln(v_k + \delta) + (v_l - \delta) \ln(v_l - \delta)$$

The optimal value for  $\delta$  is obtained when  $\frac{df_N(\mathbf{v}_\delta)}{d\delta} = 0$ , provided that  $v_k + \delta \in [0, v_k + v_l]$  and  $v_l - \delta \in [0, v_k + v_l]$ . The derivative of  $f_N$  with respect to the  $\delta$  is given by

$$\begin{aligned} \frac{\partial f_N(\mathbf{v}_\delta)}{\partial \delta} &= \frac{C}{mq^2} \sum_{j=1}^m 2 \left( q - y_j \sum_{i=0}^{2m+1} v_i G_{i,j} - y_j \delta (G_{k,j} - G_{l,j}) \right) (-y_j (G_{k,j} - G_{l,j})) + \ln \frac{v_k + \delta}{v_l - \delta} \\ &= \frac{2C}{mq^2} \sum_{j=1}^m \left[ \delta (G_{k,j} - G_{l,j})^2 + y_j (G_{k,j} - G_{l,j}) \left( y_j \sum_{i=0}^{2m+1} v_i G_{i,j} - q \right) \right] + \ln \frac{v_k + \delta}{v_l - \delta} \\ &= \frac{2C}{mq^2} \left[ \delta \sum_{j=1}^m (G_{k,j} - G_{l,j})^2 + \sum_{j=1}^m (G_{k,j} - G_{l,j}) \left( \sum_{i=0}^{2m+1} v_i G_{i,j} - q y_j \right) \right] + \ln \frac{v_k + \delta}{v_l - \delta} \\ (34) \quad &= \frac{2C}{mq^2} \left[ \delta \sum_{j=1}^m (G_{k,j} - G_{l,j})^2 + \sum_{j=1}^m (G_{k,j} - G_{l,j}) D_{\mathbf{v}}(j) \right] + \ln \frac{v_k + \delta}{v_l - \delta}, \end{aligned}$$

where  $D_{\mathbf{v}}(j) = \sum_{i=0}^{2m+1} v_i G_{i,j} - q y_j$ .

We find the optimal value for  $\delta$  with the help of a root finding method. Algorithm 2 presents the complete optimization procedure that we have used.

---

**Algorithm 2** : PBSC-N optimization procedure

---

- 1: **Initialize:**  $v_i = \frac{1}{2m+2} \quad \forall i \in \{0, \dots, m\}$  and  $D_{\mathbf{v}}(j) = -q y_j \quad \forall j \in \{1, \dots, m\}$ .
  - 2: **repeat**
  - 3:   Choose at random  $k, l \in \{0, \dots, 2m+1\}$  (with  $k \neq l$ ).
  - 4:   Find  $\delta$  given by the root of Equation (34).
  - 5:   If  $\delta > v_l$  then  $\delta \leftarrow v_l$ .
  - 6:   If  $\delta < -v_k$  then  $\delta \leftarrow -v_k$ .
  - 7:    $v_k \leftarrow v_k + \delta$  and  $v_l \leftarrow v_l - \delta$ .
  - 8:   Update  $D_{\mathbf{v}}(j) \leftarrow D_{\mathbf{v}}(j) + \delta (G_{k,j} - G_{l,j}) \quad \forall j \in \{1, \dots, m\}$ .
  - 9: **until** Convergence
-

## 4 Empirical Results

As mentioned in Section 2.4 of the main paper, this section presents all the empirical results that we have obtained in our experiments.

Table 1 of the main paper only presents the results for a subset of the data. Due to a lack of space, we have removed from Table 1 of the main paper the data sets having the smallest number of examples because they were less likely to show significative differences between the different algorithms. Here is the complete table of the results of our experiments.

Table 1: Empirical risk measured on the testing set  $T$  for the five different algorithms.

Dataset			Rbf kernel					Sigmoid kernel	
Name	$ T $	$ S $	SVM	RLSC	PBSC-A	PBSC-N	LINEAR	SVM	PBSC-A
Adult	10000	1809	0.158	0.157	<b>0.156</b>	0.160	0.193	0.163	<b>0.157</b>
BreastC	340	343	<b>0.038</b>	<b>0.038</b>	0.044	<b>0.038</b>	0.144	<b>0.038</b>	<b>0.038</b>
Credit-A	300	353	0.190	0.160	<b>0.140</b>	0.173	0.200	0.190	<b>0.170</b>
Glass	107	107	<b>0.150</b>	<b>0.150</b>	<b>0.150</b>	0.168	0.187	<b>0.355</b>	0.411
Haberman	150	144	<b>0.267</b>	0.327	0.280	<b>0.267</b>	<b>0.267</b>	<b>0.273</b>	<b>0.273</b>
Heart	147	150	<b>0.197</b>	0.211	0.204	0.218	0.238	<b>0.184</b>	0.197
Ionosphere	175	176	0.057	<b>0.023</b>	0.040	0.040	0.326	0.126	<b>0.091</b>
Letter:AB	1055	500	<b>0.001</b>	0.002	<b>0.001</b>	<b>0.001</b>	0.038	0.009	<b>0.005</b>
Letter:DO	1058	500	0.014	0.015	<b>0.011</b>	0.012	0.069	<b>0.022</b>	0.028
Letter:OQ	1036	500	0.016	<b>0.011</b>	0.016	0.014	0.123	<b>0.018</b>	0.039
Liver	175	170	0.286	0.291	<b>0.280</b>	0.286	0.349	<b>0.400</b>	<b>0.400</b>
Mnist:0vs8	1916	500	<b>0.003</b>	0.009	0.004	0.004	0.031	0.007	<b>0.003</b>
Mnist:1vs7	1922	500	0.014	0.008	0.008	<b>0.007</b>	0.161	0.012	<b>0.007</b>
Mnist:1vs8	1936	500	0.011	<b>0.010</b>	<b>0.010</b>	0.011	0.292	<b>0.014</b>	0.015
Mnist:2vs3	1905	500	0.020	0.022	<b>0.019</b>	0.020	0.114	<b>0.025</b>	0.031
Mushroom	4062	4062	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	0.022	<b>0.000</b>	0.010
Ringnorm	3700	3700	0.015	0.017	<b>0.013</b>	<b>0.013</b>	0.103	<b>0.020</b>	0.035
sonar	104	104	0.154	0.250	<b>0.125</b>	0.192	0.490	0.250	<b>0.183</b>
Tic-tac-toe	479	479	<b>0.015</b>	0.019	0.019	0.052	0.365	<b>0.023</b>	0.159
Usvotes	200	235	0.075	<b>0.065</b>	<b>0.065</b>	<b>0.065</b>	0.140	0.070	<b>0.065</b>
Waveform	4000	4000	0.068	0.067	0.068	<b>0.066</b>	0.143	<b>0.067</b>	<b>0.067</b>
Wdbc	284	285	<b>0.042</b>	0.067	0.049	0.074	0.180	<b>0.366</b>	<b>0.366</b>

Note that Table 2 below is exactly the same as in the main paper.

Table 2: Mean and standard deviation (in parentheses) of the empirical risk across 20 partitions.

Dataset	Linear SVM		k-NN		PBSC-A	
Aural Sonar	<b>0.1425</b>	(0.694)	0.1825	(0.597)	0.1500	(0.827)
Voting	0.0534	(0.193)	0.0546	(0.174)	<b>0.0529</b>	(0.184)
Yeast-5-7	<b>0.2688</b>	(0.622)	0.3063	(0.580)	0.2975	(0.668)
Yeast-5-12	<b>0.1075</b>	(0.482)	0.1275	(0.439)	0.1088	(0.598)

## 5 Another PAC-Bayes bound without $KL(Q\|P)$ (not stated in main the paper)

The following theorem can be viewed as a generalization of PAC-Bayes bound of Seeger (2002) to our setting. As for Theorem 5, the following version is free of any Kullback-Leibler divergence. This bound is tighter than the one in Theorem 5. We have, nevertheless, decided to state Theorem 5 in the main paper because it has a simpler statement and because we have  $l \cdot \deg \zeta = 2$  for the proposed learning algorithms—a case where the bound of Theorem 5 is already quite tight.

**Theorem 2.** *For any  $D$ , for any family  $(\mathcal{H}^S)_{S \in D^m}$  of sets of sc-classifiers of size at most  $l$ , for any prior  $P$ , for any margin loss function  $\zeta$  such that  $l \cdot \deg(\zeta) < m$ , and for for any  $\delta \in (0, 1]$ , we have*

$$\Pr_{S \sim D^m} \left( \begin{array}{l} \forall Q \text{ aligned on } P : \\ \text{kl}^+ \left( \frac{m}{m-l \cdot \deg \zeta} \left[ \frac{1}{2} \left( 1 + \frac{1}{\zeta(1)} \zeta_S^Q \right) + \frac{ld}{m} \right] \left\| \frac{1}{2} \left[ 1 + \frac{1}{\zeta(1)} \zeta_D^Q \right] \right) \leq \frac{\ln \frac{2\sqrt{m}}{\delta}}{m-l \cdot \deg \zeta} \end{array} \right) \geq 1 - \delta$$

where  $\text{kl}(q\|p) \stackrel{\text{def}}{=} q \ln \frac{q}{p} + (1-q) \ln \frac{1-q}{1-p}$ , and where  $\text{kl}^+(q\|p) = \text{kl}(q\|p)$  if  $q \geq p$  and 0 otherwise. Moreover, if  $l = 0$ ,  $\text{kl}^+(\cdot)$  can be replaced by  $\text{kl}(\cdot)$  in the above statement—thus giving rise to both a lower and an upper bound for  $\zeta_D^Q$ .

*Proof.* The first part of the proof is very similar to the one of Theorem 5. We thus use here the same definitions for  $d$ ,  $\bar{h} = \overline{h_1..h_k}$  (with  $k \in \{0, \dots, d\}$ ),  $\bar{R}_D(\bar{h})$ ,  $\bar{R}_S(\bar{h})$ ,  $\bar{\mathcal{H}}^S$ ,  $\bar{P}$ ,  $\bar{Q}$ ,  $\zeta_D^Q$ , and  $\zeta_S^Q$ . However, we will instead consider the following (and quite different) Laplace transform

$$(35) \quad X_{\bar{P}} \stackrel{\text{def}}{=} \mathbf{E}_{\bar{h} \sim \bar{P}} e^{(m-|\mathbf{i}_{\bar{h}}|)\text{kl}(\bar{R}_S(\bar{h}), \bar{R}_D(\bar{h}))},$$

where  $\tilde{R}_S(\bar{h})$  is the abstract empirical risk computed on the examples of  $S$  that are not in the compression sequence of  $\bar{h}$ . More formally,

$$\tilde{R}_S(\overline{h_1..h_k}) \stackrel{\text{def}}{=} \frac{1}{m - |\mathbf{i}_{h_1..h_k}|} \sum_{j=1}^m I(\neg \bigvee_{i=1}^k (h_i(x_j) \neq y_j)) I((x_j, y_j) \notin \mathbf{i}_{h_1..h_k}).$$

As in the proof of Theorem 5, we can show the following claim.

**Claim :** for any posterior  $Q$  aligned on  $P$ , we have

$$(36) \quad X_{\bar{P}} = \mathbf{E}_{\bar{h} \sim \bar{Q}} e^{(m-|\mathbf{i}_{\bar{h}}|)\text{kl}(\bar{R}_S(\bar{h}), \bar{R}_D(\bar{h}))}.$$

Now again, as in the proof of Theorem 5, by Markov's inequality, we have

$$\Pr_{S \sim D^m} \left( X_{\bar{P}} \leq \frac{1}{\delta} \mathbf{E}_{S \sim D^m} X_{\bar{P}} \right) \geq 1 - \delta.$$

Thus, by applying the claim and by taking the logarithm on each side of the innermost inequality, we obtain

$$(37) \quad \Pr_{S \sim D^m} \left( \begin{array}{l} \forall Q \text{ aligned on } P : \\ \ln \left[ \mathbf{E}_{\bar{h} \sim \bar{Q}} e^{(m-|\mathbf{i}_{\bar{h}}|)\text{kl}(\bar{R}_S(\bar{h}), \bar{R}_D(\bar{h}))} \right] \leq \ln \left[ \frac{1}{\delta} \mathbf{E}_{S \sim D^m} \mathbf{E}_{\bar{h} \sim \bar{P}} e^{(m-|\mathbf{i}_{\bar{h}}|)\text{kl}(\bar{R}_S(\bar{h}), \bar{R}_D(\bar{h}))} \right] \end{array} \right) \geq 1 - \delta.$$

Jensen's inequality applied to the concave  $\ln(x)$  gives

$$(38) \quad \ln \left[ \mathbf{E}_{\bar{h} \sim \bar{Q}} e^{(m-|\mathbf{i}_{\bar{h}}|) \text{kl}(\tilde{R}_S(\bar{h}) \| \bar{R}_D(\bar{h}))} \right] \geq \mathbf{E}_{\bar{h} \sim \bar{Q}} (m - |\mathbf{i}_{\bar{h}}|) \text{kl}(\tilde{R}_S(\bar{h}) \| \bar{R}_D(\bar{h})).$$

Again from the Jensen's inequality, applied to the convex function  $\text{kl}(\cdot \| \cdot)$ , together with the definitions of  $\zeta_D^Q$  and  $\zeta_S^Q$  (see Equation (26)) and the fact that  $m - |\mathbf{i}_{\bar{h}}| \geq m - l \cdot d$ , we obtain

$$(39) \quad \mathbf{E}_{\bar{h} \sim \bar{Q}} (m - |\mathbf{i}_{\bar{h}}|) \text{kl}(\tilde{R}_S(\bar{h}) \| \bar{R}_D(\bar{h})) \geq (m - ld) \text{kl} \left( \mathbf{E}_{\bar{h} \sim \bar{Q}} \tilde{R}_S(\bar{h}) \| \mathbf{E}_{\bar{h} \sim \bar{Q}} \bar{R}_D(\bar{h}) \right)$$

Let us now analyse the value of  $\mathbf{E}_{S \sim D^m} X_{\bar{P}}$ . Let  $\mathbf{i}^c$  be the vector of indices of  $\mathcal{I}$  that are not in the vector  $\mathbf{i}$ , and note that

$$\mathbf{E}_{S \sim D^m} \mathbf{E}_{\bar{h} \sim \bar{P}} e^{(m-|\mathbf{i}_{\bar{h}}|) \text{kl}(\tilde{R}_S(\bar{h}), \bar{R}_D(\bar{h}))} = \mathbf{E}_{\mathbf{i} \sim \bar{P}_{\mathcal{I}}} \mathbf{E}_{S_{\mathbf{i}} \sim D^{|\mathbf{i}|}} \mathbf{E}_{\mu \sim \bar{P}_{S_{\mathbf{i}}}} \mathbf{E}_{S_{\mathbf{i}^c} \sim D^{m-|\mathbf{i}|}} e^{|\mathbf{i}^c| \text{kl}(\tilde{R}_S(\bar{h}_{\mathbf{i}^c}^{\mu}), \bar{R}_D(\bar{h}_{\mathbf{i}^c}^{\mu}))}.$$

Since  $\tilde{R}_S(\bar{h}_{\mathbf{i}^c}^{\mu})$  is an arithmetic mean of iid random variables, one can apply Lemma 0 with  $M(X)$  replaced by  $\tilde{R}_S(\bar{h}_{\mathbf{i}^c}^{\mu})$ ,  $n$  replaced by  $m - |\mathbf{i}|$ , and  $\nu$  replaced by  $\bar{R}_D(\bar{h}_{\mathbf{i}^c}^{\mu})$  to obtain

$$(40) \quad \mathbf{E}_{S_{\mathbf{i}^c} \sim D^{m-|\mathbf{i}|}} e^{(m-|\mathbf{i}|) \text{kl}(\tilde{R}_S(\bar{h}_{\mathbf{i}^c}^{\mu}), \bar{R}_D(\bar{h}_{\mathbf{i}^c}^{\mu}))} \leq 2\sqrt{m - |\mathbf{i}|} \leq 2\sqrt{m}.$$

By rearranging Equation (37), and by using Equations (40), (38) and (39), we have

$$(41) \quad \mathbf{E}_{S \sim D^m} (m - ld) \text{kl} \left( \mathbf{E}_{\bar{h} \sim \bar{Q}} \tilde{R}_S(\bar{h}) \| \mathbf{E}_{\bar{h} \sim \bar{Q}} \bar{R}_D(\bar{h}) \right) \leq \ln \frac{2\sqrt{m}}{\delta}.$$

Finally, observe that for any classifier  $\bar{h} \in \bar{\mathcal{H}}^S$ , we have

$$(42) \quad \begin{aligned} \tilde{R}_S(\bar{h}) &\leq \left( \bar{R}_S(\bar{h}) + \frac{ld}{m} \right) \frac{m}{m - |\mathbf{i}|} \\ &\leq \left( \bar{R}_S(\bar{h}) + \frac{ld}{m} \right) \frac{m}{m - ld}. \end{aligned}$$

Consider the following two cases.

*case 1* :  $l = 0$ . In that case we have  $\mathbf{E}_{\bar{h} \sim \bar{Q}} \tilde{R}_S(\bar{h}) = \mathbf{E}_{\bar{h} \sim \bar{Q}} \bar{R}_S(\bar{h})$ . Hence we have

$$\text{kl} \left( \mathbf{E}_{\bar{h} \sim \bar{Q}} \tilde{R}_S(\bar{h}) \| \mathbf{E}_{\bar{h} \sim \bar{Q}} \bar{R}_D(\bar{h}) \right) = \text{kl} \left( \frac{m}{m - ld} \mathbf{E}_{\bar{h} \sim \bar{Q}} \bar{R}_S(\bar{h}) + \frac{ld}{m} \| \mathbf{E}_{\bar{h} \sim \bar{Q}} \bar{R}_D(\bar{h}) \right)$$

*case 2* :  $l > 0$ . In that case, following Equation (??), we can show that

$$\text{kl} \left( \mathbf{E}_{\bar{h} \sim \bar{Q}} \tilde{R}_S(\bar{h}) \| \mathbf{E}_{\bar{h} \sim \bar{Q}} \bar{R}_D(\bar{h}) \right) \geq \text{kl}^+ \left( \frac{m}{m - ld} \mathbf{E}_{\bar{h} \sim \bar{Q}} \bar{R}_S(\bar{h}) + \frac{ld}{m} \| \mathbf{E}_{\bar{h} \sim \bar{Q}} \bar{R}_D(\bar{h}) \right)$$

In each case, the result then follows from Equation (26), Equation (41) and straightforward calculations.  $\square$

## References

- Germain, Pascal, Lacoste, Alexandre, Laviolette, François, Marchand, Mario, and Shanian, Sara. A PAC-Bayes Sample Compression Approach to Kernel Methods. In *Proceedings of the 28th International Conference on Machine Learning*, Bellevue, WA, USA, June 2011. URL <http://www.icml-2011.org/papers.php>.
- Maurer, Andreas. A note on the pac bayesian theorem. *CoRR*, cs.LG/0411099, 2004.
- Schapire, Robert E., Freund, Yoav, Bartlett, Peter, and Lee, Wee Sun. Boosting the margin: A new explanation for the effectiveness of voting methods. *The Annals of Statistics*, 26:1651–1686, 1998.
- Seeger, Matthias. PAC-Bayesian generalization bounds for gaussian processes. *Journal of Machine Learning Research*, 3:233–269, 2002.