

Annexe de Une analyse PAC-Bayésienne de l'adaptation de domaine

Pascal Germain¹, Amaury Habrard², François Laviolette¹, et Emilie Morvant³

¹Département d'informatique et de génie logiciel, Université Laval, Québec, Canada

²Université Jean Monnet de Saint-Étienne, Laboratoire Hubert Curien, UMR CNRS 5516

³Aix-Marseille Univ., LIF-QARMA, UMR CNRS 7279

27 mai 2013

Preuve détaillée du Théorème 3

Nous rappelons le th. 3

Théorème 3. *Pour toutes marginales D_S et D_T sur X , pour tout espace d'hypothèses \mathcal{H} , pour toute distribution prior π sur \mathcal{H} , pour tout $\delta \in (0, 1]$, avec une probabilité d'au moins $1 - \delta$ sur le choix de $S \times T \sim (D_S \times D_T)^{m'}$, pour toute distribution ρ sur \mathcal{H} , on a,*

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

où $\text{dis}_\rho(S, T)$ est le désaccord empirique.

Démonstration. Tout d'abord, nous bornons :

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

Pour cela, nous considérons un classifieur "abstrait" $\hat{h} \stackrel{\text{def}}{=} (h, h') \in \mathcal{H}^2$ choisi selon une distribution $\hat{\rho}$, telle que $\hat{\rho}(\hat{h}) = \rho(h)\rho(h')$. Notons qu'avec $\hat{\pi}(\hat{h}) = \pi(h)\pi(h')$, on obtient $\text{KL}(\hat{\rho} \parallel \hat{\pi}) = 2\text{KL}(\rho \parallel \pi)$. En effet :

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

Nous définissons la fonction perte "abstraite" de \hat{h} sur une paire d'exemples $(\mathbf{x}^s, \mathbf{x}^t) \sim (D_S \times D_T)$ par :

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

Le risque "abstrait" de \hat{h} sur la distribution jointe est défini par : $R_{D_{S \times T}}^{(1)}(\hat{h}) = \mathbf{E}_{\mathbf{x}^s \sim D_S} \mathbf{E}_{\mathbf{x}^t \sim D_T} \ell_{d^{(1)}}(\hat{h}, \mathbf{x}^s, \mathbf{x}^t)$.

Ainsi, le risque du classifieur de Gibbs associé à cette perte est :

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

Les versions empiriques de ces deux quantités sont :

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

et

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

Puis que $\ell_{d^{(1)}}$ retourne des valeurs entre $[0, 1]$, on peut borner $R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}})$ en suivant le principe de la preuve du th. 2 (avec $c = 2\alpha$). Pour ce faire, nous définissons la fonction convexe

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

et nous considérons la variable aléatoire non négative suivante :

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

Nous appliquons ensuite l'inégalité de Markov (voir le th. 9). Pour tout $\delta \in (0, 1]$, avec une probabilité d'au moins $1 - \delta$ sur le choix de $S \times T \sim (D_{S \times T})^m$, on a :

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

En prenant le logarithme de chaque côté de l'inégalité, pour tout $\delta \in (0, 1]$, avec une probabilité d'au moins $1 - \delta$ sur le choix de $S \times T \sim (D_{S \times T})^m$, et pour toute distribution posterior $\hat{\rho}$, on a :

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

Le choix de \mathcal{F} implique :

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

On a :

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

On insère $\text{KL}(\rho \parallel \pi)$ dans la partie gauche de l'inégalité précédente et on trouve une borne inférieure en appliquant deux fois l'inégalité de Jensen (voir th. 10) : tout d'abord sur la fonction logarithme (concave), puis sur la fonction convexe \mathcal{F} . On obtient :

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

On a :

$$\begin{aligned} \mathcal{F}(R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}})) & \leq 2\alpha R_{S \times T}^{(1)}(G_{\hat{\rho}}) + \frac{2\text{KL}(\rho \parallel \pi) + \ln \frac{1}{\delta}}{m} \\ \Leftrightarrow R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}}) & \leq \frac{1}{1 - e^{-2\alpha}} \left[1 - e^{-\left(2\alpha R_{S \times T}^{(1)}(G_{\hat{\rho}}) + \frac{2\text{KL}(\rho \parallel \pi) + \ln \frac{1}{\delta}}{m} \right)} \right], \end{aligned}$$

puis l'inégalité $1 - e^{-x} \leq x$, donne :

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

Ainsi, puisque $d^{(1)} = 2R_{D_{S \times T}}^{(1)}(G_{\hat{\rho}}) - 1$, on obtient une borne sur $d^{(1)}$ à partir de son estimation empirique (notée $d_{S \times T}^{(1)}$). Nous obtenons donc avec une probabilité d'au moins $1 - \frac{\delta}{2}$ sur le choix de $S \times T \sim (D_S \times D_T)^m$,

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

Ensuite, nous bornons $d^{(2)} \stackrel{\text{def}}{=} \mathbf{E}_{h, h' \sim \rho^2} [R_{D_T}(h, h') - R_{D_S}(h, h')]$ en utilisant la même méthode :

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

Notons que $|d^{(1)}| = |d^{(2)}| = \text{dis}_{\rho}(D_S, D_T)$ et $|d_{S \times T}^{(1)}| = |d_{S \times T}^{(2)}| = \text{dis}_{\rho}(S, T)$. Ainsi, la borne sur $\text{dis}_{\rho}(D_S, D_T)$ est obtenue en prenant le maximum de la borne sur $d^{(1)}$ et de la borne $d^{(2)}$. Finalement en appliquant la borne de l'union, avec une probabilité d'au moins $1 - \delta$ sur le choix de $S \times T \sim (D_S \times D_T)^m$, on a

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

□

Bornes PAC-Bayésienne si $m \neq m'$

Nous rappelons tout d'abord la borne PAC-Bayésienne proposée par [McA03], énoncée sans terme de contrôle du compromis complexité/risque.

Théorème 6 ([McA03]). *Pour tout P_S sur $X \times Y$, pour tout ensemble \mathcal{H} , et pour toute distribution prior π sur \mathcal{H} , pour tout $\delta \in (0, 1]$, avec une probabilité d'au moins $1 - \delta$ sur le choix de $S \sim (P_S)^m$, pour toute distribution posterior ρ sure \mathcal{H} , on a*

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

Nous pouvons maintenant prouver la borne de consistance suivante pour $\text{dis}_\rho(D_S, D_T)$, pour le cas $m \neq m'$.

Théorème 7. *Pour toute distribution marginale D_S et D_T sur X , pour tout ensemble \mathcal{H} , pour toute distribution prior π sur \mathcal{H} , pour tout $\delta \in (0, 1]$, avec une probabilité d'au moins $1 - \delta$ sur le choix de $S \sim (D_S)^m$ et $T \sim (D_T)^{m'}$, pour toute distribution posterior ρ sur \mathcal{H} , on a :*

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

Démonstration. Considérons la variable aléatoire non négative :

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

On applique l'inégalité de Markov (voir th. 9), pour tout $\delta \in (0, 1]$, avec une probabilité d'au moins $1 - \delta$ sur le choix de $S \sim (D_S)^m$, on a :

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

En prenant le logarithme de chaque côté de cette inégalité, pour tout $\delta \in (0, 1]$, avec une probabilité d'au moins $1 - \delta$ sur le choix $S \sim (D_S)^m$ et pour toute distribution ρ , on a :

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

Puis que $\ln(\cdot)$ est une fonction concave, on applique l'inégalité de Jensen (voir th. 10). Ainsi, pour tout $\delta \in (0, 1]$, avec une probabilité d'au moins $1 - \delta$ sur le choix $S \sim (D_S)^m$ et pour toute distribution ρ , on a :

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

D'après les propriétés de la KL-divergence (éq. (8)), on a :

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

Ainsi :

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

Pour tout $\delta \in (0, 1]$, avec une probabilité d'au moins $1 - \delta$ sur le choix $S \sim (D_S)^m$ et pour toute distribution ρ , on a :

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

Puisque $2(a - b)^2$ est une fonction convexe, on applique encore une fois l'inégalité de Jensen :

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

Donc pour tout $\delta \in (0, 1]$, avec une probabilité d'au moins $1 - \delta$ sur le choix $S \sim (D_S)^m$ et pour toute distribution ρ , on a :

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

Il ne reste plus qu'à borner :

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

On a :

$$\begin{aligned} & \mathbf{E}_{S \sim (D_S)^m} \mathbf{E}_{h, h' \sim \pi^2} e^{\frac{m}{2} (R_{D_S}(h, h') - R_S(h, h'))^2} \\ & = \mathbf{E}_{h, h' \sim \pi^2} \mathbf{E}_{S \sim (D_S)^m} e^{\frac{m}{2} (R_{D_S}(h, h') - R_S(h, h'))^2} \quad (9) \end{aligned}$$

$$\leq \mathbf{E}_{h, h' \sim \pi^2} \mathbf{E}_{S \sim (D_S)^m} e^{\frac{m}{4} \text{kl}(R_{D_S}(h, h') \parallel R_S(h, h'))} \quad (10)$$

$$\leq \sqrt{m}. \quad (11)$$

La ligne (9) provient de l'indépendance entre D_S et ρ . L'inégalité $2(q - p)^2 \leq \text{kl}(q \parallel p)$ pour tout $p, q \in [0, 1]$ implique la ligne (10). La dernière ligne (11) est due au lemme de Maurer (voir lem. 1).

Ainsi, pour tout $\delta \in (0, 1]$, avec une probabilité d'au moins $1 - \delta$ sur le choix $S \sim (D_S)^m$ et pour toute

distribution ρ , on a :

$$\begin{aligned}
& \frac{m}{2} \left(\mathbf{E}_{h,h' \sim \rho^2} R_{D_S}(h, h') - \mathbf{E}_{h,h' \sim \rho^2} R_S(h, h') \right)^2 \\
& \leq 2\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \\
& \Leftrightarrow \left(\mathbf{E}_{h,h' \sim \rho^2} R_{D_S}(h, h') - \mathbf{E}_{h,h' \sim \rho^2} R_S(h, h') \right)^2 \\
& \leq \frac{2}{m} \left[2\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \\
& \Leftrightarrow \left| \mathbf{E}_{h,h' \sim \rho^2} R_{D_S}(h, h') - \mathbf{E}_{h,h' \sim \rho^2} R_S(h, h') \right| \\
& \leq \sqrt{\frac{2}{m} \left[2\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right]}. \tag{12}
\end{aligned}$$

En suivant le même principe de preuve, on borne

$$\left| \mathbf{E}_{h,h' \sim \rho^2} R_{D_T}(h, h') - \mathbf{E}_{h,h' \sim \rho^2} R_T(h, h') \right|.$$

Pour tout $\delta \in (0, 1]$, avec une probabilité d'au moins $1 - \delta$ sur le choix $S \sim (D_S)^m$ et pour toute distribution ρ , on a :

$$\begin{aligned}
& \left| \mathbf{E}_{h,h' \sim \rho^2} R_{D_T}(h, h') - \mathbf{E}_{h,h' \sim \rho^2} R_T(h, h') \right| \\
& \leq \sqrt{\frac{2}{m'} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m'}}{\delta} \right]}. \tag{13}
\end{aligned}$$

On remplace δ par $\frac{\delta}{2}$ dans les Inégalités (12) et (13). Puis en appliquant la borne de l'union, les deux résultats sont vrais simultanément avec une probabilité d'au moins $1 - \delta$, et nous les réunissons¹ pour borner :

$$\left| \mathbf{E}_{h,h' \sim \rho^2} [R_{D_T}(h, h') - R_{D_S}(h, h')] \right| = \text{dis}_\rho(D_S, D_T). \quad \square$$

On obtient alors la borne suivante.

Théorème 8. *Pour tout domaine P_S et P_T (de marginale respective D_S et D_T) sur $X \times Y$, et pour tout ensemble \mathcal{H} , pour toute distribution π sur \mathcal{H} , pour tout $\delta \in (0, 1]$, avec une probabilité d'au moins $1 - \delta$ sur le choix de $S_1 \sim (D_S)^{m_1}$, $S_2 \sim (D_S)^{m_2}$, et $T \sim (D_T)^{m'}$,*

1. Si $-c'_1 \leq a_1 - b_1 \leq c_1$ et $-c'_2 \leq a_2 - b_2 \leq c_2$, alors $-(c'_1 + c'_2) \leq (a_1 - a_2) - (b_1 - b_2) \leq c'_1 + c'_2$.

pour toute distribution posterior ρ sure \mathcal{H} , on a :

$$\begin{aligned}
R_{P_T}(G_\rho) - R_{P_T}(G_{\rho_T^*}) & \leq R_S(G_\rho) + \text{dis}_\rho(S, T) + \lambda_\rho \\
& + \sqrt{\frac{2}{m_1} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{8\sqrt{m_1}}{\delta} \right]} \\
& + \sqrt{\frac{2}{m_2} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{8\sqrt{m_2}}{\delta} \right]} \\
& + \sqrt{\frac{2}{m'} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m'}}{\delta} \right]}.
\end{aligned}$$

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

Démonstration. Le résultat est directement obtenu en insérant les ths. 6, 7 (avec $\delta := \frac{\delta}{2}$) dans Th. 4. \square

Quelques outils

Théorème 9 (Inégalité de Markov). *Soit Z une variable aléatoire et $\epsilon \geq 0$, alors :*

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

Théorème 10 (Inégalité de Jensen). *Soit X une variable aléatoire réelle intégrable et $g(\cdot)$ convexe, alors :*

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

Lemme 1 (des inégalités (1) et (2) de [Mau04]). *Soit $m \geq 8$, soit $X = (X_1, \dots, X_m)$ un vecteur i.i.d. de variables aléatoires, $0 \leq X_i \leq 1$. Alors :*

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

où $\text{kl}(a \parallel b) = a \ln \frac{a}{b} + (1 - a) \ln \frac{1-a}{1-b}$.

Références

- [Mau04] A. Maurer. A note on the pac bayesian theorem. *CoRR*, cs.LG/0411099, 2004.
- [McA03] D. McAllester. PAC-Bayesian stochastic model selection. *Mach. Learn.*, 51 :5–21, 2003.