

# Mathématiques du Machine Learning ERM2

## Contexte Contrôle de l'erreur d'estimation

$$\sup_{f \in F} |\hat{M}(f) - M(f)| \leq \sup_{f \in F} \hat{M}(f) - M(f) + \sup_{f \in F} M(f) - \hat{M}(f).$$

## I. Lien entre Forte Proba - Espérance.

Hypothèse:  $\forall f \in F, \forall (x, y) \in \text{Support}, 0 \leq P(y, f(\cdot)) \leq P_0$ .

Membre:  $\forall f \in F, \hat{M}(f) = \frac{1}{n} \sum_{i=1}^n P(y_i, f(\cdot))$ .

en changeant un couple  $(x_i, y_i)$  par n'importe quel autre couple dans le support,  $\hat{M}(f)$  change au plus de  $P_0/n$ .

Donc  $\sup_{f \in F} \hat{M}(f) - M(f)$  change également au plus de  $P_0/n$

et  $\sup_{f \in F} M(f) - \hat{M}(f)$  change également au plus de  $P_0/n$ .

On peut donc appliquer la méthode des différences bornées :

$\forall \delta \in (0, 1)$ , chacune des inégalités suivantes est vérifiée avec probabilité au moins  $1 - \delta$  :

$$\left| \sup_{f \in F} (M(f) - \hat{M}(f)) - \mathbb{E} \left( \sup_{f \in F} (M(f) - \hat{M}(f)) \right) \right| \leq \frac{P_0}{\sqrt{2n}} \sqrt{\log \left( \frac{1}{\delta} \right)}$$

$$\left| \sup_{f \in F} (\hat{M}(f) - M(f)) - \mathbb{E} \left( \sup_{f \in F} (\hat{M}(f) - M(f)) \right) \right| \leq \frac{P_0}{\sqrt{2n}} \sqrt{\log \left( \frac{1}{\delta} \right)}$$

Donc : Pour obtenir des bornes avec facile proba, on peut se contenter de travailler sur bornes les espérances

$$\mathbb{E} \left( \sup_{f \in F} (M(f) - \hat{M}(f)) \right) \text{ et } \mathbb{E} \left( \sup_{f \in F} (\hat{M}(f) - M(f)) \right).$$

\*

## II. Complexité de Macdonald

Notons  $\forall i, z_i = (x_i, y_i)$  et  $z = (x, y)$ .

Symétrisation Soient  $z'_1, \dots, z'_n$  iid, de même loi que  $z$  et indépendantes de  $z_1, \dots, z_n, z$  (des copies indépendantes des données).

$$\begin{aligned} \text{Nous avons, } \forall f, \quad \hat{M}(f) - M(f) &= \frac{1}{n} \sum_{i=1}^n P(y_i, f(x_i)) - \mathbb{E}_z (P(y, f(x))) \\ &= \frac{1}{n} \sum_{i=1}^n P(y_i, f(x_i)) - \mathbb{E}_{z'_1, \dots, z'_n} \left( \frac{1}{n} \sum_{i=1}^n P(y'_i, f(x'_i)) \right) \\ &= \mathbb{E}_{z'_1, \dots, z'_n} \left( \underbrace{\frac{1}{n} \sum_{i=1}^n P(y_i, f(x_i))}_{\text{}} - \underbrace{\frac{1}{n} \sum_{i=1}^n P(y'_i, f(x'_i))}_{\text{}} \right) \end{aligned}$$

$$\text{donc, } \sup_{f \in F} (\hat{M}(f) - M(f)) = \sup_{f \in F} //$$

$$\leq \mathbb{E}_{z'_1, \dots, z'_n} \left( \sup_{f \in F} \left( \frac{1}{n} \sum_{i=1}^n P(y_i, f(x_i)) - \frac{1}{n} \sum_{i=1}^n P(y'_i, f(x'_i)) \right) \right)$$

$$\text{et enfin, } (*) \leq \mathbb{E}_{\substack{z_1, \dots, z_n \\ z'_1, \dots, z'_n}} \left( \sup_{f \in F} \underbrace{\frac{1}{n} \sum_{i=1}^n P(y_i, f(x_i))}_{\text{}} - \underbrace{\frac{1}{n} \sum_{i=1}^n P(y'_i, f(x'_i))}_{\text{}} \right)$$

Loi symétrique

Nous avons quasiment terminé! Soient  $\varepsilon_1, \dots, \varepsilon_n$  des variables aléatoires iid de loi  $\text{Rad}(1/2)$  (=1 avec proba  $\frac{1}{2}$ , =-1 avec proba  $\frac{1}{2}$ ) qui sont indépendantes des autres quantités du cours.

Remarquons que  $\forall i, \underbrace{p(y_i, f(x_i)) - p(y_i', f(x_i'))}_{\alpha_i}$  a une loi symétrique.

Alors,  $\forall i, \alpha_i$  a la même loi que  $\varepsilon_i \alpha_i$ .

$$\begin{aligned} \text{Alors } (*) &\leq \mathbb{E}_{\substack{z_1, \dots, z_n \\ z_1', \dots, z_n' \\ \varepsilon_1, \dots, \varepsilon_n}} \left( \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \left( p(y_i, f(x_i)) - p(y_i', f(x_i')) \right) \right) \\ &\leq \mathbb{E}_{\substack{z_1, \dots, z_n \\ \varepsilon_1, \dots, \varepsilon_n}} \left( \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i p(y_i, f(x_i)) \right) \\ &\quad + \mathbb{E}_{\substack{z_1', \dots, z_n' \\ \varepsilon_1, \dots, \varepsilon_n}} \left( \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i p(y_i', f(x_i')) \right) \\ &= 2 \mathbb{E}_{\substack{z_1, \dots, z_n \\ \varepsilon_1, \dots, \varepsilon_n}} \left( \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i p(y_i, f(x_i)) \right). \end{aligned}$$

De manière générale, si  $\mathcal{H}$  est une classe de fonctions  $h$  et étant donné une observation  $S = (z_1, \dots, z_n)$ , on définit

$$R_S(\mathcal{H}) = \mathbb{E}_{\varepsilon_1, \dots, \varepsilon_n} \left( \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i h(z_i) \right).$$

On veut de prouver le résultat suivant:

Théorème (Symétrisation):

$$0 \leq \mathbb{E} \left( \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i h(z_i) \right) \leq \mathbb{E} \left( \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n h(z_i) \right)$$

Si  $H$  est une famille de fonctions,  $z_1, \dots, z_n$  iid,

$$\mathbb{E}_{(z_1, \dots, z_n)} \sup_{h \in H} \left( \frac{1}{n} \sum_{i=1}^n h(z_i) - \mathbb{E} h(z) \right) \leq 2 \underbrace{\mathbb{E}_{(z_1, \dots, z_n)} \left( \underbrace{M_{(z_1, \dots, z_n)}(H)}_{:= M_n(H)} \right)}_{:= M_n(H)}$$

et

$$\mathbb{E}_{(z_1, \dots, z_n)} \sup_{h \in H} \left( \mathbb{E} h(z) - \frac{1}{n} \sum_{i=1}^n h(z_i) \right) \leq 2 M_n(H).$$

$M_S(H)$  = Complexité de Rademacher de  $H$  conditionnellement à  $S$

$M_n(H)$  = Complexité de Rademacher de  $H$  sous la loi des données.

Proposition : Si  $g$  est  $G$ -Lipschitz,  $M(g(H)) \leq G M(H)$   
avec  $g(H) = \{g \circ h; h \in H\}$ .

Preuve : Admis, voir Buch 2024.

### III. Analyse des modèles linéaires

Modèles linéaires contraindre:  $F = \left\{ f_\theta : x \mapsto \theta^T \varphi(x), \Omega(\theta) \leq D \right\}$   
 $\Omega$  : norme sur  $\mathbb{R}^d$ .

$$M_n(F) = \mathbb{E}_{\substack{z_1, \dots, z_n \\ \varepsilon_1, \dots, \varepsilon_n}} \left( \sup_{\substack{\theta \\ \Omega(\theta) \leq D}} \left( \frac{1}{n} \sum_{i=1}^n \varepsilon_i \theta^T \varphi(z_i) \right) \right)$$

$\varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$

$$= \mathbb{E} \left( \sup_{\Omega(\phi) \leq 0} \frac{1}{n} \varepsilon^T \Phi \phi \right) \text{ avec } \begin{matrix} c = (\varepsilon_i) \\ \Phi = \begin{pmatrix} \ell(a_1)^T \\ \vdots \\ \ell(a_n)^T \end{pmatrix} \end{matrix}$$

Case  $\Omega = \mathbb{H}_2$ :

$$\mu_n(f) = \frac{D}{n} \mathbb{E} \left( \|\Phi^T \varepsilon\|_2 \right)$$

$$\leq \frac{D}{n} \sqrt{\mathbb{E} \left( \left\| \Phi^T \Sigma \right\|_2^2 \right)}$$

$$= \frac{D}{n} \sqrt{\mathbb{E} \left( \varepsilon^T \Phi^T \Phi \varepsilon \right)}$$

$$= \frac{D}{n} \sqrt{\mathbb{E} \left( \text{tr} \left( \varepsilon^T \Phi^T \Phi \varepsilon \right) \right)}$$

$$= \frac{D}{n} \sqrt{\mathbb{E} \left( \text{tr} \left( \mathbb{I}^T \mathbb{I} \Sigma \Sigma^T \right) \right)}$$

$$= \frac{D}{n} \int \mathbb{E} \left( t_k \left( \Phi^T \Phi \mathbb{E} \left( \Sigma \Sigma^T \mid \Phi \right) \right) \right)$$

$$= \frac{0}{n} \sqrt{\mathbb{E}(t_r(\Phi^T \Phi))} = I$$

$$= \frac{D}{n} \sqrt{\mathbb{E} \left( \sum_i \|\varphi(z_i)\|^2 \right)}$$

$$= \frac{0}{\sqrt{n}} \sqrt{\mathbb{E} \|\varphi(\omega)\|^2}$$

$\partial \bar{\partial} = -\partial^2 = -\bar{\partial}^2$

Conclusion: Si  $\forall g$  dans le support,  $\ell(g, \cdot)$  et  $G$  lipschitz,

$$\mathbb{E}(\mathcal{M}(\hat{P}_\theta)) \leq \inf_{\| \theta \|_2 \leq D} \mathcal{M}(P_\theta) + \frac{4GD}{\sqrt{n}} \sqrt{\mathbb{E} \|\ell(z)\|^2}.$$