

Inégalités Oracles pour le hold-out

Joseph Rynkiewicz

11 janvier 2021

1 Introduction

La technique du hold-out (sélection de modèle sur un ensemble de validation) est largement répandue. Cela revient à choisir un meilleur modèle parmi un nombre fini de modèle.

Nous allons voir, que pour la classification, le hold-out vérifie des inégalités oracles qui justifient son importance dans la pratique.

Notons, que bien que nous ne parlerons que de classification, on pourrait généraliser de nombreuses inégalités à d'autres fonctions de perte bornées.

2 Cadre de la classification

Dans toute la suite on considère une suite de variables aléatoires indépendantes et identiquement distribuées (i.i.d.)

$$\left(\left(\begin{array}{c} X_1 \\ Y_1 \end{array} \right), \dots, \left(\begin{array}{c} X_n \\ Y_n \end{array} \right) \right)$$

où $Y_i \in \{1, \dots, K\}$, est une variable qualitative à expliquer (par exemple la classe d'une image) et $X_i \in \mathbb{R}^d$ la variables explicative (par exemple les pixels de l'image).

Soit \mathcal{F} , l'ensemble des fonctions de $\mathbb{R}^d \rightarrow \{1, \dots, K\}$ et f une fonction de classification, on définit la perte théorique $L(f) = \mathbb{E}(\mathbf{1}_{\{Y_i \neq f(X_i)\}})$, la probabilité de mauvaise classification de la fonction f . On notera f^* la meilleure fonction de classification :

$$f^* = \arg \min_{f \in \mathcal{F}} L(f).$$

f^* est l'oracle que fait des prédictions parfaites à partir de $\left(\left(\begin{array}{c} X_1 \\ Y_1 \end{array} \right), \dots, \left(\begin{array}{c} X_n \\ Y_n \end{array} \right) \right)$.

Cette fonction n'est pas facilement atteignable, pour l'approximer on se donne un ensemble de fonctions \mathcal{G} (par exemple un ensemble de réseaux profonds) et on cherche à estimer

$$\tilde{f} = \arg \min_{f \in \mathcal{G}} L(f).$$

Si l'ensemble \mathcal{G} est bien choisi, on peut espérer que la perte $L(\hat{f})$ soit proche de la perte de l'oracle $L(f^*)$.

Evidemment, $L(f)$ n'est pas calculable directement mais elle peut être approximée par le risque empirique

$$\hat{L}_n(f) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{Y_i \neq f(X_i)\}}.$$

Si l'ensemble de fonctions \mathcal{G} est trop grand, il se peut que le risque empirique et le risque théorique soient très différents (c'est le phénomène de sur-apprentissage). Le hold-out résout ce problème en coupant le problème en deux. On cherche d'abord de bonnes fonctions de classification grâce à l'ensemble d'apprentissage, on obtient ainsi un ensemble fini de fonctions candidates. Ensuite, sur un ensemble indépendant (l'ensemble de validation) on choisit la meilleure fonction de classification parmi cet ensemble fini de fonctions.

Plus formellement, si on dispose d'une collection finie de fonctions $\{f_1, \dots, f_N\}$ on choisira comme fonction :

$$f_{\hat{k}} = \arg \min_{k \in \{1, \dots, N\}} \hat{L}_n(f_k).$$

Si on savait calculer la perte théorique, on aurait plutôt choisi

$$f_{\tilde{k}} = \arg \min_{k \in \{1, \dots, N\}} L(f_k).$$

On va donc étudier la différence qu'il peut y avoir entre $L(f_{\hat{k}})$ et $L(f_{\tilde{k}})$. Mais pour cela on a besoin de quelques résultats simples et utilisés constamment dans toute la suite.

2.1 Quelques astuces de calcul

On va donner des astuces dans le cas particulier des variables bornée par 1, mais des généralisations sont souvent simples. On rappelle que, si Y est une variable aléatoire positive, bornée par 1, on aura, par Fubini :

$$E(Y) = \int_0^1 P(Y > t) dt.$$

— Soit $\lambda > 0$ et $C \geq 1$, si on a une variable aléatoire Z , avec $|Z| \leq 1$, qui, pour tout $t > 0$, vérifie l'inégalité

$$P(Z > t) \leq C \exp(-\lambda t) = \exp(\ln C - \lambda t)$$

ainsi

$$P\left(Z - \frac{\ln C}{\lambda} > \varepsilon\right) = P\left(Z > \varepsilon + \frac{\ln C}{\lambda}\right) \leq \exp(\ln C - \lambda\varepsilon - \ln(C)) = \exp(-\lambda\varepsilon).$$

On aura donc,

$$E\left(Z - \frac{\ln C}{\lambda}\right) \leq E\left(\max\left(Z - \frac{\ln C}{\lambda}, 0\right)\right) = \int_0^1 P\left(\max\left(Z - \frac{\ln C}{\lambda}, 0\right) > t\right) dt \leq \int_0^1 \exp(-\lambda t) dt = \left[-\frac{1}{\lambda} \exp(-\lambda t)\right]_0^1 = \frac{1}{\lambda} (1 - \exp(-\lambda)) \leq \frac{1}{\lambda}.$$

On en déduit que

$$E(Z) \leq \frac{1}{\lambda} (\ln C + 1) \quad (1)$$

— Soit $\lambda > 0$ et $C \geq 1$, si on a une variable aléatoire $Z > 0$, avec $Z \leq 1$, qui, pour tout $t > 0$, vérifie l'inégalité

$$P(Z > t) \leq C \exp(-\lambda t^2) = \exp(\ln C - \lambda t^2)$$

ainsi

$$\begin{aligned} P(Z^2 > t) &\leq \exp(\ln C - \lambda t) \\ P\left(Z^2 - \frac{\ln C}{\lambda} > \varepsilon\right) &= P\left(Z^2 > \varepsilon + \frac{\ln C}{\lambda}\right) \leq \\ &\exp(\ln C - \lambda \varepsilon - \ln C) = \exp(-\lambda \varepsilon). \end{aligned}$$

On aura donc,

$$E\left(Z^2 - \frac{\ln C}{\lambda}\right) \leq E\left(\max\left(Z^2 - \frac{\ln C}{\lambda}, 0\right)\right) = \int_0^1 P\left(\max\left(Z^2 - \frac{\ln C}{\lambda}, 0\right) > t\right) dt \leq \int_0^1 \exp(-\lambda t) dt = \left[-\frac{1}{\lambda} \exp(-\lambda t)\right]_0^1 = \frac{1}{\lambda} (1 - \exp(-\lambda)) \leq \frac{1}{\lambda}.$$

On en déduit que

$$E(Z^2) \leq \frac{1}{\lambda} (\ln C + 1)$$

et comme, par l'inégalité de Jensen, $(E(Z))^2 \leq E(Z^2)$, on aura :

$$E(Z) \leq \sqrt{\frac{1}{\lambda} (\ln C + 1)} \quad (2)$$

$$E(Z) \leq \sqrt{\frac{1}{\lambda} (\ln C + 1)}.$$

— Si on a une collection Z_1, \dots, Z_N de variables aléatoires i.i.d., on a la borne de l'union (une borne plutôt grossière, mais utilisée constamment) :

$$P\left(\max_{k \in \{1, \dots, N\}} Z_k > \varepsilon\right) \leq P\left(\cup_{k \in \{1, \dots, N\}} \{Z_k > \varepsilon\}\right) \leq NP(Z_k > \varepsilon) \quad (3)$$

3 Un premier Résultat en $O\left(\frac{1}{\sqrt{n}}\right)$

3.1 Inégalités exponentielles

Commençons par un résultat de Hoeffding (1963), on pourra en trouver une démonstration dans [Lugosi, 2002] :

Lemme 1 Soit $S(f) = \sum_{i=1}^n \mathbf{1}_{\{Y_i \neq f(X_i)\}}$, avec $\mathbb{E}(S(f)) = nL(f) \geq 0$, on a $|S(f) - nL(f)| \leq n$ et l'inégalité de Hoeffding permet d'écrire :

$$P(S(f) > nL(f) + nt) \leq \exp(-2nt^2), \quad (4)$$

soit

$$P\left(\frac{1}{n}S(f) - L(f) > t\right) \leq \exp(-2nt^2), \quad (5)$$

et

$$P\left(L(f) - \frac{1}{n}S(f) > t\right) \leq \exp(-2nt^2), \quad (6)$$

Ce résultat permettrait d'avoir la borne sur l'espérance de $\frac{1}{n}S(f)$, en utilisant les résultats de la section 2.1, on aura

$$E\left(\sup_{k \in \{1, \dots, N\}} \frac{1}{n}S(f_k) - L(f_k)\right) \leq \sqrt{\frac{\ln N + 1}{2n}}.$$

On peut faire un peu mieux grâce à un résultat que l'on trouvera par exemple dans [Massart, 2003], p. 196. :

Lemme 2 Soit $(f_k)_{1 \leq k \leq N}$, une collection finie de fonctions. Notons $Z(f_k) = \sum_{i=1}^n (\mathbf{1}_{\{Y_i \neq f_k(X_i)\}} - E(\mathbf{1}_{\{Y_i \neq f_k(X_i)\}}))$, on a $E(Z(f_k)) = 0$, et $-nE(\mathbf{1}_{\{Y_i \neq f_k(X_i)\}}) \leq Z(f_k) \leq n - nE(\mathbf{1}_{\{Y_i \neq f_k(X_i)\}})$, alors

$$E\left(\sup_{k \in \{1, \dots, N\}} \frac{1}{n}Z(f_k)\right) \leq \sqrt{\frac{\ln N}{2n}}.$$

3.2 Première inégalité oracle

Soit (f_1, \dots, f_N) une collection finie de fonctions obtenues sur l'ensemble d'apprentissage (indépendant de l'ensemble de validation). On note

$$\tilde{k} = \arg \min_{k \in \{1, \dots, N\}} L(f_k) \text{ et } \hat{k} = \arg \min_{k \in \{1, \dots, N\}} \hat{L}_n(f_k).$$

On a

$$L(f_{\hat{k}}) - L(f_{\tilde{k}}) = L(f_{\hat{k}}) - \hat{L}_n(f_{\hat{k}}) + \hat{L}_n(f_{\hat{k}}) - \hat{L}_n(f_{\tilde{k}}) + \hat{L}_n(f_{\tilde{k}}) - L(f_{\tilde{k}}) \quad (7)$$

Puisque, $\hat{L}_n(f_{\hat{k}}) = \min_{k \in \{1, \dots, N\}} (\hat{L}_n(f_k))$, alors

$$L(f_{\hat{k}}) - \hat{L}_n(f_{\hat{k}}) \leq \max_{k \in \{1, \dots, N\}} (L(f_k) - \hat{L}_n(f_k)) \quad (8)$$

De la même façon, puisque $\tilde{k} = \arg \min_{k \in \{1, \dots, N\}} L(f_k)$, alors

$$\hat{L}_n(f_{\tilde{k}}) - L(f_{\tilde{k}}) \leq \max_{k \in \{1, \dots, N\}} (\hat{L}_n(f_k) - L(f_k)) \quad (9)$$

Pour $\varepsilon > 0$, la borne de l'union donnera

$$\begin{aligned} P\left(\max_{k \in \{1, \dots, N\}} \left(L(f_k) - \hat{L}_n(f_k)\right) > \varepsilon\right) &\leq \\ N \times P\left(\left(L(f_k) - \hat{L}_n(f_k)\right) > \varepsilon\right) &\leq Ne^{-2n\varepsilon^2} \end{aligned}$$

et, symétriquement :

$$\begin{aligned} P\left(\max_{k \in \{1, \dots, N\}} \left(\hat{L}_n(f_k) - L(f_k)\right) > \varepsilon\right) &\leq \\ N \times P\left(\left(L(f_k) - \hat{L}_n(f_k)\right) > \varepsilon\right) &\leq Ne^{-2n\varepsilon^2} \end{aligned}$$

Finalement, par l'équation (4), puisque $\hat{L}_n(f_{\hat{k}}) - \hat{L}_n(f_{\bar{f}}) < 0$, on aura :

$$P\left(L(f_{\hat{k}}) - L(f_{\bar{f}}) > \varepsilon\right) \leq 2Ne^{-2n\varepsilon^2}.$$

Et utilisant le Lemme 2, on conclut que :

$$\begin{aligned} E\left(L(f_{\hat{k}}) - L(f_{\bar{f}})\right) &\leq E\left(\max_{k \in \{1, \dots, N\}} \left(L(f_k) - \hat{L}_n(f_k)\right)\right) + \\ E\left(\max_{k \in \{1, \dots, N\}} \left(\hat{L}_n(f_k) - L(f_k)\right) > \varepsilon\right) &\leq 2\sqrt{\frac{\ln N}{2n}}. \end{aligned}$$

Maintenant, si on considère l'oracle f^* , on aura :

$$E\left(L(f_{\hat{k}}) - L(f^*)\right) \leq L(f_{\bar{f}}) - L(f^*) + 2\sqrt{\frac{\ln N}{2n}}.$$

Le terme $L(f_{\bar{f}}) - L(f^*)$ est un terme de biais qui dépend du choix de l'ensemble $\{f_1, \dots, f_N\}$. Le terme $2\sqrt{\frac{\ln N}{2n}}$ est le terme de variance provenant du choix de $f_{\hat{k}}$ avec le risque empirique sur l'ensemble de validation.

Remarque sur le nombre N de fonctions. On peut voir que le terme de variance de l'inégalité est en $O\left(\sqrt{\frac{\log N}{2n}}\right)$, on aurait donc pu prendre un nombre de fonctions N qui soit une fonction du nombre d'observations n , par exemple $N = V \log(n)$, où V est une constante positive. On aurait alors une borne de la forme $O\left(\sqrt{\frac{V \log n}{2n}}\right)$, qui converge vers 0 quand $n \rightarrow \infty$. C'est exactement ce qui arrive si l'ensemble de fonctions n'est pas fini, mais est une classe de Vapnik. La constante V est alors liée à la dimension de Vapnik de l'ensemble de fonctions (cf [Lugosi, 2002]).

4 Convergence rapide

On peut obtenir des vitesses de convergence plus rapides grâce à l'inégalité de Bernstein :

Lemme 3 Soit $S(f) = \sum_{i=1}^n \mathbf{1}_{\{Y_i \neq f(X_i)\}}$, avec $\mathbb{E}(S(f)) = nL(f) \geq 0$, on a $|S(f) - nL(f)| \leq n$ et l'inégalité de Bernstein (cf [Lugosi, 2002]) :

Pour $t > 0$, on aura

$$P(S_n - nL(f) > t) \leq \exp\left(-\frac{t^2}{2n\sigma^2 + \frac{2}{3}t}\right) \quad (10)$$

et

$$P(nL(f) - S_n > t) \leq \exp\left(-\frac{t^2}{2n\sigma^2 + \frac{2}{3}t}\right) \quad (11)$$

4.1 Application pour un risque $L(f)$ petit.

Si, on augmente par un petit facteur $\alpha L(f)$ l'écart entre le risque empirique et le risque théorique, on peut obtenir une borne plus rapide, en $O\left(\frac{1}{n}\right)$.

Avec les notations précédentes, $E\left(n\hat{L}_n(f)\right) = nL(f) \geq 0$, $Var\left(n\hat{L}_n(f)\right) = \sigma^2 = nL(f)(1-L(f)) \leq nL(f)$ et $\left|n\hat{L}_n(f) - nL(f)\right| \leq n$. Les inégalités de Bernstein permettent d'écrire :

$$\begin{aligned} P\left(L(f) - \hat{L}_n(f) \geq t\right) &\leq \exp\left(-\frac{n^2 t^2}{2(\sigma^2 + \frac{2}{3}t)}\right) \leq \exp\left(-\frac{n^2 t^2}{2(nL(f) + \frac{2}{3}t)}\right) \\ P\left(\hat{L}_n(f) - L(f) \geq t\right) &\leq \exp\left(-\frac{n^2 t^2}{2(\sigma^2 + \frac{2}{3}t)}\right) \leq \exp\left(-\frac{n^2 t^2}{2(nL(f) + \frac{2}{3}t)}\right) \end{aligned}$$

Pour $0 < \alpha < 1$ and $1 \geq \varepsilon > 0$, posons $t = \alpha L(f) + \varepsilon$, alors $2\left(nL(f) + \frac{n}{3}t\right) \leq 2\left(\frac{1}{\alpha} + \frac{1}{3}\right)nt$ et

$$\frac{n^2 t^2}{2\left(\frac{1}{\alpha} + \frac{1}{3}\right)nt} = \frac{n(\alpha L(f) + \varepsilon)}{2\left(\frac{1}{\alpha} + \frac{1}{3}\right)} \leq \frac{n(\alpha L(f) + \varepsilon)}{2\left(\frac{3+\alpha}{3\alpha}\right)} \geq \frac{3\alpha(\alpha L(f) + 1)}{2(3 + \alpha)}n\varepsilon \geq \frac{3\alpha}{8}n\varepsilon.$$

Ainsi

$$\begin{aligned} P\left(\frac{1}{1+\alpha}\hat{L}_n(f) - L(f) > \varepsilon\right) &= \\ P\left(\hat{L}_n(f) - L(f) > \alpha L(f) + (1 + \alpha)\varepsilon\right) &\leq \exp\left(-\frac{3\alpha(1+\alpha)}{8}n\varepsilon\right) \leq \exp\left(-\frac{3\alpha(1-\alpha)}{8}n\varepsilon\right) \end{aligned}$$

De la même façon :

$$\begin{aligned} P\left(L(f) - \frac{1}{1-\alpha}\hat{L}_n(f) > \varepsilon\right) &= \\ P\left(L(f) - \hat{L}_n(f) > \alpha L(f) + (1 - \alpha)\varepsilon\right) &\leq \exp\left(-\frac{3\alpha(1-\alpha)}{8}n\varepsilon\right). \end{aligned}$$

On aura de plus :

$$\begin{aligned} L(f_{\hat{k}}) - L(f_{\hat{f}}) &= L(f_{\hat{k}}) - \frac{1}{(1-\alpha)}\hat{L}_m(f_{\hat{k}}) + \frac{1}{(1-\alpha)}\hat{L}_n(f_{\hat{k}}) - \frac{1}{(1-\alpha)}\hat{L}_n(f_{\hat{k}}) + \\ &\frac{1}{1+\alpha}\hat{L}_n(f_{\hat{k}}) - L(f_{\hat{f}}) + 2\frac{\alpha}{1-\alpha^2}\hat{L}_n(f_{\hat{k}}) \end{aligned}$$

Ainsi

$$E \left(L(f_{\hat{k}}) - L(f_{\bar{f}}) \right) \leq E \left(\max_{k \in \{1, \dots, N\}} \left(L(f_k) - \frac{1}{1-\alpha} \hat{L}_n(f_k) \right) \right) + E \left(\max_{k \in \{1, \dots, N\}} \left(\frac{1}{1+\alpha} \hat{L}_n(f_k) - L(f_k) \right) \right) + 2 \frac{\alpha}{1-\alpha^2} L(f_{\bar{f}}),$$

et pour $\alpha \in]0; 1[$,

$$E \left(L(f_{\hat{k}}) - L(f_{\bar{f}}) \right) \leq 2 \left(\frac{8(\log N + 1)}{3\alpha(1-\alpha)n} + \frac{\alpha}{1-\alpha^2} L(f_{\bar{f}}) \right).$$

Maintenant, si on considère l'oracle f^* , on aura :

$$E \left(L(f_{\hat{k}}) - L(f^*) \right) \leq \left(1 + \frac{2\alpha}{1-\alpha^2} \right) \left(L(f_{\bar{f}}) - L(f^*) \right) + 2 \left(\frac{8(\log N + 1)}{3\alpha(1-\alpha)n} + \frac{\alpha}{1-\alpha^2} L(f^*) \right).$$

4.2 Convergence rapide sous condition de bruit

Si on suppose que le bruit (l'erreur de classification pour l'oracle) a de bonne propriété, on peut obtenir des vitesses de convergence encore plus rapide.

Une application immédiate de l'inégalité de Bernstein est que si Z une variable aléatoire de moyenne μ , de variance σ^2 et bornée : $|Z - \mu| \leq b$, alors :

$$P \left(\mu - \frac{1}{n} \sum_{t=1}^n Z_i > t \right) = P \left(n\mu - \sum_{t=1}^n Z_i > nt \right) \leq \exp \left(- \frac{n^2 t^2}{2(n\sigma^2 + b \frac{n}{3} t)} \right).$$

Ainsi, pour $1 > \delta > 0$,

$$\begin{aligned} & P \left(n\mu - \sum_{t=1}^n Z_i > \sqrt{n} \sqrt{2\sigma^2 \log \left(\frac{1}{\delta} \right)} + \frac{2b \log \left(\frac{1}{\delta} \right)}{3} \right) \leq \\ & \exp \left(- \frac{\left(\sqrt{n} \sqrt{2\sigma^2 \log \left(\frac{1}{\delta} \right)} + \frac{2b \log \left(\frac{1}{\delta} \right)}{3} \right)^2}{2 \left(n\sigma^2 + \frac{b}{3} \left(\sqrt{n} \sqrt{2\sigma^2 \log \left(\frac{1}{\delta} \right)} + \frac{2b \log \left(\frac{1}{\delta} \right)}{3} \right) \right)} \right) = \\ & \exp \left(- \frac{2n\sigma^2 \log \left(\frac{1}{\delta} \right) + 2\sqrt{n} \sqrt{2\sigma^2 \log \left(\frac{1}{\delta} \right)} \frac{2b \log \left(\frac{1}{\delta} \right)}{3} + \left(\frac{2b \log \left(\frac{1}{\delta} \right)}{3} \right)^2}{2n\sigma^2 + \frac{2b}{3} \sqrt{n} \sqrt{2\sigma^2 \log \left(\frac{1}{\delta} \right)} + \frac{4b^2 \log \left(\frac{1}{\delta} \right)}{9}} \right) = \\ & \exp \left(- \log \left(\frac{1}{\delta} \right) \frac{2n\sigma^2 + \frac{4b}{3} \sqrt{n} \sqrt{2\sigma^2 \log \left(\frac{1}{\delta} \right)} + \frac{4b^2 \log \left(\frac{1}{\delta} \right)}{9}}{2n\sigma^2 + \frac{2b}{3} \sqrt{n} \sqrt{2\sigma^2 \log \left(\frac{1}{\delta} \right)} + \frac{4b^2 \log \left(\frac{1}{\delta} \right)}{9}} \right) \leq \delta \end{aligned}$$

et finalement :

$$P \left(\mu - \frac{1}{n} \sum_{t=1}^n Z_i \leq \sqrt{\frac{2\sigma^2 \log\left(\frac{1}{\delta}\right)}{n}} + \frac{2b \log\left(\frac{1}{\delta}\right)}{3n} \right) \geq 1 - \delta \quad (12)$$

On aura alors le théorème suivant :

Théorème Soit $(f_k)_{1 \leq k \leq N}$ un ensemble de fonctions obtenues grâce à un ensemble d'apprentissage. On définit :

$$\tilde{k} = \arg \min_{k \in (1, \dots, N)} L(f_k).$$

et

$$\hat{k} = \arg \min_{k \in (1, \dots, N)} \hat{L}_n(f_k).$$

On rappelle que f^* est la meilleure fonction de classification (l'oracle). Soit la fonction $w(\cdot)$ telle que pour toute fonction de classification f ,

$$\sqrt{\text{Var}(\mathbf{1}_{\{f(X) \neq f^*(X)\}})} \leq w(L(f) - L(f^*))$$

et telle que $\frac{w(x)}{\sqrt{x}}$ ne soit pas croissante. Soit τ^* la plus petite solution positive de $w(\varepsilon) = \sqrt{n\varepsilon}$, si $\theta \in]0, 1[$, alors

$$E(L(f_{\hat{k}}) - L(f^*)) \leq (1 + \theta) \left((L(f_{\tilde{k}}) - L(f^*)) + \left(\frac{8}{3n} + \frac{4\tau^*}{\theta} \right) (\log(N) + 1) \right).$$

Preuve Puisque $\left| \hat{L}_n(f_k) - \hat{L}_n(f^*) - (L(f_k) - L(f^*)) \right| \leq 2$, par l'équation (12) et la borne de l'union, avec probabilité au moins $1 - \delta$, pour toutes fonctions f_k ,

$$L(f_k) - L(f^*) \leq \hat{L}_n(f_k) - \hat{L}_n(f^*) + \sqrt{\frac{2 \log \frac{N}{\delta}}{n}} w(L(f_k) - L(f^*)) + \frac{4 \log\left(\frac{N}{\delta}\right)}{3n},$$

et

$$L(f^*) - L(f_{\tilde{k}}) \leq \hat{L}_n(f^*) - \hat{L}_n(f_{\tilde{k}}) + \sqrt{\frac{2 \log \frac{N}{\delta}}{n}} w(L(f_{\tilde{k}}) - L(f^*)) + \frac{4 \log\left(\frac{N}{\delta}\right)}{3n}.$$

En additionnant les deux inégalités, puisque, pour tout $k \in \{1, \dots, N\}$,

$\text{Var}(\mathbf{1}_{\{f_{\tilde{k}}(X) \neq f^*(X)\}}) \leq \text{Var}(\mathbf{1}_{\{f_k(X) \neq f^*(X)\}})$, on obtient, avec probabilité au moins $1 - \delta$,

$$L(f_k) - L(f_{\tilde{k}}) \leq \hat{L}_n(f_k) - \hat{L}_n(f_{\tilde{k}}) + 2\sqrt{\frac{2 \log \frac{N}{\delta}}{n}} w(L(f_k) - L(f^*)) + \frac{8 \log\left(\frac{N}{\delta}\right)}{3n}.$$

Comme, par définition, $\hat{L}_n(f_{\hat{k}}) - \hat{L}_n(f_k) \leq 0$, avec probabilité plus grande que $1 - \delta$,

$$L(f_{\hat{k}}) - L(f_k) \leq 2\sqrt{\frac{2\log \frac{N}{\delta}}{n}} w(L(f_{\hat{k}}) - L(f^*)) + \frac{8\log(\frac{N}{\delta})}{3m}.$$

Soit τ^* , comme défini dans le théorème. Si $L(f_{\hat{k}}) - L(f^*) \geq \tau^*$, alors

$$\frac{w(L(f_{\hat{k}}) - L(f^*))}{\sqrt{n}} \leq \sqrt{L(f_{\hat{k}}) - L(f^*)} \tau^*,$$

et on aura, avec probabilité au moins $1 - \delta$,

$$L(f_{\hat{k}}) - L(f_k) \leq 2\sqrt{2\log \frac{N}{\delta}} \times \sqrt{\tau^*} \times \sqrt{L(f_{\hat{k}}) - L(f^*)} + \frac{8\log(\frac{N}{\delta})}{3n}.$$

On a pour tout $\theta \in \mathbb{R}$,

$$\frac{\theta^2}{2} (L(f_{\hat{k}}) - L(f^*)) - 2\sqrt{2\log \frac{N}{\delta}} \times \sqrt{\tau^*} \times \sqrt{L(f_{\hat{k}}) - L(f^*)} \theta + 4\tau^* \log \frac{N}{\delta} \geq 0,$$

car le discriminant réduit de l'équation du second degré est nul.

On aura ainsi

$$2\sqrt{2\log \frac{N}{\delta}} \times \sqrt{\tau^*} \times \sqrt{L(f_{\hat{k}}) - L(f^*)} \leq \frac{\theta}{2} (L(f_{\hat{k}}) - L(f^*)) + \frac{4\tau^*}{\theta} \log \frac{N}{\delta},$$

ce qui implique :

$$L(f_{\hat{k}}) - L(f_k) \leq 2\sqrt{2\log \frac{N}{\delta}} \times \sqrt{\tau^*} \times \sqrt{L(f_{\hat{k}}) - L(f^*)} + \frac{8\log(\frac{N}{\delta})}{3n} \leq \frac{\theta}{2} (L(f_{\hat{k}}) - L(f^*)) + \frac{4\tau^*}{\theta} \log \frac{N}{\delta} + \frac{8\log(\frac{N}{\delta})}{3n}.$$

Ainsi, avec probabilité au moins $1 - \delta$,

$$\left(1 - \frac{\theta}{2}\right) (L(f_{\hat{k}}) - L(f^*)) \leq L(f_{\hat{k}}) - L(f^*) + \frac{4\tau^*}{\theta} \log \frac{N}{\delta} + \frac{8\log(\frac{N}{\delta})}{3n},$$

et

$$\begin{aligned} & (L(f_{\hat{k}}) - L(f^*)) \leq \\ & \frac{1}{(1-\frac{\theta}{2})} \left(L(f_{\hat{k}}) - L(f^*) + \frac{4\tau^*}{\theta} (\log N + \log(\frac{1}{\delta})) + \frac{8(\log N + \log(\frac{1}{\delta}))}{3n} \right), \end{aligned}$$

donc

$$\begin{aligned} & (L(f_{\hat{k}}) - L(f^*)) \leq \\ & \frac{1}{(1-\frac{\theta}{2})} \left(L(f_{\hat{k}}) - L(f^*) + \frac{4\tau^*}{\theta} (\log N) + \frac{8(\log N)}{3n} + \left(\frac{4\tau^*}{\theta} + \frac{8}{3n}\right) \log(\frac{1}{\delta}) \right). \end{aligned}$$

En prenant l'espérance (car \hat{k} est aléatoire),

$$E(L(f_{\hat{k}}) - L(f^*)) \leq \frac{1}{(1 - \frac{\theta}{2})} \left(L(f_{\hat{k}}) - L(f^*) + \left(\frac{4\tau^*}{\theta} + \frac{8}{3n} \right) (\log N + 1) \right)$$

Finalement, puisque si $\theta \in]0, 1[$, $\frac{1}{(1 - \frac{\theta}{2})} \leq 1 + \theta$

$$E(L(f_{\hat{k}}) - L(f^*)) \leq (1 + \theta) \left(L(f_{\hat{k}}) - L(f^*) + \left(\frac{4\tau^*}{\theta} + \frac{8}{3n} \right) (\log N + 1) \right)$$

■

Remarque Dans le cadre d'une classification binaire : $Y_k \in \{0, 1\}$, on peut expliciter la fonction w et la solution τ^* .

D'abord, définissons la condition de bruit de Mammem Tsybakov (cf par exemple [Boucheron et al., 2005]). Soit $\alpha \in [0, 1]$, et soit f^* la meilleure fonction de classification. La condition de Mammem Tsybakov est définie par : $\exists \beta > 0, \forall g \in \{0, 1\}^{\mathcal{X}}, E(\mathbf{1}_{g(X) \neq g^*(X)}) \leq \beta (L(g) - L(g^*))^\alpha$. De plus, dans le cas de la classification binaire, si $\eta(X) = E(Y = 1|X)$, et il existe $s > 0$ tel que $|2\eta(X) - 1| > s$ presque sûrement, alors le cas $\alpha = 1$ est réalisé.

Si une telle condition est vérifiée par le bruit, le théorème précédent implique le corollaire suivant :

Corollaire Supposons que condition de Mammem-Tsybakov soit vérifiéz avec un nombre α . Alors, on peut choisir $w(r) = \left(\frac{r}{h}\right)^{\frac{\alpha}{2}}$ pour une constante positive h et $\tau^* = \left(\frac{1}{nh^\alpha}\right)^{-\frac{1}{2-\alpha}}$. On déduit alors du précédent théorème :

$$E(L(f_{\hat{k}}) - L(f^*)) \leq (1 + \theta) \left((L(f_{\hat{k}}) - L(f^*)) + \left(\frac{8}{3n} + \frac{4}{\theta(nh^\alpha)^{\frac{1}{2-\alpha}}} \right) (\log(N) + 1) \right)$$

et la convergence rapide en $\frac{1}{n}$ est atteinte si $\alpha = 1$. Ce résultat permet de dire que le hold-out s'adapte au conditions de bruit (il est adaptatif), cf [Blanchard and Massart, 2006].

Références

- [Blanchard and Massart, 2006] Blanchard, G. and Massart, P. (2006). Discussion : Local rademacher complexities and oracle inequalities in risk minimization. *Annals of statistics*, 34 :2664–2671.
- [Boucheron et al., 2005] Boucheron, S., Bousquet, O., and Lugosi, G. (2005). Theory of classification : A survey of some recent advances. *ESAIM : PS*, 9 :323–375.
- [Lugosi, 2002] Lugosi, G. (2002). Pattern classification and learning theory. In Györfi, L., editor, *Principles of Nonparametric Learning*. Springer-Verlag.
- [Massart, 2003] Massart, P. (2003). *Concentration inequalities and model selection*. Ecole d'Été de Probabilités de Saint-Flour XXXIII. Springer-Verlag.