

Mathématiques du Machine Learning ERM1

Contexte: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ iid $\in X \times Y$.

$$\underbrace{R(f) = \mathbb{E} \left(\ell(y, f(x)) \right)}_{\text{Risque théorique}} \quad \underbrace{\hat{R}(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))}_{\text{Risque empirique}}$$

Minimisation du risque empirique : $\hat{f} \in \underset{f \in \mathcal{F}}{\operatorname{argmin}} \hat{R}(f)$

Question: Que peut-on dire sur $R(\hat{f}) - \left\{ R^* = \inf_f R(f) \right\}$?

I. Convexification du risque

Pour un problème de classification binaire ($Y = \{-1, 1\}$), la fonction de perte usuelle est celle du 0-1 :

$$\ell(y, f) = \mathbb{1}_{y \neq f}$$

⚠ Le problème devient donc combinatoire.

Solution Pour éviter ce problème, on commence par rajouter de la régularité en effectuant une convexification du risque

• Étape 1 : Reparamétrisation en un test de signe

On reparamétrise $f: X \rightarrow \{-1, 1\}$ par $g: X \rightarrow \mathbb{R}$:

$$\forall x \in X, f(x) = \operatorname{sgn}(g(x)) \text{ où } \operatorname{sgn}(g) = \begin{cases} 1 & \text{si } g > 0 \\ -1 & \text{si } g < 0 \\ \text{Unif}(\{-1, 1\}) & \text{si } g = 0 \end{cases}$$

remarque: Pour gérer le cas en 0, on a remplacé une fonction par une fonction

1. alternative. En pratique, ce ne sera pas un problème car on ne devrait jamais tomber sur 0.

$$\begin{aligned}
 \text{On a donc } \mathcal{R}(f) &= \mathbb{E}(\mathbb{1}_{y(x) \neq 0} \mathbb{1}_{f(x) \neq y}) + \mathbb{E}(\mathbb{1}_{y(x)=0} \mathbb{1}_{f(x) \neq y}) \\
 &= \mathbb{E}(\mathbb{1}_{y(x)y < 0}) + \frac{1}{2} \mathbb{E}(\mathbb{1}_{y(x)=0}) \\
 &= \mathbb{E}(\Phi_{0-1}(y y(x))) \text{ ou } \Phi_{0-1}(u) = \begin{cases} 1 & \text{si } u < 0 \\ \frac{1}{2} & \text{si } u = 0 \\ 0 & \text{si } u > 0 \end{cases}
 \end{aligned}$$

• Étape 2: On remplace Φ_{0-1} par une fonction avec de meilleures propriétés numériques.

"perte Hinge"
 $\Phi(u) = (1-u)_+$

$\Phi_{0-1}(u)$



II. Décomposition du risque

$$M(\hat{f}) - M^* = \underbrace{\left\{ M(\hat{f}) - \inf_{f' \in F} M(f') \right\}}_{\text{Erreur d'estimation}} + \underbrace{\left\{ \inf_{f' \in F} M(f') - M^* \right\}}_{\text{Erreur d'approximation}}$$

Question: Comment se comporte l'erreur d'estimation ?

Notons $f_{\text{opt}} \in \arg \min_{f' \in F} M(f')$: meilleur estimateur dans la classe F .

$$M(\hat{f}) - M(f_{\text{opt}}) = \underbrace{\left\{ M(\hat{f}) - \hat{M}(\hat{f}) \right\}}_{\text{Erreur d'estimation}} + \underbrace{\left\{ \hat{M}(\hat{f}) - \hat{M}(f_{\text{opt}}) \right\}}_{\text{Erreur d'approximation}} + \underbrace{\left\{ \hat{M}(f_{\text{opt}}) - M(f_{\text{opt}}) \right\}}_{\text{Erreur d'optimisation}}$$
$$\leq \boxed{2 \sup_{f \in F} |\hat{M}(f) - M(f)|} + \text{Erreur d'optimisation} \quad (*)$$

L'erreur d'optimisation dépend de l'algorithme d'entraînement utilisé.
Il peut souvent être aussi petit que l'on veut. Le reste de cette leçon étudie le premier terme.

III. Concentration : Méthode des différences bornées

Une manière de contrôler (*) est d'utiliser la méthode générique des différences bornées:

Soient X_1, \dots, X_n des variables aléatoires indépendantes sur \mathcal{X} et $f: \mathcal{X}^n \rightarrow \mathbb{R}$ tq $\underbrace{X_i}_{X_i}$ $\underbrace{X_i'}_{X_i'}$
 $\forall i, \forall (x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n), (x_1, \dots, x_{i-1}, x_i', x_{i+1}, \dots, x_n)$ dans le support de la distribution de (X_1, \dots, X_n) ,

$$|f(x_i) - f(x_i')| \leq c_i.$$

Question: Que peut-on dire de $f(X_1, \dots, X_n)$ vs $\mathbb{E}(f(X_1, \dots, X_n))$?

1) Inégalité de Markov.

Soit A une variable positive dans L^1 et $a > 0$.

$$a \mathbb{1}_{A \geq a} \leq A \text{ (p.s.)} \Rightarrow \mathbb{E}(a \mathbb{1}_{A \geq a}) \leq \mathbb{E}(A) \\ \Rightarrow \mathbb{P}(A \geq a) \leq \frac{\mathbb{E}(A)}{a}$$

2) Transformation en martingale

$$\text{Notons } V_i, V_i = \mathbb{E}(f(X_1, \dots, X_n) | X_1, \dots, X_i) - \underbrace{\mathbb{E}(f(X_1, \dots, X_n) | X_1, \dots, X_{i-1})}_{= \mathbb{E}(f(X_1, \dots, X_n)) \text{ si } i=1}$$

$$\text{Notons } M_i = \sum_{j=1}^i V_j.$$

• M_i est une martingale pour la filtration $(\sigma(X_1, \dots, X_{i-1}))_i$.

— $M_1 = f(X_1) - \mathbb{E}(f(X_1, \dots, X_n))$ — $M_2 = f(X_1, X_2) - \mathbb{E}(f(X_1, \dots, X_n))$ — $M_3 = f(X_1, X_2, X_3) - \mathbb{E}(f(X_1, \dots, X_n))$ —

$$\begin{aligned}
 \text{En effet, } E(V_i | X_1, \dots, X_{i-1}) &= E(E(f(X_{1:n}) | X_{1:i}) | X_{1:i-1}) \\
 &\quad - E(E(f(X_{1:n}) | X_{1:i-1}) | X_{1:i-1}) \\
 &= E(f(X_{1:n}) | X_{1:i-1}) - E(f(X_{1:n}) | X_{1:i-1}) \\
 &= 0.
 \end{aligned}$$

$$\underline{f(X_1, \dots, X_n) - E(f(X_1, \dots, X_n)) = \sum_i V_i}$$

3) Fonction génératrice des moments et règle de récurrence

$$\begin{aligned}
 P(f - E(f) \geq t) &\leq P(\exp(f - E(f)) \geq \exp(st)) \\
 &\stackrel{\text{Markov}}{\leq} e^{-st} E\left(\exp\left(s \sum_i V_i\right)\right) \\
 &= e^{-st} E\left(\prod_i \exp(s V_i)\right) \\
 &= e^{-st} E\left(\left(\prod_{i=1}^{n-1} \exp(s V_i)\right) \exp(s V_n)\right) \\
 &= e^{-st} E\left(\left(\prod_{i=1}^{n-1} \exp(s V_i)\right) E(\exp(s V_n) | X_{1:n-1})\right)
 \end{aligned}$$

Gn va travailler là dessus (*)
et généraliser par récurrence.

4) Borner les incréments

Lemme : $\forall i, \exists \alpha_i \leq \beta_i$ t.q. $V_i \in [\alpha_i, \beta_i]$ p.s. et $|\beta_i - \alpha_i| \leq c_i$.
 \swarrow
 $X_{1:i-1}$ mesurables

Preuve : $V_i = E(f | X_{1:i}) - E(f | X_{1:i-1})$

$$\begin{aligned}
&= \int \int (x_1, \dots, x_i, x_{i+1}, \dots, x_n) P_{|X_{1:i}}(dx_{i+1:n}) \\
&\quad - \int \int (x_1, \dots, x_{i-1}, x_i, \dots, x_n) P_{|X_{1:i-1}}(dx_{i:n}) \\
&\stackrel{(\text{Indépendance})}{=} \int \underbrace{\left(\int (x_{1:i}, x_{i+1:n}) - \int (x_{1:i-1}, x_{i:n}) P(dx_i) \right)}_{:= \Delta_{|X_{1:i-1}}(x_i)} P(dx_{i+1:n})
\end{aligned}$$

et notons $\alpha_i := \sup_{x \in \text{supp}(X_i)} \Delta_{|X_{1:i-1}}(x)$

et $\beta_i := \inf_{x \in \text{supp}(X_i)} \Delta_{|X_{1:i-1}}(x)$.

Alors clairement $\alpha_i \leq V_i \leq \beta_i$

de plus,

$$|\beta_i - \alpha_i| = \sup_{x, x' \in \text{supp}(X_i)} \underbrace{\int (x_{1:i-1}, x, x_{i+1:n}) - \int (x_{1:i-1}, x', x_{i+1:n})}_{\leq c_i} \underbrace{P(dx_{i+1}) \dots P(dx_n)}_{\leq c_i}$$

$\leq c_i$ □

Corollaire : $E\left(\prod_i \exp(s V_i)\right) \leq \exp\left(s^2 \frac{\sum_i c_i^2}{8}\right)$

Preuve : Comme expliqué plus haut, on se focalise sur (*), le reste étant obtenu par récurrence

Montrons que $E(\exp(s V_n) | X_{1:n-1}) \leq \exp(s^2 c_i^2 / 8)$.

Notons, $\varphi(s) = P_n(E(\exp(s V_n) | X_{1:n-1}))$

Alors,
$$\varphi'(s) = \frac{\mathbb{E}(V_n \exp(s V_n) | X_{1:n-1})}{\mathbb{E}(\exp(s V_n) | X_{1:n-1})}$$

et,
$$\varphi''(s) = \frac{\mathbb{E}(V_n^2 \exp(s V_n) | X_{1:n-1})}{\mathbb{E}(e^{s V_n} | X_{1:n-1})} - \left(\frac{\mathbb{E}(V_n \exp(s V_n) | X_{1:n-1})}{\mathbb{E}(\exp(s V_n) | X_{1:n-1})} \right)^2$$

et on reconnaît l'expression de la variance de V_n sous la loi conditionnelle à $X_{1:n}$ de densité $v \mapsto \exp(s v) / \mathbb{E}(\exp(s V_n) | X_{1:n-1})$.

d'où
$$\varphi''(s) = \arg \min_{\xi} \int (V_n - \xi)^2 \times \text{densité modifiée}(V_n) dV_n$$

$$\leq \frac{(\beta_n - \alpha_n)^2}{4} \text{ en prenant } \xi = \frac{\beta_n + \alpha_n}{2}$$

donc, d'après Taylor:
$$\varphi(s) \leq \frac{(\beta_i - \alpha_i)^2}{8} s^2 \quad \square$$

4) Conclusion

$$\mathbb{P}\left(\sum_i V_i \geq t\right) \leq \inf_{s > 0} \exp(-st) \exp\left(s^2 \frac{\sum_i c_i^2}{8}\right)$$

et avec $s^2 = \frac{4t}{\sum_i c_i^2}$, on obtient

$$\mathbb{P}\left(\sum_i V_i - \mathbb{E}\left(\sum_i V_i\right) \geq t\right) \leq \exp\left(-\frac{2t^2}{\sum_i c_i^2}\right)$$

IV. Application : ERM pour un ensemble d'hypothèses fini

Hypothèse 1: $P(y, f(z)) \in [0, 1] \quad \forall f \in F, \forall (z, y) \in S_{\text{app.}}$

Hypothèse 2: $|F| < +\infty$.

Par borne d'union,

$$\begin{aligned} P\left(\sup_{f \in F} |\hat{M}(f) - M(f)| \geq t\right) &\leq \sum_{f \in F} P(|\hat{M}(f) - M(f)| \geq t) \\ &\leq \sum_{f \in F} 2 \exp\left(-\frac{2nt^2}{L^2}\right) \\ &\leq \underbrace{2|F|}_{= \delta} \exp\left(-\frac{2nt^2}{L^2}\right) \end{aligned}$$

et donc, avec probabilité au moins $1 - \delta$,

$$\sup_{f \in F} |\hat{M}(f) - M(f)| \leq \frac{L}{\sqrt{2n}} \sqrt{\log\left(\frac{2|F|}{\delta}\right)}$$