
Supplementary Material to A PAC-Bayesian Approach for Domain Adaptation with Specialization to Linear Classifiers

Pascal Germain

PASCAL.GERMAIN@IFT.ULAVAL.CA

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

Amaury Habrard

AMAURY.HABRARD@UNIV-ST-ETIENNE.FR

Laboratoire Hubert Curien UMR CNRS 5516, Université Jean Monnet, 42000 St-Etienne, France

François Laviolette

FRANCOIS.LAVIOLETTE@IFT.ULAVAL.CA

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

Emilie Morvant

EMILIE.MORVANT@LIF.UNIV-MRS.FR

Aix-Marseille Univ., LIF-QARMA, CNRS, UMR 7279, 13013, Marseille, France

In this document, Section 1 contains some lemmas used in subsequent proofs, Section 2 presents an extended proof of the bound on the domain disagreement $\text{dis}_\rho(D_S, D_T)$ (Theorem 3 of the main paper), Section 3 introduces other PAC-Bayesian bounds for $\text{dis}_\rho(D_S, D_T)$ and $R_{P_T}(G_\rho)$, Section 4 shows equations and implementation details about PBDA (our proposed learning algorithm for PAC-Bayesian DA tasks).

1. Some tools

Lemma 1 (Markov's inequality). *Let Z be a random variable and $t \geq 0$, then,*

$$P(|Z| \geq t) \leq \mathbf{E} (|Z|) / t.$$

Lemma 2 (Jensen's inequality). *Let Z be an integrable real-valued random variable and $g(\cdot)$ any function.*

If $g(\cdot)$ is convex, then,

$$g(\mathbf{E} [Z]) \leq \mathbf{E} [g(Z)].$$

If $g(\cdot)$ is concave, then,

$$g(\mathbf{E} [Z]) \geq \mathbf{E} [g(Z)].$$

Lemma 3 (Maurer (2004)). *Let $X = (X_1, \dots, X_m)$ be a vector of i.i.d. random variables, $0 \leq X_i \leq 1$, with $\mathbf{E} X_i = \mu$. Denote $X' = (X'_1, \dots, X'_m)$, where X'_i is the unique Bernoulli ($\{0, 1\}$ -valued) random variable with $\mathbf{E} X'_i = \mu$. If $f : [0, 1]^n \rightarrow \mathbb{R}$ is convex, then,*

$$\mathbf{E} [f(X)] \leq \mathbf{E} [f(X')].$$

Lemma 4 (from Inequalities (1) and (2) of Maurer (2004)). *Let $m \geq 8$, and $X = (X_1, \dots, X_m)$ be a vector of i.i.d. random variables, $0 \leq X_i \leq 1$. Then,*

$$\sqrt{m} \leq \mathbf{E} \exp \left(m \text{kl} \left(\frac{1}{m} \sum_{i=1}^n X_i \parallel \mathbf{E} [X_i] \right) \right) \leq 2\sqrt{m},$$

where, $\text{kl}(a \parallel b) \stackrel{\text{def}}{=} a \ln \frac{a}{b} + (1-a) \ln \frac{1-a}{1-b}$. (7)

2. Detailed Proof of Theorem 3

We recall the Theorem 3 of the main paper.

Theorem 3. *For any distributions D_S and D_T over X , any set of hypothesis \mathcal{H} , any prior distribution π over \mathcal{H} , any $\delta \in (0, 1]$, and any real number $\alpha > 0$, with a probability at least $1 - \delta$ over the choice of $S \times T \sim (D_S \times D_T)^m$, for every ρ on \mathcal{H} , we have,*

$$\text{dis}_\rho(D_S, D_T) \leq \frac{2\alpha \left[\text{dis}_\rho(S, T) + \frac{2\text{KL}(\rho \parallel \pi) + \ln \frac{2}{\delta}}{m \times \alpha} + 1 \right] - 1}{1 - e^{-2\alpha}},$$

where $\text{dis}_\rho(S, T)$ is the empirical domain disagreement.

Proof. Firstly, we propose to upper-bound,

$$d^{(1)} \stackrel{\text{def}}{=} \mathbf{E}_{(h, h') \sim \rho^2} [R_{D_S}(h, h') - R_{D_T}(h, h')],$$

by its empirical counterpart,

$$d_{S \times T}^{(1)} \stackrel{\text{def}}{=} \mathbf{E}_{(h, h') \sim \rho^2} [R_S(h, h') - R_T(h, h')].$$

and some extra terms related to the Kullback-Leibler divergence between the posterior and the prior.

To do that, we consider an “abstract” classifier $\hat{h} \stackrel{\text{def}}{=} (h, h') \in \mathcal{H}^2$ chosen according a distribution $\hat{\rho}$, with $\hat{\rho}(\hat{h}) = \rho(h)\rho(h')$. Notice that with $\hat{\pi}(\hat{h}) = \pi(h)\pi(h')$, we obtain that $\text{KL}(\hat{\rho}||\hat{\pi}) = 2\text{KL}(\rho||\pi)$,

$$\begin{aligned} \text{KL}(\hat{\rho}||\hat{\pi}) &= \mathbf{E}_{(h, h') \sim \rho^2} \ln \frac{\rho(h)\rho(h')}{\pi(h)\pi(h')} \\ &= \mathbf{E}_{h \sim \rho} \ln \frac{\rho(h)}{\pi(h)} + \mathbf{E}_{h' \sim \rho} \ln \frac{\rho(h')}{\pi(h')} \\ &= 2 \mathbf{E}_{h \sim \rho} \ln \frac{\rho(h)}{\pi(h)} = 2\text{KL}(\rho||\pi). \end{aligned} \quad (8)$$

Let us define the “abstract” loss of \hat{h} on a pair of examples $(\mathbf{x}^s, \mathbf{x}^t) \sim D_{S \times T} = D_S \times D_T$ by,

$$\mathcal{L}_{d^{(1)}}(\hat{h}, \mathbf{x}^s, \mathbf{x}^t) \stackrel{\text{def}}{=} \frac{1 + \mathcal{L}_{0.1}(h(\mathbf{x}^s), h'(\mathbf{x}^s)) - \mathcal{L}_{0.1}(h(\mathbf{x}^t), h'(\mathbf{x}^t))}{2}.$$

Therefore, the “abstract” risk of \hat{h} on the joint distribution is defined as,

$$R_{D_{S \times T}}^{(1)}(\hat{h}) = \mathbf{E}_{\mathbf{x}^s \sim D_S} \mathbf{E}_{\mathbf{x}^t \sim D_T} \mathcal{L}_{d^{(1)}}(\hat{h}, \mathbf{x}^s, \mathbf{x}^t),$$

and the error of the related Gibbs classifier associated with this loss is,

$$R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}}) = \mathbf{E}_{\hat{h} \sim \hat{\rho}} R_{D_{S \times T}}^{(1)}(\hat{h}).$$

The empirical counterparts of these two quantities are,

$$R_{S \times T}^{(1)}(\hat{h}) = \mathbf{E}_{(\mathbf{x}^s, \mathbf{x}^t) \sim S \times T} \mathcal{L}_{d^{(1)}}(\hat{h}, \mathbf{x}^s, \mathbf{x}^t)$$

and,

$$R_{S \times T}^{(1)}(G_{\hat{\rho}}) = \mathbf{E}_{\hat{h} \sim \hat{\rho}} R_{S \times T}^{(1)}(\hat{h}).$$

It is easy to show that,

$$d^{(1)} = 2R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}}) - 1, \quad (9)$$

$$d_{S \times T}^{(1)} = 2R_{S \times T}^{(1)}(G_{\hat{\rho}}) - 1. \quad (10)$$

As $\mathcal{L}_{d^{(1)}}$ lies in $[0, 1]$, we can bound the true $R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}})$ following the proof process of Th. 2 of the main paper (with $c = 2\alpha$). To do so, we define the convex function,

$$\mathcal{F}(p) \stackrel{\text{def}}{=} -\ln[1 - (1 - e^{-2\alpha})p], \quad (11)$$

and consider the non-negative random variable,

$$\mathbf{E}_{\hat{h} \sim \hat{\pi}} e^{m(\mathcal{F}(R_{D_{S \times T}}^{(1)}(\hat{h})) - 2\alpha R_{S \times T}^{(1)}(\hat{h}))}.$$

We apply Markov’s inequality (Lemma 1 of this Supp. Material). For every $\delta \in (0, 1]$, with a probability at

least $1 - \delta$ over the choice of $S \times T \sim (D_{S \times T})^m$, we have,

$$\begin{aligned} &\mathbf{E}_{\hat{h} \sim \hat{\pi}} e^{m(\mathcal{F}(R_{D_{S \times T}}^{(1)}(\hat{h})) - 2\alpha R_{S \times T}^{(1)}(\hat{h}))} \\ &\leq \frac{1}{\delta} \mathbf{E}_{S \times T \sim (D_{S \times T})^m} \mathbf{E}_{\hat{h} \sim \hat{\pi}} e^{m(\mathcal{F}(R_{D_{S \times T}}^{(1)}(\hat{h})) - 2\alpha R_{S \times T}^{(1)}(\hat{h}))}. \end{aligned}$$

By taking the logarithm on each side of the previous inequality, and transforming the expectation over $\hat{\pi}$ into an expectation over $\hat{\rho}$, we obtain that,

$$\begin{aligned} &\ln \left[\mathbf{E}_{\hat{h} \sim \hat{\rho}} \frac{\hat{\pi}(\hat{h})}{\hat{\rho}(\hat{h})} e^{m(\mathcal{F}(R_{D_{S \times T}}^{(1)}(\hat{h})) - 2\alpha R_{S \times T}^{(1)}(\hat{h}))} \right] \\ &\leq \ln \left[\frac{1}{\delta} \mathbf{E}_{S \times T \sim (D_{S \times T})^m} \mathbf{E}_{\hat{h} \sim \hat{\pi}} e^{m(\mathcal{F}(R_{D_{S \times T}}^{(1)}(\hat{h})) - 2\alpha R_{S \times T}^{(1)}(\hat{h}))} \right] \\ &= \ln \left[\frac{1}{\delta} \mathbf{E}_{\hat{h} \sim \hat{\rho}} e^{m\mathcal{F}(R_{D_{S \times T}}^{(1)}(\hat{h}))} \mathbf{E}_{S \times T \sim (D_{S \times T})^m} e^{-2m\alpha R_{S \times T}^{(1)}(\hat{h})} \right]. \end{aligned} \quad (12)$$

For a classifier \hat{h} , let us define a random variable $X_{\hat{h}}$ that follows a binomial distribution of m trials with a probability of success $R_{D_{S \times T}}^{(1)}(\hat{h})$ denoted by $B(m, R_{D_{S \times T}}^{(1)}(\hat{h}))$. Lemma 3 gives,

$$\begin{aligned} &\mathbf{E}_{S \times T \sim (D_{S \times T})^m} e^{-2m\alpha R_{S \times T}^{(1)}(\hat{h})} \\ &\leq \mathbf{E}_{X_{\hat{h}} \sim B(m, R_{D_{S \times T}}^{(1)}(\hat{h}))} e^{-2\alpha X_{\hat{h}}} \\ &= \sum_{k=0}^m \Pr_{X_{\hat{h}} \sim B(m, R_{D_{S \times T}}^{(1)}(\hat{h}))} (X_{\hat{h}} = k) e^{-2\alpha k} \\ &= \sum_{k=0}^m \binom{m}{k} (R_{S \times T}^{(1)}(\hat{h}))^k (1 - R_{S \times T}^{(1)}(\hat{h}))^{m-k} e^{-2\alpha k} \\ &= \sum_{k=0}^m \binom{m}{k} (R_{S \times T}^{(1)}(\hat{h}) e^{-2\alpha})^k (1 - R_{S \times T}^{(1)}(\hat{h}))^{m-k} \\ &= \left[R_{S \times T}^{(1)}(\hat{h}) e^{-2\alpha} + (1 - R_{S \times T}^{(1)}(\hat{h})) \right]^m. \end{aligned}$$

The last line result, together with the choice of \mathcal{F} (Eq. (11)), leads to,

$$\begin{aligned} &\mathbf{E}_{\hat{h} \sim \hat{\pi}} e^{m\mathcal{F}(R_{D_{S \times T}}^{(1)}(\hat{h}))} \mathbf{E}_{S \times T \sim (D_{S \times T})^m} e^{-2m\alpha R_{S \times T}^{(1)}(\hat{h})} \\ &\leq \mathbf{E}_{\hat{h} \sim \hat{\pi}} e^{m\mathcal{F}(R_{D_{S \times T}}^{(1)}(\hat{h}))} \left[R_{S \times T}^{(1)}(\hat{h}) e^{-2\alpha} + (1 - R_{S \times T}^{(1)}(\hat{h})) \right]^m \\ &= \mathbf{E}_{\hat{h} \sim \hat{\pi}} 1 = 1. \end{aligned}$$

We can now upper bound Eq. (12) simply by,

$$\ln \left[\mathbf{E}_{\hat{h} \sim \hat{\rho}} \frac{\hat{\pi}(\hat{h})}{\hat{\rho}(\hat{h})} e^{m(\mathcal{F}(R_{D_{S \times T}}^{(1)}(\hat{h})) - 2\alpha R_{S \times T}^{(1)}(\hat{h}))} \right] \leq \ln \frac{1}{\delta}.$$

Let us insert the term $\text{KL}(\rho\|\pi)$ in the left-hand side of the last inequality and find a lower bound by using Jensen's inequality (Lemma 2) twice, first on the concave logarithm function and then on the convex function \mathcal{F} ,

$$\begin{aligned} & \ln \left[\mathbf{E}_{\substack{\hat{h} \sim \hat{\rho} \\ \hat{\rho}(\hat{h})}} \frac{\hat{\pi}(\hat{h})}{\hat{\rho}(\hat{h})} e^{m(\mathcal{F}(R_{D_{S \times T}}^{(1)}(\hat{h})) - 2\alpha R_{S \times T}^{(1)}(\hat{h}))} \right] \\ &= \ln \left[\mathbf{E}_{\hat{h} \sim \hat{\rho}} e^{m(\mathcal{F}(R_{D_{S \times T}}^{(1)}(\hat{h})) - 2\alpha R_{S \times T}^{(1)}(\hat{h}))} \right] - 2\text{KL}(\rho\|\pi) \\ &\geq \mathbf{E}_{\hat{h} \sim \hat{\rho}} m \left(\mathcal{F}(R_{D_{S \times T}}^{(1)}(\hat{h})) - 2\alpha R_{S \times T}^{(1)}(\hat{h}) \right) - 2\text{KL}(\rho\|\pi) \\ &\geq m\mathcal{F}(\mathbf{E}_{\hat{h} \sim \hat{\rho}} R_{D_{S \times T}}^{(1)}(\hat{h})) - 2m\alpha \mathbf{E}_{\hat{h} \sim \hat{\rho}} R_{S \times T}^{(1)}(\hat{h}) - 2\text{KL}(\rho\|\pi) \\ &= m\mathcal{F}(R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}})) - 2m\alpha R_{S \times T}^{(1)}(G_{\hat{\rho}}) - 2\text{KL}(\rho\|\pi). \end{aligned}$$

We then have,

$$m\mathcal{F}(\mathbf{E}_{\hat{h} \sim \hat{\rho}} R_{D_{S \times T}}^{(1)}(\hat{h})) - 2m\alpha \mathbf{E}_{\hat{h} \sim \hat{\rho}} R_{S \times T}^{(1)}(\hat{h}) - 2\text{KL}(\rho\|\pi) \leq \ln \frac{1}{\delta}.$$

This, in turn, implies that,

$$\mathcal{F}(R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}})) \leq 2\alpha R_{S \times T}^{(1)}(G_{\hat{\rho}}) + \frac{2\text{KL}(\rho\|\pi) + \ln \frac{1}{\delta}}{m}.$$

Now, by isolating $R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}})$, we obtain,

$$R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}}) \leq \frac{1}{1 - e^{-2\alpha}} \left[1 - e^{-\left(2\alpha R_{S \times T}^{(1)}(G_{\hat{\rho}}) + \frac{2\text{KL}(\rho\|\pi) + \ln \frac{1}{\delta}}{m}\right)} \right],$$

and from the inequality $1 - e^{-x} \leq x$,

$$R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}}) \leq \frac{1}{1 - e^{-2\alpha}} \left[2\alpha R_{S \times T}^{(1)}(G_{\hat{\rho}}) + \frac{2\text{KL}(\rho\|\pi) + \ln \frac{1}{\delta}}{m} \right].$$

It then follows from Equations (9) and (10) that, with probability at least $1 - \frac{\delta}{2}$ over the choice of $S \times T \sim (D_S \times D_T)^m$, we have,

$$\frac{d^{(1)} + 1}{2} \leq \frac{2\alpha}{1 - e^{-2\alpha}} \left[\frac{d_{S \times T}^{(1)} + 1}{2} + \frac{2\text{KL}(\rho\|\pi) + \ln \frac{1}{\delta}}{m \times 2\alpha} \right],$$

We now bound $d^{(2)} \stackrel{\text{def}}{=} \mathbf{E}_{(h, h') \sim \rho^2} [R_{D_T}(h, h') - R_{D_S}(h, h')]$

using exactly the same argument as for $d^{(1)}$ except that we instead consider the following ‘‘abstract’’ loss of \hat{h} on a pair of examples $(\mathbf{x}^s, \mathbf{x}^t) \sim D_{S \times T} = D_S \times D_T$:

$$\mathcal{L}_{d^{(1)}}(\hat{h}, \mathbf{x}^s, \mathbf{x}^t) \stackrel{\text{def}}{=} \frac{1 + \mathcal{L}_{0-1}(h(\mathbf{x}^t), h'(\mathbf{x}^t)) - \mathcal{L}_{0-1}(h(\mathbf{x}^s), h'(\mathbf{x}^s))}{2}.$$

We then obtain that, with probability at least $1 - \frac{\delta}{2}$ over the choice of $S \times T \sim (D_S \times D_T)^m$,

$$\frac{d^{(2)} + 1}{2} \leq \frac{2\alpha}{1 - e^{-2\alpha}} \left[\frac{d_{S \times T}^{(2)} + 1}{2} + \frac{2\text{KL}(\rho\|\pi) + \ln \frac{1}{\delta}}{m \times 2\alpha} \right].$$

To finish the proof, note that by definition, we have that $d^{(1)} = -d^{(2)}$, hence

$$|d^{(1)}| = |d^{(2)}| = \text{dis}_{\rho}(D_S, D_T),$$

and,

$$|d_{S \times T}^{(1)}| = |d_{S \times T}^{(2)}| = \text{dis}_{\rho}(S, T).$$

Then, the maximum of the bound on $d^{(1)}$ and the bound on $d^{(2)}$ gives a bound on $\text{dis}_{\rho}(D_S, D_T)$.

Finally, by the union bound, we have that, with probability $1 - \delta$ over the choice of $S \times T \sim (D_S \times D_T)^m$, we have,

$$\frac{|d^{(1)}| + 1}{2} \leq \frac{\alpha}{1 - e^{-2\alpha}} \left[|d_{S \times T}^{(1)}| + 1 + \frac{2\text{KL}(\rho\|\pi) + \ln \frac{2}{\delta}}{m \times \alpha} \right],$$

or, which is equivalent,

$$\text{dis}_{\rho}(D_S, D_T) \leq \frac{2\alpha \left[\text{dis}_{\rho}(S, T) + \frac{2\text{KL}(\rho\|\pi) + \ln \frac{2}{\delta}}{m \times \alpha} + 1 \right] - 1}{1 - e^{-2\alpha}},$$

and we are done. \square

3. Other PAC-Bayesian Bounds

3.1. PAC-Bayesian Bounds with the kl term

Let us recall the PAC-Bayesian bound proposed by Seeger (2002), in which the trade-off between the complexity and the risk is handled by the kl function defined by Equation (7) in this supplementary materials.

Theorem 6 (Seeger (2002)). *For any domain P_S over $X \times Y$, any set of hypothesis \mathcal{H} , and any prior distribution π over \mathcal{H} , any $\delta \in (0, 1]$, with a probability at least $1 - \delta$ over the choice of $S \sim (P_S)^m$, for every ρ over \mathcal{H} , we have,*

$$\text{kl}\left(R_S(G_{\rho}) \parallel R_{P_S}(G_{\rho})\right) \leq \frac{1}{m} \left[\text{KL}(\rho\|\pi) + \ln \frac{2\sqrt{m}}{\delta} \right].$$

Here is a ‘‘Seeger’s type’’ PAC-Bayesian bound for our domain disagreement dis_{ρ} .

Theorem 7. *For any distributions D_S and D_T over X , any set of hypothesis \mathcal{H} , and any prior distribution π over \mathcal{H} , any $\delta \in (0, 1]$, with a probability at least $1 - \delta$ over the choice of $S \times T \sim (D_S \times D_T)^m$, for every ρ on \mathcal{H} , we have,*

$$\text{kl}\left(\frac{\text{dis}_{\rho}(S, T) + 1}{2} \parallel \frac{\text{dis}_{\rho}(D_S, D_T) + 1}{2}\right) \leq \frac{1}{m} \left[2\text{KL}(\rho\|\pi) + \ln \frac{2\sqrt{m}}{\delta} \right].$$

Proof. Similarly as in the proof of Theorem 3, we will first bound,

$$d^{(1)} \stackrel{\text{def}}{=} \mathbf{E}_{(h,h') \sim \rho^2} [R_{D_S}(h, h') - R_{D_T}(h, h')],$$

by its empirical counterpart,

$$d_{S \times T}^{(1)} \stackrel{\text{def}}{=} \mathbf{E}_{(h,h') \sim \rho^2} [R_S(h, h') - R_T(h, h')],$$

and some extra terms related to the Kullback-Leibler divergence between the posterior and the prior. However, a notable difference with the proof of Theorem 3 is that the obtained bound will be simultaneously valid as an upper and a lower bound. Because of this, there will no need here to redo the all the proof to bound

$$d^{(2)} \stackrel{\text{def}}{=} \mathbf{E}_{(h,h') \sim \rho^2} [R_{D_T}(h, h') - R_{D_S}(h, h')],$$

and also, the present proof will not require the use of the union bound argument.

Again, we consider “abstract” classifiers $\hat{h} \in \mathcal{H}^2$ whose loss on a pair of examples $(\mathbf{x}^s, \mathbf{x}^t) \sim D_{S \times T}$ is defined by,

$$\mathcal{L}_{d^{(1)}}(\hat{h}, \mathbf{x}^s, \mathbf{x}^t) \stackrel{\text{def}}{=} \frac{1 + \mathcal{L}_{0 \rightarrow 1}(h(\mathbf{x}^s), h'(\mathbf{x}^s)) - \mathcal{L}_{0 \rightarrow 1}(h(\mathbf{x}^t), h'(\mathbf{x}^t))}{2}.$$

Note that, again, $\mathcal{L}_{d^{(1)}}$ lies in $[0, 1]$, and that $R_{S \times T}^{(1)}(\hat{h})$ and $R_{D_{S \times T}}^{(1)}(\hat{h})$ are as defined in the proof of Theorem 3.

Now, let us consider the non-negative random variable,

$$\mathbf{E}_{\hat{h} \sim \hat{\pi}} e^{m \text{kl}(R_{S \times T}^{(1)}(\hat{h}) \| R_{D_{S \times T}}^{(1)}(\hat{h}))}.$$

We apply Markov’s inequality (Lemma 1). For every $\delta \in (0, 1]$, with a probability at least $1 - \delta$ over the choice of $S \times T \sim (D_{S \times T})^m$, we have,

$$\begin{aligned} & \mathbf{E}_{\hat{h} \sim \hat{\pi}} e^{m \text{kl}(R_{S \times T}^{(1)}(\hat{h}) \| R_{D_{S \times T}}^{(1)}(\hat{h}))} \\ & \leq \frac{1}{\delta} \mathbf{E}_{S \times T \sim (D_{S \times T})^m} \mathbf{E}_{\hat{h} \sim \hat{\pi}} e^{m \text{kl}(R_{S \times T}^{(1)}(\hat{h}) \| R_{D_{S \times T}}^{(1)}(\hat{h}))}. \end{aligned}$$

By taking the logarithm on each side of the previous inequality, and transforming the expectation over $\hat{\pi}$ into an expectation over $\hat{\rho}$, we then obtain that,

$$\begin{aligned} & \ln \left[\mathbf{E}_{\hat{h} \sim \hat{\rho}} \frac{\hat{\pi}(\hat{h})}{\hat{\rho}(\hat{h})} e^{m \text{kl}(R_{S \times T}^{(1)}(\hat{h}) \| R_{D_{S \times T}}^{(1)}(\hat{h}))} \right] \quad (13) \\ & \leq \ln \left[\frac{1}{\delta} \mathbf{E}_{S \times T \sim (D_{S \times T})^m} \mathbf{E}_{\hat{h} \sim \hat{\rho}} e^{m \text{kl}(R_{S \times T}^{(1)}(\hat{h}) \| R_{D_{S \times T}}^{(1)}(\hat{h}))} \right] \\ & \leq \ln \frac{2\sqrt{m}}{\delta}. \end{aligned}$$

The last inequality comes from the Maurer’s lemma (Lemma 4).

Let us now re-write a part of the equation as $\text{KL}(\rho \| \pi)$ and let us then find a lower bound by using twice the Jensen’s inequality (Lemma 2), first on the concave logarithm function, and then on the convex function kl ,

$$\begin{aligned} & \ln \left[\mathbf{E}_{\hat{h} \sim \hat{\rho}} \frac{\hat{\pi}(\hat{h})}{\hat{\rho}(\hat{h})} e^{m \text{kl}(R_{S \times T}^{(1)}(\hat{h}) \| R_{D_{S \times T}}^{(1)}(\hat{h}))} \right] \\ & = \ln \left[\mathbf{E}_{\hat{h} \sim \hat{\rho}} e^{m \text{kl}(R_{S \times T}^{(1)}(\hat{h}) \| R_{D_{S \times T}}^{(1)}(\hat{h}))} \right] - 2\text{KL}(\rho \| \pi) \\ & \geq \mathbf{E}_{\hat{h} \sim \hat{\rho}} m \text{kl} \left(R_{S \times T}^{(1)}(\hat{h}) \| R_{D_{S \times T}}^{(1)}(\hat{h}) \right) - 2\text{KL}(\rho \| \pi) \\ & \geq m \text{kl} \left(\mathbf{E}_{\hat{h} \sim \hat{\rho}} R_{S \times T}^{(1)}(\hat{h}) \| \mathbf{E}_{\hat{h} \sim \hat{\rho}} R_{D_{S \times T}}^{(1)}(\hat{h}) \right) - 2\text{KL}(\rho \| \pi) \\ & \geq m \text{kl} \left(R_{S \times T}^{(1)}(G_{\hat{\rho}}) \| R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}}) \right) - 2\text{KL}(\rho \| \pi). \end{aligned}$$

This implies that,

$$\text{kl} \left(R_{S \times T}^{(1)}(G_{\hat{\rho}}) \| R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}}) \right) \leq \frac{1}{m} \left[2\text{KL}(\rho \| \pi) + \ln \frac{2\sqrt{m}}{\delta} \right].$$

Since, as in the proof of Theorem 3 for $d^{(1)}$, we have: $d^{(1)} = 2R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}}) - 1$ and $d_{S \times T}^{(1)} = 2R_{S \times T}^{(1)}(G_{\hat{\rho}}) - 1$, the previous line directly implies a bound on $d^{(1)}$ from its empirical counterpart $d_{S \times T}^{(1)}$. Hence, with probability at least $1 - \delta$ over the choice of $S \times T \sim (D_S \times D_T)^m$, we have,

$$\text{kl} \left(\frac{d_{S \times T}^{(1)} + 1}{2} \left\| \frac{d^{(1)} + 1}{2} \right. \right) \leq \frac{1}{m} \left[2\text{KL}(\rho \| \pi) + \ln \frac{2\sqrt{m}}{\delta} \right]. \quad (14)$$

We claim that we also have,

$$\text{kl} \left(\frac{|d_{S \times T}^{(1)}| + 1}{2} \left\| \frac{|d^{(1)}| + 1}{2} \right. \right) \leq \frac{1}{m} \left[2\text{KL}(\rho \| \pi) + \ln \frac{2\sqrt{m}}{\delta} \right], \quad (15)$$

which, since

$$|d^{(1)}| = \text{dis}_{\rho}(D_S, D_T) \quad \text{and} \quad |d_{S \times T}^{(1)}| = \text{dis}_{\rho}(S, T),$$

implies the result. Hence to finish the proof, let us prove the claim of Equation (15). There are four cases to consider.

Case 1: $d_{S \times T}^{(1)} \geq 0$ and $d^{(1)} \geq 0$. There is nothing to prove since in that case, Equations (14) and (15) coincide.

Case 2: $d_{S \times T}^{(1)} \leq 0$ and $d^{(1)} \leq 0$. This case reduces to Case 1 because of the following property of $\text{kl}(\cdot \| \cdot)$:

$$\text{kl} \left(\frac{a+1}{2} \left\| \frac{b+1}{2} \right. \right) = \text{kl} \left(\frac{-a+1}{2} \left\| \frac{-b+1}{2} \right. \right). \quad (16)$$

Case 3: $d_{S \times T}^{(1)} \leq 0$ and $d^{(1)} \geq 0$. From straightforward calculations, one can show that,

$$\begin{aligned}
 & \text{kl}\left(\frac{|d_{S \times T}^{(1)}|+1}{2} \parallel \frac{|d^{(1)}|+1}{2}\right) - \text{kl}\left(\frac{d_{S \times T}^{(1)}+1}{2} \parallel \frac{d^{(1)}+1}{2}\right) \\
 &= \text{kl}\left(\frac{-d_{S \times T}^{(1)}+1}{2} \parallel \frac{d^{(1)}+1}{2}\right) - \text{kl}\left(\frac{d_{S \times T}^{(1)}+1}{2} \parallel \frac{d^{(1)}+1}{2}\right) \\
 &= \left(\frac{-d_{S \times T}^{(1)}+1}{2} - \frac{d_{S \times T}^{(1)}+1}{2}\right) \ln\left(\frac{1}{\frac{d^{(1)}+1}{2}}\right) \\
 &\quad + \left(\left(1 - \frac{-d_{S \times T}^{(1)}+1}{2}\right) - \left(1 - \frac{d_{S \times T}^{(1)}+1}{2}\right)\right) \ln\left(\frac{1}{1 - \frac{d^{(1)}+1}{2}}\right) \\
 &= (-d_{S \times T}^{(1)}) \ln\left(\frac{1}{\frac{d^{(1)}+1}{2}}\right) + (d_{S \times T}^{(1)}) \ln\left(\frac{1}{1 - \frac{d^{(1)}+1}{2}}\right) \\
 &= (-d_{S \times T}^{(1)}) \ln\left(\frac{1}{\frac{d^{(1)}+1}{2}}\right) + (d_{S \times T}^{(1)}) \ln\left(\frac{1}{\frac{-d^{(1)}+1}{2}}\right) \\
 &= d_{S \times T}^{(1)} \ln\left(\frac{d^{(1)}+1}{-d^{(1)}+1}\right) \\
 &\leq 0. \tag{17}
 \end{aligned}$$

The last inequality follows from the fact that we have $d_{S \times T}^{(1)} \leq 0$ and $d^{(1)} \geq 0$.

Hence, from Equations (17) and (14), we have,

$$\begin{aligned}
 \text{kl}\left(\frac{|d_{S \times T}^{(1)}|+1}{2} \parallel \frac{|d^{(1)}|+1}{2}\right) &\leq \text{kl}\left(\frac{d_{S \times T}^{(1)}+1}{2} \parallel \frac{d^{(1)}+1}{2}\right) \\
 &\leq \frac{1}{m} \left[2\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right],
 \end{aligned}$$

as wanted.

Case 4: $d_{S \times T}^{(1)} \geq 0$ and $d^{(1)} \leq 0$. Again because of Equation (16), this case reduces to Case 3, and we are done. \square

From the preceding ‘‘Seeger’s type’’ results, one can then obtain the following PAC-Bayesian DA-bound.

Theorem 8. *For any domains P_S and P_T (respectively with marginals D_S and D_T) over $X \times Y$, any set of hypothesis \mathcal{H} , and any prior distribution π over \mathcal{H} , any $\delta \in (0, 1]$, with a probability at least $1 - \delta$ over the choice of $S \times T \sim (P_S \times P_T)^m$, we have,*

$$R_{P_T}(G_\rho) - R_{P_T}(G_{\rho_T^*}) \leq \sup \mathcal{R}_\rho + \sup \mathcal{D}_\rho + \lambda_\rho,$$

where $\lambda_\rho \stackrel{\text{def}}{=} R_{D_T}(G_\rho, G_{\rho_T^*}) + R_{D_S}(G_\rho, G_{\rho_T^*})$ and,

$$\begin{aligned}
 \mathcal{R}_\rho &\stackrel{\text{def}}{=} \left\{ r : \text{kl}(R_S(G_\rho) \parallel r) \leq \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right] \right\}, \\
 \mathcal{D}_\rho &\stackrel{\text{def}}{=} \left\{ d : \text{kl}\left(\frac{\text{dis}_\rho(S, T)+1}{2} \parallel \frac{d+1}{2}\right) \leq \frac{1}{m} \left[2\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right] \right\}.
 \end{aligned}$$

Proof. The result is obtained by inserting Ths. 6 and 7 (with $\delta := \frac{\delta}{2}$) in Th. 4 of the main paper. \square

3.2. PAC-Bayesian Bounds when $m \neq m'$

In the main paper, for the sake of simplicity, we restrict to the case where m (the size of the source set S) and m' (the size of the target set T) are equal. All the results generalize to the $m \neq m'$ case. In this subsection, we will show how it can be done from a ‘‘McAllester’s type’’ of bound (Similar results can be achieved for ‘‘Catoni’s type’’ or ‘‘Seeger’s type’’).

First we recall the PAC-Bayesian bound proposed by [McAllester \(2003\)](#), which is stated without a term allowing to control the trade-off between the complexity and the risk.

Theorem 9 ([McAllester \(2003\)](#)). *For any domain P_S over $X \times Y$, any set of hypothesis \mathcal{H} , and any prior distribution π over \mathcal{H} , any $\delta \in (0, 1]$, with a probability at least $1 - \delta$ over the choice of $S \sim (P_S)^m$, for every ρ over \mathcal{H} , we have,*

$$\left| R_{P_S}(G_\rho) - R_S(G_\rho) \right| \leq \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right]}.$$

Now we can prove the following consistency bound for $\text{dis}_\rho(D_S, D_T)$, when $m \neq m'$.

Theorem 10. *For any marginal distributions D_S and D_T over X , any set of hypothesis \mathcal{H} , any prior distribution π over \mathcal{H} , any $\delta \in (0, 1]$, with a probability at least $1 - \delta$ over the choice of $S \sim (D_S)^m$ and $T \sim (D_T)^{m'}$, for every ρ over \mathcal{H} , we have,*

$$\begin{aligned}
 \left| \text{dis}_\rho(D_S, D_T) - \text{dis}_\rho(S, T) \right| &\leq \sqrt{\frac{1}{2m} \left[2\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right]} \\
 &\quad + \sqrt{\frac{1}{2m'} \left[2\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m'}}{\delta} \right]}.
 \end{aligned}$$

Proof. Let us consider the non-negative random variable,

$$\mathbf{E}_{(h, h') \sim \pi^2} e^{2m(R_{D_S}(h, h') - R_S(h, h'))^2}.$$

We apply Markov’s inequality (Lemma 1). For every $\delta \in (0, 1]$, with a probability at least $1 - \delta$ over the choice of $S \sim (D_S)^m$, we have,

$$\begin{aligned}
 &\mathbf{E}_{(h, h') \sim \pi^2} e^{2m(R_{D_S}(h, h') - R_S(h, h'))^2} \\
 &\leq \frac{1}{\delta} \mathbf{E}_{S \sim (D_S)^m} \mathbf{E}_{(h, h') \sim \pi^2} e^{2m(R_{D_S}(h, h') - R_S(h, h'))^2}.
 \end{aligned}$$

By taking the logarithm on each side of the previous inequality and transforming the expectation over π^2 into an expectation over ρ^2 , we obtain that for every $\delta \in (0, 1]$, with a probability at least $1 - \delta$ over the choice of $S \sim (D_S)^m$, and for every posterior distribution ρ , we have,

$$\begin{aligned} & \ln \left[\mathbf{E}_{(h,h') \sim \rho^2} \frac{\pi(h)\pi(h')}{\rho(h)\rho(h')} e^{2m(R_{D_S}(h,h') - R_S(h,h'))^2} \right] \\ & \leq \ln \left[\frac{1}{\delta} \mathbf{E}_{S \sim (D_S)^m} \mathbf{E}_{(h,h') \sim \pi^2} e^{2m(R_{D_S}(h,h') - R_S(h,h'))^2} \right]. \end{aligned}$$

Since $\ln(\cdot)$ is a concave function, we can apply the Jensen's inequality (Lemma 2). Then, for every $\delta \in (0, 1]$, with a probability at least $1 - \delta$ over the choice of $S \sim (D_S)^m$, and for every posterior distribution ρ , we have,

$$\begin{aligned} & \mathbf{E}_{(h,h') \sim \rho^2} \ln \left[\frac{\pi(h)\pi(h')}{\rho(h)\rho(h')} e^{2m(R_{D_S}(h,h') - R_S(h,h'))^2} \right] \\ & \leq \ln \left[\frac{1}{\delta} \mathbf{E}_{S \sim (D_S)^m} \mathbf{E}_{(h,h') \sim \pi^2} e^{2m(R_{D_S}(h,h') - R_S(h,h'))^2} \right]. \end{aligned}$$

By the Equation (8),

$$\mathbf{E}_{(h,h') \sim \rho^2} \ln \left[\frac{\pi(h)\pi(h')}{\rho(h)\rho(h')} \right] = -2\text{KL}(\rho \parallel \pi).$$

For every $\delta \in (0, 1]$, with a probability at least $1 - \delta$ over the choice of $S \sim (D_S)^m$, and for every posterior distribution ρ , we have,

$$\begin{aligned} & -2\text{KL}(\rho \parallel \pi) + \mathbf{E}_{(h,h') \sim \rho^2} m 2(R_{D_S}(h,h') - R_S(h,h'))^2 \\ & \leq \ln \left[\frac{1}{\delta} \mathbf{E}_{S \sim (D_S)^m} \mathbf{E}_{(h,h') \sim \pi^2} e^{2m(R_{D_S}(h,h') - R_S(h,h'))^2} \right]. \end{aligned}$$

Since $2(a - b)^2$ is a convex function, we again apply Jensen inequality,

$$\begin{aligned} & \left(\mathbf{E}_{(h,h') \sim \rho^2} (R_{D_S}(h,h') - R_S(h,h')) \right)^2 \\ & \leq \mathbf{E}_{(h,h') \sim \rho^2} (R_{D_S}(h,h') - R_S(h,h'))^2. \end{aligned}$$

Thus, for every $\delta \in (0, 1]$, with a probability at least $1 - \delta$ over the choice of $S \sim (D_S)^m$, and for every posterior distribution ρ , we have,

$$\begin{aligned} 2m \left(\mathbf{E}_{(h,h') \sim \rho^2} R_{D_S}(h,h') - \mathbf{E}_{h,h' \sim \rho^2} R_S(h,h') \right)^2 & \leq 2\text{KL}(\rho \parallel \pi) \\ & + \ln \left[\frac{1}{\delta} \mathbf{E}_{S \sim (D_S)^m} \mathbf{E}_{(h,h') \sim \pi^2} e^{2m(R_{D_S}(h,h') - R_S(h,h'))^2} \right]. \end{aligned}$$

Let us now bound,

$$\ln \left[\frac{1}{\delta} \mathbf{E}_{S \sim (D_S)^m} \mathbf{E}_{(h,h') \sim \pi^2} e^{2m(R_{D_S}(h,h') - R_S(h,h'))^2} \right].$$

To do so, we have,

$$\begin{aligned} & \mathbf{E}_{S \sim (D_S)^m} \mathbf{E}_{(h,h') \sim \pi^2} e^{2m(R_{D_S}(h,h') - R_S(h,h'))^2} \\ & = \mathbf{E}_{(h,h') \sim \pi^2} \mathbf{E}_{S \sim (D_S)^m} e^{2m(R_{D_S}(h,h') - R_S(h,h'))^2} \quad (18) \end{aligned}$$

$$\leq \mathbf{E}_{(h,h') \sim \pi^2} \mathbf{E}_{S \sim (D_S)^m} e^{\text{kl}(R_S(h,h') \parallel R_{D_S}(h,h'))} \quad (19)$$

$$\leq 2\sqrt{m}. \quad (20)$$

Line (18) comes from the independence between D_S and π^2 . The Pinsker's inequality,

$$2(q - p)^2 \leq \text{kl}(q \parallel p) \quad \text{for any } p, q \in [0, 1],$$

gives Line (19). The last Line (20) comes from the Maurer's lemma (Lemma 4).

Thus for every $\delta \in (0, 1]$, with a probability at least $1 - \delta$ over the choice of $S \sim (D_S)^m$, and for every posterior distribution ρ , we obtain,

$$\begin{aligned} & 2m \left(\mathbf{E}_{(h,h') \sim \rho^2} R_{D_S}(h,h') - \mathbf{E}_{(h,h') \sim \rho^2} R_S(h,h') \right)^2 \\ & \leq 2\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \\ & \Leftrightarrow \left(\mathbf{E}_{(h,h') \sim \rho^2} R_{D_S}(h,h') - \mathbf{E}_{(h,h') \sim \rho^2} R_S(h,h') \right)^2 \\ & \leq \frac{1}{2m} \left[2\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \\ & \Leftrightarrow \left| \mathbf{E}_{(h,h') \sim \rho^2} R_{D_S}(h,h') - \mathbf{E}_{(h,h') \sim \rho^2} R_S(h,h') \right| \\ & \leq \sqrt{\frac{1}{2m} \left[2\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right]}. \quad (21) \end{aligned}$$

Following the same proof process for bounding $\left| \mathbf{E}_{(h,h') \sim \rho^2} R_{D_T}(h,h') - \mathbf{E}_{(h,h') \sim \rho^2} R_T(h,h') \right|$, we obtain the following result.

For every $\delta \in (0, 1]$, with a probability at least $1 - \delta$ over the choice of $T \sim (D_T)^{m'}$, and for every posterior distribution ρ ,

$$\begin{aligned} & \left| \mathbf{E}_{(h,h') \sim \rho^2} R_{D_T}(h,h') - \mathbf{E}_{(h,h') \sim \rho^2} R_T(h,h') \right| \\ & \leq \sqrt{\frac{1}{2m'} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m'}}{\delta} \right]}. \quad (22) \end{aligned}$$

Finally, let us substitute δ by $\frac{\delta}{2}$ in Inequalities (21) and (22). This, together with the union bound that assure that both results hold simultaneously, gives the

result because,

$$\begin{aligned} \left| \mathbf{E}_{(h,h') \sim \rho^2} [R_{D_T}(h, h') - R_{D_S}(h, h')] \right| &= \text{dis}_\rho(D_S, D_T), \\ \left| \mathbf{E}_{(h,h') \sim \rho^2} [R_T(h, h') - R_S(h, h')] \right| &= \text{dis}_\rho(S, T), \end{aligned}$$

and because if $|a_1 - b_1| \leq c_1$ and $|a_2 - b_2| \leq c_2$, then $|(a_1 - a_2) - (b_1 - b_2)| \leq c'_1 + c'_2$. \square

Then we can obtain the following PAC-Bayesian DA-bound.

Theorem 11. *For any domains P_S and P_T (respectively with marginals D_S and D_T) over $X \times Y$, and for any set \mathcal{H} of hypothesis, for any prior distribution π over \mathcal{H} , any $\delta \in (0, 1]$, with a probability at least $1 - \delta$ over the choice of $S_1 \sim (D_S)^m$, $S_2 \sim (D_S)^{m'}$, and $T \sim (D_T)^{m'}$, for every ρ over \mathcal{H} , we have,*

$$\begin{aligned} R_{P_T}(G_\rho) - R_{P_T}(G_{\rho_T^*}) &\leq R_S(G_\rho) + \text{dis}_\rho(S, T) + \lambda_\rho \\ &\quad + \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right]} \\ &\quad + \sqrt{\frac{1}{2m} \left[2\text{KL}(\rho \parallel \pi) + \ln \frac{8\sqrt{m}}{\delta} \right]} \\ &\quad + \sqrt{\frac{1}{2m'} \left[2\text{KL}(\rho \parallel \pi) + \ln \frac{8\sqrt{m'}}{\delta} \right]}. \end{aligned}$$

where $\lambda_\rho \stackrel{\text{def}}{=} R_{D_T}(G_\rho, G_{\rho_T^*}) + R_{D_S}(G_\rho, G_{\rho_T^*})$.

Proof. The result is obtained by inserting Ths. 9 and 10 (with $\delta := \frac{\delta}{2}$) in Th. 4 of the main paper. \square

4. PBDA Algorithm Details

4.1. Objective function and gradient

Given a source sample $S = \{(\mathbf{x}_i^s, y_i^s)\}_{i=1}^m$, a target sample $T = \{(\mathbf{x}_i^t)\}_{i=1}^m$, and fixed parameters $A > 0$ and $C > 0$, the learning algorithm PBDA consists in finding the weight vector \mathbf{w} minimizing,

$$\begin{aligned} \frac{\|\mathbf{w}\|^2}{2} + C \sum_{i=1}^m \Phi_{\text{cvx}} \left(y_i^s \frac{\mathbf{w} \cdot \mathbf{x}_i^s}{\|\mathbf{x}_i^s\|} \right) \\ + A \left| \sum_{i=1}^m \Phi_{\text{dis}} \left(\frac{\mathbf{w} \cdot \mathbf{x}_i^s}{\|\mathbf{x}_i^s\|} \right) - \Phi_{\text{dis}} \left(\frac{\mathbf{w} \cdot \mathbf{x}_i^t}{\|\mathbf{x}_i^t\|} \right) \right|, \quad (23) \end{aligned}$$

where, **Erf** being the Gauss error function,

$$\begin{aligned} \Phi(a) &\stackrel{\text{def}}{=} \frac{1}{2} \left[1 - \text{Erf} \left(\frac{a}{\sqrt{2}} \right) \right], \\ \Phi_{\text{cvx}}(a) &\stackrel{\text{def}}{=} \max \left[\Phi(a), \frac{1}{2} - \frac{a}{\sqrt{2\pi}} \right], \\ \Phi_{\text{dis}}(a) &\stackrel{\text{def}}{=} 2 \times \Phi(a) \times \Phi(-a). \end{aligned}$$

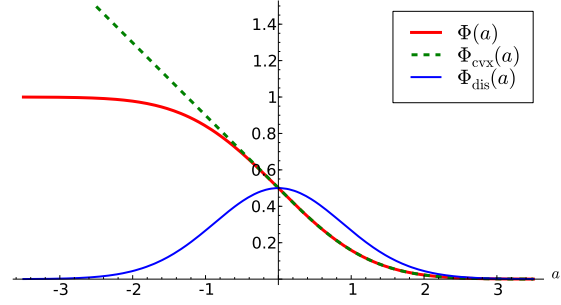


Figure 1. Behaviour of functions $\Phi(\cdot)$, $\Phi_{\text{cvx}}(\cdot)$ and $\Phi_{\text{dis}}(\cdot)$.

Figure 1 illustrates these three functions.

The gradient of the Equation (23) is given by,

$$\begin{aligned} \mathbf{w} + C \sum_{i=1}^m \Phi'_{\text{cvx}} \left(\frac{y_i^s \mathbf{w} \cdot \mathbf{x}_i^s}{\|\mathbf{x}_i^s\|} \right) \frac{y_i^s \mathbf{x}_i^s}{\|\mathbf{x}_i^s\|} \\ + s \times A \left[\sum_{i=1}^m \Phi'_{\text{dis}} \left(\frac{\mathbf{w} \cdot \mathbf{x}_i^s}{\|\mathbf{x}_i^s\|} \right) \frac{\mathbf{x}_i^s}{\|\mathbf{x}_i^s\|} - \Phi'_{\text{dis}} \left(\frac{\mathbf{w} \cdot \mathbf{x}_i^t}{\|\mathbf{x}_i^t\|} \right) \frac{\mathbf{x}_i^t}{\|\mathbf{x}_i^t\|} \right], \end{aligned}$$

where $\Phi'_{\text{cvx}}(a)$ and $\Phi'_{\text{dis}}(a)$ are respectively the derivatives of functions Φ_{cvx} and Φ_{dis} evaluated at point a , and $s = \text{sgn} \left[\sum_{i=1}^m \Phi_{\text{dis}} \left(\frac{\mathbf{w} \cdot \mathbf{x}_i^s}{\|\mathbf{x}_i^s\|} \right) - \Phi_{\text{dis}} \left(\frac{\mathbf{w} \cdot \mathbf{x}_i^t}{\|\mathbf{x}_i^t\|} \right) \right]$.

4.2. Using a kernel function

The kernel trick allows us to work with dual weight vector $\boldsymbol{\alpha} \in \mathbb{R}^{2m}$ that is a linear classifier in an augmented space. Given a kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, we have,

$$h_{\mathbf{w}}(\mathbf{x}) = \sum_{i=1}^m \alpha_i k(\mathbf{x}_i^s, \mathbf{x}) + \sum_{i=1}^m \alpha_{i+m} k(\mathbf{x}_i^t, \mathbf{x}).$$

Let us denote K the kernel matrix of size $2m \times 2m$ such as,

$$K_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$$

where,

$$\mathbf{x}_\# = \begin{cases} \mathbf{x}_\#^s & \text{if } \# \leq m \\ \mathbf{x}_{\#-m}^t & \text{otherwise.} \end{cases}$$

In that case, the objective function of Equation (23) is rewritten in term of the vector $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_{2m})$ as,

$$\begin{aligned} \frac{1}{2} \sum_{i=1}^{2m} \sum_{j=1}^{2m} \alpha_i \alpha_j K_{i,j} + C \sum_{i=1}^m \Phi_{\text{cvx}} \left(y_i^s \frac{\sum_{j=1}^{2m} \alpha_j K_{i,j}}{\sqrt{K_{i,i}}} \right) \\ + A \left| \sum_{i=1}^m \Phi_{\text{dis}} \left(\frac{\sum_{j=1}^{2m} \alpha_j K_{i,j}}{\sqrt{K_{i,i}}} \right) - \Phi_{\text{dis}} \left(\frac{\sum_{j=1}^{2m} \alpha_j K_{i+m,j}}{\sqrt{K_{i+m,i+m}}} \right) \right|. \end{aligned}$$

The gradient of the latter equation is given by the vector $\alpha' = (\alpha'_1, \alpha'_2, \dots, \alpha'_{2m})$, with $\alpha'_{\#}$ equals to,

$$\begin{aligned} & \sum_{j=1}^{2m} \alpha_j K_{i,\#} + C \sum_{i=1}^m \Phi_{\text{cvx}} \left(y_i^s \frac{\sum_{j=1}^{2m} \alpha_j K_{i,j}}{\sqrt{K_{i,i}}} \right) \frac{y_i^s K_{i,\#}}{\sqrt{K_{i,i}}} \\ & + s \times A \left[\sum_{i=1}^m \Phi_{\text{dis}} \left(\frac{\sum_{j=1}^{2m} \alpha_j K_{i,j}}{\sqrt{K_{i,i}}} \right) \frac{K_{i,\#}}{\sqrt{K_{i,i}}} \right. \\ & \quad \left. - \Phi_{\text{dis}} \left(\frac{\sum_{j=1}^{2m} \alpha_j K_{i+m,j}}{\sqrt{K_{i+m,i+m}}} \right) \frac{K_{i+m,\#}}{\sqrt{K_{i+m,i+m}}} \right], \end{aligned}$$

where,

$$s = \text{sgn} \left[\sum_{i=1}^m \Phi_{\text{dis}} \left(\frac{\sum_{j=1}^{2m} \alpha_j K_{i,j}}{\sqrt{K_{i,i}}} \right) - \Phi_{\text{dis}} \left(\frac{\sum_{j=1}^{2m} \alpha_j K_{i+m,j}}{\sqrt{K_{i+m,i+m}}} \right) \right].$$

4.3. Implementation details

For our experiments, we minimize the objective function using a *Broyden-Fletcher-Goldfarb-Shanno method (BFGS)* implemented in the *scipy* python library¹. We made our code available at the following URL:

<http://graal.ift.ulaval.ca/pbda/>

When selecting hyperparameters by reverse cross-validation, we search on a 20×20 parameter grid for a A between 0.01 and 10^6 and a parameter C between 1.0 and 10^8 , both on a logarithm scale.

References

- Maurer, A. A note on the PAC Bayesian theorem. *CoRR*, cs.LG/0411099, 2004.
- McAllester, D. PAC-Bayesian stochastic model selection. *Machine Learning*, 51:5–21, 2003.
- Seeger, M. PAC-Bayesian generalization bounds for gaussian processes. *Journal of Machine Learning Research*, 3:233–269, 2002.

¹Available at <http://www.scipy.org/>