

①

Leçon n°3: EMMI; Risques convexifiés; Décomposition du risque;  
Inégalités de concentration; Applications pour EMM

Contexte:  $(x_1, y_1), \dots, (x_n, y_n)$  i.i.d e  $X \times Y$

$$R(f) = \mathbb{E}(P(y, f(x))) \quad (x, y) \perp \sigma((x_1, y_1), \dots, (x_{n-1}, y_{n-1}))$$

de même distribution que  $(x_1, y_1)$ .

EMM:  $\hat{f} \in \operatorname{argmin}_f \left( \hat{R}(f) = \frac{1}{n} \sum_{i=1}^n P(x_i, y_i) \right)$

$\uparrow$   
 $\approx R(f)$

Question: Peut-on contrôler  $R(\hat{f}) - R^*$  ?

I. Convexification du risque

Pour des problèmes structurés (par exemple la classification binaire  $y = \{-1, 1\}$ ), le problème de minimisation du risque empirique ~~adec~~ par la perte du 0-1 (i.e.  $P(y, f) = \mathbb{1}_{y \neq f}$ ) ~~est~~ est souvent NP difficile.

Pour cette raison (et potentiellement des considérations d'optimisation et de stabilité), il peut être judicieux de changer la fonction de perte  $P$  de sorte à satisfaire certaines propriétés.

② Propriétés désirables:

- Différentiab. p.l. (SGD, ...)
- Smoothness
- Convexité
- ...

Étape 1: On reparamétrise  $f: X \rightarrow \{-1, 1\}$  par  $g: X \rightarrow \mathbb{R}$  de la manière suivante

$$\forall x \in X, f(x) = \text{sgn}(g(x))$$

$$\text{et } \text{sgn}(y) = \begin{cases} 1 & \text{si } y > 0 \\ -1 & \text{si } y < 0 \\ \underbrace{U_{\text{unif}}(\{-1, 1\})}_{\perp \text{ des autres qualités}} & \text{si } y = 0 \end{cases}$$

remarque:  $\text{sgn}(\cdot)$  est en fait un noyau conditionnel de probas.

On voit alors que  $R(g) \equiv R(f)$   
abus de notations

$$= \mathbb{E} \left( \mathbb{1}_{g(x) \neq 0} \mathbb{1}_{f(x) \neq y} \right) + \mathbb{E} \left( \mathbb{1}_{g(x) = 0} \mathbb{1}_{f(x) \neq y} \right)$$

$$= \mathbb{E} \left( \mathbb{1}_{y g(x) < 0} \right) + \frac{1}{2} \mathbb{E} \left( \mathbb{1}_{g(x) = 0} \right)$$

$$= \mathbb{E} \left( \Phi_{0,1}(y g(x)) \right)$$

③

$$\text{où } \Phi_{0,1}(u) = \begin{cases} 1 & \text{si } u < 0 \\ \frac{1}{2} & \text{si } u = 0 \\ 0 & \text{si } u > 0 \end{cases}$$

Étape 2: On remplace  $\Phi_{0,1}$  par une fonction avec de meilleures propriétés numériques

$$M_{\Phi}(y) = \mathbb{E}(\Phi(yg(\cdot)))$$

Exemples:

• Perceptron quadratique  $\Phi(u) = (u-1)^2$

obs:  $\Phi(yg(z)) = (yg(z) - 1)^2$   
 $= (yg(z) - y)^2$   
 $= y^2 (g(z) - y)^2$   
 $= (g(z) - y)^2$

• Hinge  $\Phi(u) = \max(1-y, 0)$  (c.f. SVM)

• Exponential Loss:  $\Phi(u) = \exp(-u)$  (c.f. AdaBoost)

• Logistic Loss:  $\Phi(u) = \log(1 + e^{-u})$

$\hookrightarrow \Phi(yg(z)) = \log(1 + e^{-yg(z)})$   
 $= -\log\left(\frac{1}{1 + e^{yg(z)}}\right)$   
 $= -\log(\sigma(yg(z)))$

⑨

avec  $\sigma(v) = \frac{1}{1+e^{-v}}$  : fonction sigmoïde.

remarque:  $\mathbb{P}(y|z) = -\log(\underbrace{\sigma(yg(z))}_{\text{vraisemblance conditionnelle}})$  peut

être vu comme l'opposé de la vraisemblance conditionnelle dans le modèle probabiliste

$$\mathbb{P}(y=1|z) = \sigma(g(z))$$

$$\mathbb{P}(y=0|z) = 1 - \sigma(g(z))$$

Pourquoi  $M_{\mathbb{P}}(g) \approx M_{\mathbb{P}_{0,1}}(g)$  ?

↳ voir Bach 2024 (section 4.1).

## II - Décomposition du risque

$$\hat{f} \in \underset{f \in \mathcal{F}}{\operatorname{argmin}} \left( \hat{M}(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(z_i)) \right)$$

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

⑤

• L'erreur d'approximation :

- Est déterministe
- Dépend de la régularité du pb et de la classe  $F$

• L'erreur d'estimation :

- Est de nature stochastique

Décomposition de l'erreur d'estimation :

Notons  $p_{opt}$  e argmin  $M(p)$   $p \in F$

Alors  $M(\hat{p}) - M(p_{opt})$

$$\leq \underbrace{(M(\hat{p}) - \hat{M}(\hat{p})) + (\hat{M}(\hat{p}) - \hat{M}(p_{opt}))}_{\text{erreur d'opt.}} + \underbrace{(\hat{M}(p_{opt}) - M(p_{opt}))}_{\checkmark}$$

$$\leq 2 \sup_{p \in F} |\hat{M}(p) - M(p)| + \text{Erreur d'optimisation.}$$

Remarque :  $\hat{M}(p) - M(p) = \hat{M}(p) - E(M(p))$

$$= \frac{1}{n} \sum \dots$$

⑥

### III. Concentration : Méthode des différences bornées.

- Soient  $X_1, \dots, X_n$  des V.A.s indépendantes (dans  $\mathcal{X}$ )

- Soit  $f$  une fonction de  $\mathcal{X}^n$  dans  $\mathbb{R}$  tq,

$\forall (x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n), (x_1, \dots, x_{i-1}, x_i', x_{i+1}, \dots, x_n)$   
dans le support de la distribution de  $(X_1, \dots, X_n)$ ,

$$|f(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n) - f(x_1, \dots, x_{i-1}, x_i', x_{i+1}, \dots, x_n)| \leq c_i$$

Que peut-on dire de  $f(X_1, \dots, X_n)$  vs  $E(f(X_1, \dots, X_n))$ .

#### 1) Inégalité de Markov

Sait  $A$  une VA <sup>positive</sup> dans  $L^1$ . Alors  $\forall a > 0$ ,

~~$$E(A) = E(A \mathbb{1}_{A < a}) + E(A \mathbb{1}_{A \geq a})$$~~

$$E(\mathbb{1}_{A \geq a}) \leq E(A) \quad (\text{car p.s. } a \mathbb{1}_{A \geq a} \leq A)$$

$$\text{i.e. } P(A \geq a) \leq \frac{E(A)}{a}$$

#### 2) Transformation et somme de martingale

$$\text{Notons } V_i = E(f(X_1, \dots, X_n) | X_1, \dots, X_i)$$

$$- E(f(X_1, \dots, X_n) | X_1, \dots, X_{i-1})$$

⑦ Alors:

• Lemme:  $(\sum_{k=1}^i V_k)$  est une martingale pour la filtration  $(\mathcal{F}(X_1, \dots, X_{i-1}))$ .

preuve:  $E(V_{i+1} | X_1, \dots, X_i)$

~~$E(\sum_{k=1}^{i+1} V_k | X_1, \dots, X_i)$~~

$$= E(E(\sum_{k=1}^{i+1} V_k | X_1, \dots, X_{i+1}) - E(\sum_{k=1}^i V_k | X_1, \dots, X_i) | X_1, \dots, X_i)$$

$$= E(E(\sum_{k=1}^{i+1} V_k | X_1, \dots, X_{i+1}) | X_1, \dots, X_i)$$

$$- E(E(\sum_{k=1}^i V_k | X_1, \dots, X_i) | X_1, \dots, X_i)$$

$$= E(\sum_{k=1}^{i+1} V_k | X_1, \dots, X_i) - E(\sum_{k=1}^i V_k | X_1, \dots, X_i)$$

$$= 0$$

□

• De plus,  $\boxed{\sum_{i=1}^n V_i - E(\sum_{i=1}^n V_i) = \sum_{i=1}^n V_i - \sum_{i=1}^n E(V_i)}$

3) Fonction génératrice des moments et règle de récurrence

Soit  $t \geq 0$ , Soit  $s \geq 0$ , exp. croissante

$$P(\sum_{i=1}^n V_i \geq t) \leq P(\exp(s(\sum_{i=1}^n V_i)) \geq e^{st})$$

$$\leq e^{-st} E(e^{s \sum_{i=1}^n V_i})$$

Markov

③

De plus,

$$\begin{aligned} \mathbb{E} \left( e^{s \sum_{i=1}^n V_i} \right) &= \mathbb{E} \left( \prod_{i=1}^n e^{s V_i} \right) \\ &= \mathbb{E} \left( \mathbb{E} \left( \prod_{i=1}^{n-1} e^{s V_i} \right) e^{s V_n} \mid X_1, \dots, X_{n-1} \right) \\ &= \mathbb{E} \left( \prod_{i=1}^{n-1} e^{s V_i} \right) \mathbb{E} \left( e^{s V_n} \mid X_1, \dots, X_n \right) \end{aligned}$$

④ 9) Borne des incréments

Lemme:  $\forall i, \exists \alpha_i \leq \beta_i$  t.q.  $V_i \in [\alpha_i, \beta_i]$  p.s. et  $\beta_i - \alpha_i \leq \epsilon_i$ .

Preuve:

$$\begin{aligned} V_i &= \mathbb{E}(P \mid X_1, \dots, X_i) - \mathbb{E}(P \mid X_1, \dots, X_{i-1}) \\ &= \mathbb{E} \left( \int P(X_1, \dots, X_i, X_{i+1}, \dots, X_n) P^{(X_1, \dots, X_i)}(dx_{i+1}, \dots, dx_n) \right. \\ &\quad \left. - \int P(X_1, \dots, X_{i-1}, X_i, \dots, X_n) P^{(X_1, \dots, X_{i-1})}(dx_i, \dots, dx_n) \right) \\ &\stackrel{\text{Indep}}{\rightarrow} \mathbb{E} \left( \int P(X_1, \dots, X_i, X_{i+1}, \dots, X_n) P(dx_{i+1}) \dots P(dx_n) \right. \\ &\quad \left. - \int P(X_1, \dots, X_{i-1}, X_i, \dots, X_n) P(dx_i) \dots P(dx_n) \right) \\ &\quad \left. \mid X_1, \dots, X_{i-1} \right) \end{aligned}$$



⑨

$$= E \left( \int \left( P(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n) - \int P(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n) P(dx_i) P(dx_{i+1}) \dots P(dx_n) \right) \middle| x_1, \dots, x_{i-1} \right)$$

De plus,  $f$  est bornée sur  $\text{Supp}(X_1) \times \dots \times \text{Supp}(X_n)$   
(conséquence directe de l'hypothèse sur  $f$ .)

Alors en notant

$$\beta_i = E \left( \sup_{z \in \text{Supp}(X_i)} \int P(x_1, \dots, x_{i-1}, z, x_{i+1}, \dots, x_n) \right)$$

$$\text{et } \alpha_i = \inf_{z \in \text{Supp}(X_i)} \int P(x_1, \dots, x_{i-1}, z, x_{i+1}, \dots, x_n)$$

nous avons  $-\infty < \alpha_i \leq \beta_i < +\infty$  et  $\beta_i \leq c_i$ .

$$(\beta_i - \alpha_i) = E \left( \sup_{z \in \text{Supp}(X_i)} \int \left( P(x_1, \dots, x_{i-1}, z, x_{i+1}, \dots, x_n) - \int P(x_1, \dots, x_{i-1}, z', x_{i+1}, \dots, x_n) P(dx_i) \dots P(dx_n) \right) \middle| x_1, \dots, x_{i-1} \right)$$

$\leq c_i$ .

□

(10) Condition:  $E(\prod_i e^{sV_i}) \leq \exp(s^2 \sum_{i=1}^n c_i^2 / 8)$

preuve: Soit  $i$ . Notés  $\varphi(s) = E(e^{sV_i} | V_1, \dots, V_{i-1})$ .

Alors  $\varphi'(s) = \frac{E(V_i e^{sV_i} | V_1, \dots, V_{i-1})}{E(e^{sV_i} | V_1, \dots, V_{i-1})}$

et  $\varphi''(s) = \frac{E(V_i^2 e^{sV_i} | V_1, \dots, V_{i-1})}{E(e^{sV_i} | \dots)} - \left( \frac{E(V_i e^{sV_i} | \dots)}{E(e^{sV_i} | \dots)} \right)^2$

et on reconnaît l'expression de la variance de  $V_i$  conditionnelle pour densité conditionnelle à  $X_1, \dots, X_{i-1}$ :  $x \mapsto \frac{e^{sx}}{E(e^{sV_i} | X_1, \dots, X_{i-1})}$

d'où:  $\varphi''(s) = \int (V_i - \bar{s})^2 \times \text{densité conditionnelle}(V_i, \dots, V_{i-1})$   
 $\leq \frac{(\beta_i - \alpha_i)^2}{4}$  (en prenant  $\bar{s} = \text{milieu du segment}$   
 $= \frac{\beta_i + \alpha_i}{2}$ ).

\* Le résultat est alors immédiat par la formule de récursion

\* d'où, d'après Taylor Lagrange:

$$\varphi(s) \leq \frac{(\beta_i - \alpha_i)^2 s^2}{8}$$

(11)

c) Conclusion

$$P(\sum_i V_i \geq t) \leq \inf_{s>0} e^{-st} e^{s^2 \frac{\sum_i (\beta_i - \alpha_i)^2}{8}}$$

et minimiser et  $s \left( s^2 = \frac{4t}{\sum_i (\beta_i - \alpha_i)^2} \right)$  donc

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

IV. Application I: EM M avec F fini.

Hypothèse 1:  $P(y, f(x)) \in [0, P_\infty] \forall f \in F, \forall (x, y) \in \text{Supp}$

Hypothèse 2:  $|F| \leq +\infty$

Alors, par borne d'union,

$$P\left(\sup_{f \in F} |\hat{M}(f) - M(f)| \geq t\right) \leq \sum_{f \in F} P(|\hat{M}(f) - M(f)| \geq t)$$

$$\leq \sum_{f \in F} 2 \exp\left(-2nt^2/P_\infty^2\right)$$

McDermid

$$\leq 2|F| \exp\left(-2nt^2/P_\infty^2\right)$$

(12)

Donc, on note  $\delta = 2|f| \exp(-2nt^2/P_\infty^2)$ , avec proba  $1-\delta$ ,

$$\sup_{f \in F} |\hat{M}(f) - M(f)| \leq t = \frac{P_\infty}{\sqrt{2n}} \sqrt{P_y \frac{2|f|}{\delta}}$$


---

$$= \frac{P_\infty}{\sqrt{2n}} \sqrt{P_y (2|f|) + P_y \frac{1}{\delta}}$$

$$\leq P_\infty \sqrt{\frac{P_y (2|f|)}{2n}} + \frac{P_\infty}{\sqrt{2n}} \sqrt{P_y \frac{1}{\delta}}$$


---