

Apprentissage Statistique
M2 Probabilités et Statistiques, Université Paris-Sud

Cours 2 : Théorie de l'apprentissage statistique: de Vapnik à la localisation (2/2)

SYLVAIN ARLOT ET FRANCIS BACH
NOTES DE COURS INITIALEMENT PRISES PAR NICOLAS DROUGARD (2012)

TABLE DES MATIÈRES

2.4. Pénalisation	2
2.5. Majoration de la pénalité idéale "globale"	2
2.6. Bornes inférieures	6
3. Localisation et condition de marge	7
3.1. Condition de marge	7
3.2. Bornes supérieures sur le risque : vitesses rapides	8
3.3. Bornes inférieures sous condition de marge	8
3.4. Sélection de modèles sous condition de marge	9
Références	9

Rappels :

- Observations : $D_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$,
 X_i variable explicative et Y_i étiquette associée
- Prédicteur : application mesurable $t : \mathcal{X} \rightarrow \mathcal{Y}$
- \mathbb{S} l'ensemble des prédicteurs t
- Contraste : $\gamma : \begin{cases} \mathbb{S} \times (\mathcal{X} \times \mathcal{Y}) & \rightarrow & \mathbb{R} \\ t, (x, y) & \mapsto & \gamma(t, (x, y)) \end{cases}$
- Perte $P\gamma : \begin{cases} \mathbb{S} & \rightarrow & \mathbb{R} \\ t & \mapsto & \mathbb{E}_{(X,Y) \sim P} [\gamma(t, (X, Y))] \end{cases}$
- Prédicteur de Bayes :

$$s^* \in \operatorname{argmin}_{t \in \mathbb{S}} P\gamma(t)$$

- Perte relative $\ell(s^*, t) = P\gamma(t) - P\gamma(s^*)$
- Risque empirique $P_n\gamma(t) = \frac{1}{n} \sum_{i=1}^n \gamma(t, (x_i, y_i))$

2.4. Pénalisation. Nous rappelons que le minimiseur du risque empirique (ERM pour Empirical Risk Minimiseur) sur le modèle $S_m \subset \mathbb{S}$ est noté \hat{s}_m . La sélection de modèle nous amène à minimiser en $m \in \mathcal{M}$ (indices des modèles), un critère proche du risque de l'ERM. Une première idée serait de choisir le risque empirique, *i.e.* d'utiliser le critère $\text{crit}(m) = P_n \gamma(\hat{s}_m)$. Cependant, si $m_1 < m_2$, c'est-à-dire si le modèle S_{m_2} est plus complexe que S_{m_1} , le minimum du risque empirique sur S_{m_2} est plus petit que sur S_{m_1} : $P_n \gamma(\hat{s}_{m_1}) \geq P_n \gamma(\hat{s}_{m_2})$. Cela nous amènerait alors à choisir les modèles les plus complexes, et comme expliqué en section (2.2), mènerait à une explosion de l'erreur d'estimation : le problème du sur-apprentissage (overfitting).

La pénalisation consiste à contourner le problème en pénalisant les modèles les plus complexes : nous considérons les critères de la forme

$$\begin{aligned} \text{crit}(m) &= P_n \gamma(\hat{s}_m) + \text{pen}(m) \\ \Rightarrow \hat{m} &\in \arg \min_{m \in \mathcal{M}} \{ P_n \gamma(\hat{s}_m) + \text{pen}(m) \} \end{aligned}$$

$\text{pen}(m)$ corrige alors le biais de $P_n \gamma(\hat{s}_m)$ comme estimateur du risque $\mathbb{E}[P \gamma(\hat{s}_m)]$. La pénalité la plus adaptée serait la pénalité idéale

$$\text{pen}_{\text{id}}(m) = (P - P_n) \gamma(\hat{s}_m)$$

car nous minimiserions directement la perte de l'ERM.

Tout d'abord, nous pouvons remarquer l'inégalité suivante définissant la pénalité idéale globale :

$$\text{pen}_{\text{id}}(m) \leq \sup_{t \in S_m} \{ (P - P_n) \gamma(t) \} = \text{pen}_{\text{id,g}}(m)$$

Comme précisé dans la section précédente, nous cherchons un modèle satisfaisant une inégalité oracle. Si $\text{pen}(m) \geq \text{pen}_{\text{id}}(m)$ pour tout $m \in \mathcal{M}$, en utilisant le lemme "recette" avec $A(m) = \text{pen}(m) - \text{pen}_{\text{id}}(m)$ et $B(m) = 0$,

$$\text{pen}(m) - \text{pen}_{\text{id}}(m) \geq 0 \Leftrightarrow \text{crit}(m) - \text{crit}_{\text{id}}(m) \geq 0$$

$$(\text{Lemme}) \Rightarrow \ell(s^*, \hat{s}_{\hat{m}}) \leq \inf_{m \in \mathcal{M}} \{ \ell(s^*, \hat{s}_m) + \text{pen}(m) - \text{pen}_{\text{id}}(m) \}$$

où l'on rappelle que s^* est l'estimateur de Bayes.

Exemple : L'inégalité est vraie pour la pénalité idéale $\text{pen}(m) = \text{pen}_{\text{id,g}}(m)$ (la pénalité $\text{pen}_{\text{id,g}}(m)$ mesure la « capacité » de S_m).

2.5. Majoration de la pénalité idéale "globale". Référence : Section 8.2 de [6]. Nous souhaitons, comme précédemment, trouver une majoration de la pénalité idéale globale

$$\text{pen}_{\text{id,g}}(m) = \sup_{t \in S_m} \{ (P - P_n) \gamma(t) \}$$

afin de majorer l'erreur d'estimation (en espérance) et de construire une pénalité (voir propriété section 2.2).

Dans le cas où la fonction de contraste γ prend ses valeurs dans $[0,1]$ et S_m est fini, l'inégalité obtenue à l'exercice 3 du premier cours permet d'obtenir

une borne supérieure pour $\text{pen}_{\text{id,g}}(m)$: pour tout $x \geq 0$,

$$\text{pen}_{\text{id,g}}(m) < \sqrt{\frac{\ln(\text{Card}(S_m))}{2n}} + \sqrt{\frac{x}{2n}}$$

avec probabilité supérieure à $1 - \text{Card}(S_m) \exp(-x)$, si bien que

$$\mathbb{E} [\text{pen}_{\text{id,g}}(m)] \leq \sqrt{\frac{\ln(\text{Card}(S_m))}{2n}}$$

L'objectif de la démonstration ci-dessous est de généraliser ce résultat au cas où S_m est éventuellement infini, mais « se comporte essentiellement comme s'il était fini ».

- (1) Première étape : concentration. Le résultat suivant, qui peut être vu comme une généralisation de l'inégalité de Hoeffding, nous permet d'obtenir l'inégalité de concentration désirée :

Théorème 1. *Inégalité de Mc Diarmid* [6, Section 5.2]

Soient ξ_1, \dots, ξ_n des variables aléatoires indépendantes, $F : \mathbb{R}^n \rightarrow \mathbb{R}$ mesurable. Notons $Z = F(\xi_1, \dots, \xi_n)$.

$$\text{Si } \forall x, x' \in \mathbb{R}^n, \forall 1 \leq i \leq n,$$

$$|F(\xi_1, \dots, \xi_{i-1}, \xi'_i, \xi_{i+1}, \dots, \xi_n) - F(\xi)| \leq c_i$$

Alors $\forall x \leq 0$,

$$\mathbb{P}(Z - \mathbb{E}[Z] \geq x) \leq \exp\left(-\frac{2x^2}{\sum_{i=1}^n c_i^2}\right)$$

Utilisons cette inégalité avec $Z = \sup_{t \in \mathbb{S}_m} (P - P_n)\gamma(t) = F((X_1, Y_1), \dots, (X_n, Y_n))$.

Supposons que \hat{t}_i réalise $\sup_{t \in \mathbb{S}_m} (P - P_n^{(i)})\gamma(t) = F((X_1, Y_1), \dots, (X'_i, Y'_i), \dots, (X_n, Y_n))$:

$$\begin{aligned} \sup_{t \in \mathbb{S}_m} (P - P_n^{(i)})\gamma(t) - \sup_{t \in \mathbb{S}_m} (P - P_n)\gamma(t) &\leq (P - P_n^{(i)})\gamma(\hat{t}_i) - (P - P_n)\gamma(\hat{t}_i) \\ &= (P_n - P_n^{(i)})\gamma(\hat{t}_i) \\ &= \frac{1}{n}\gamma(\hat{t}_i, (X_i, Y_i)) - \frac{1}{n}\gamma(\hat{t}_i, (X_i, Y_i)) \\ &\leq \frac{1}{n} \end{aligned}$$

si l'on suppose encore que $0 \leq \gamma \leq 1$. Nous pouvons alors conclure avec l'inégalité de Mc Diarmid : si $\gamma(t; (x, y)) \in [0, 1]$,

$$\mathbb{P}\left(\text{pen}_{\text{id,g}}(m) \geq \mathbb{E}[\text{pen}_{\text{id,g}}(m)] + \sqrt{\frac{x}{2n}}\right) \leq e^{-x}$$

- (2) Deuxième étape : symétrisation. Nous calculons ici une majoration de l'espérance (sur l'échantillon) de la pénalité idéale globale. Nous posons $D_n = ((X_1, Y_1), \dots, (X_n, Y_n))$ iid \perp $D'_n = ((X'_1, Y'_1), \dots, (X'_n, Y'_n))$, et $\varepsilon_1, \dots, \varepsilon_n$ iid \perp D_n, D'_n , des variables de Rademacher, i.e. $\xi_1 \in$

$\{-1, 1\}$ et $\mathbb{P}(\varepsilon = 1) = \frac{1}{2}$.

$$\begin{aligned} \mathbb{E}[\text{pen}_{\text{id,g}}(m)] &= \mathbb{E}[\sup_{t \in \mathbb{S}_m} (P - P_n)\gamma(t)] \\ &= \mathbb{E}[\sup_{t \in \mathbb{S}_m} (\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D'_n} [\gamma(t, (X'_i, Y'_i)) - \gamma(t, (X_i, Y_i)) \mid D_n])] \\ &\leq \mathbb{E} \left[\sup_{t \in \mathbb{S}_m} (\frac{1}{n} \sum_{i=1}^n \gamma(t, (X'_i, Y'_i)) - \gamma(t, (X_i, Y_i))) \right] \end{aligned} \quad (1)$$

$$\begin{aligned} &= \mathbb{E} \left[\sup_{t \in \mathbb{S}_m} \left[\frac{1}{n} \sum_{i=1}^n \varepsilon_i (\gamma(t, (X'_i, Y'_i)) - \gamma(t, (X_i, Y_i))) \right] \right] \quad (2) \\ &\leq \mathbb{E} \left[\sup_{t \in \mathbb{S}_m} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \gamma(t, (X'_i, Y'_i)) \right] + \mathbb{E} \left[\sup_{t \in \mathbb{S}_m} \frac{1}{n} \sum_{i=1}^n -\varepsilon_i \gamma(t, (X_i, Y_i)) \right] \\ &= 2\mathbb{E} \left[\sup_{t \in \mathbb{S}_m} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \gamma(t, (X'_i, Y'_i)) \right] \end{aligned}$$

l'inégalité (1) est due à l'inégalité de Jensen version conditionnelle, car \sup_t est convexe en γ , et l'égalité (2), la "symétrisation", se comprend facilement en remarquant que les variables (X_i, Y_i) jouent le même rôle à cette étape.

Ainsi, si $\varepsilon_1, \dots, \varepsilon_n$ sont des variables aléatoires de Rademacher i.i.d.,

$$\mathbb{E}[\text{pen}_{\text{id,g}}(m)] \leq 2\mathbb{E} \left[\frac{1}{n} \sup_{t \in \mathbb{S}_m} \left\{ \sum_{i=1}^n \varepsilon_i \gamma(t; (X_i, Y_i)) \right\} \right]$$

(3) Troisième étape : trois options.

Supposons maintenant que γ est le contraste 0-1 : un Lemme et quelques définitions nous amènent à une majoration de l'espérance de l'étape précédente, conditionnellement à l'échantillon D_n . Notons pour tout $t \in \mathbb{S}$,

$$\begin{aligned} A_t &= \{x \in \mathcal{X} \text{ t.q. } t(x) = 1\} \subset \mathcal{X} \\ \mathcal{A}_m(D_n) &= \{(\gamma(t, (X_i, Y_i)))_{1 \leq i \leq n} \text{ t.q. } t \in S_m\} \\ &= \left\{ (\mathbf{1}_{t(X_i) \neq Y_i})_{1 \leq i \leq n} \text{ t.q. } t \in S_m \right\} \subset \{0, 1\}^n \end{aligned}$$

Alors,

$$\text{Card}(\mathcal{A}_m(D_n)) = \text{Card} \{A_t \cap \{X_1, \dots, X_n\} \text{ t.q. } t \in S_m\}$$

et

$$\mathbb{E} \left[\sup_{t \in \mathbb{S}_m} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \gamma(t, (X_i, Y_i)) \mid D_n \right] = \frac{1}{n} \mathbb{E} \left[\sup_{\alpha \in \mathcal{A}_m(D_n)} \sum_{i=1}^n \alpha_i \varepsilon_i \mid D_n \right]$$

Nous pouvons donc utiliser le résultat suivant relatif à la concentration de sommes de variables de Rademacher :

Lemme 2 (Section 6.1.1 de [6]). *Si $\mathcal{A} \subset \mathbb{R}^n$ est fini et $(\varepsilon_i)_{i=1,\dots,n}$ sont des variables aléatoires de Rademacher iid, alors*

$$\mathbb{E} \left[\sup_{\alpha \in \mathcal{A}} \sum_{i=1}^n \alpha_i \varepsilon_i \right] \leq \sqrt{2 \log(\text{Card}(\mathcal{A})) \sup_{\alpha \in \mathcal{A}} \sum_{i=1}^n \alpha_i^2}$$

Ce résultat est appliqué à $(\alpha_i)_{1 \leq i \leq n} = (\gamma(t, (X_i, Y_i)))_{1 \leq i \leq n} \in \mathcal{A}_m(D_n)$ pour D_n fixé. Nous remarquons que la somme des carrés des α_i est toujours inférieure à n (car nous considérons ici le contraste $0 - 1$), et donc

$$\begin{aligned} \mathbb{E} \left[\sup_{t \in \mathbb{S}_m} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \gamma(t, (X_i, Y_i)) \mid D_n \right] &= \frac{1}{n} \mathbb{E} \left[\sup_{\alpha \in \mathcal{A}_m(D_n)} \sum_{i=1}^n \alpha_i \varepsilon_i \mid D_n \right] \\ &\leq \sqrt{\frac{2 \log(\text{Card}(\mathcal{A}_m))}{n}} \end{aligned}$$

En définissant l'entropie combinatoire empirique H_m comme suit

$$H_m(X_1, \dots, X_n) = \ln(\text{Card}(\mathcal{A}_m(D_n)))$$

nous pouvons écrire la majoration combinatoire

$$\frac{2}{n} \mathbb{E} \left[\sup_{t \in \mathbb{S}_m} \left\{ \sum_{i=1}^n \varepsilon_i \gamma(t; (X_i, Y_i)) \right\} \mid D_n \right] \leq \frac{2\sqrt{2}}{\sqrt{n}} \sqrt{H_m(X_1, \dots, X_n)} \quad (3)$$

Ce sont maintenant que trois options s'offrent à nous :

- (a) première option : la majoration (3) est facilement exploitable pour le cas des classes de Vapnik-Chervonenkis.

Définition 1. Une *classe de Vapnik-Chervonenkis (VC)* est un modèle S_m tel que la quantité

$$V_m := \sup \left\{ k \in \mathbb{N} \text{ t.q. } \underbrace{\sup_{x_1, \dots, x_k \in \mathcal{X}} \exp(H_m(x_1, \dots, x_k))}_{= \sup_{D_n} \text{Card}(\mathcal{A}_m(D_n))} = 2^k \right\},$$

appelée *dimension*, est finie.

Par exemple, le modèle S_m associé à l'ensemble des demi-espaces de \mathbb{R}^d est une classe de Vapnik-Chervonenkis de dimension $d+1$. Le Lemme de Sauer, qui s'applique à ce type de classe, permet de majorer l'entropie combinatoire pour un échantillon assez grand :

Lemme 3. Si S_m est une classe VC, alors $\forall n \geq V_m$,

$$H_m(x_1, \dots, x_n) \leq V_m \ln\left(\frac{en}{V_m}\right)$$

Ainsi, en reprenant le résultat de l'étape 1 ainsi que (3), nous obtenons

$$\text{pen}_{\text{id,g}}(m) \leq \frac{2\sqrt{2}}{\sqrt{n}} \sqrt{V_m \left(1 + \ln \left(\frac{n}{V_m} \right) \right)} + \sqrt{\frac{2x_m}{n}} =: \text{pen}(m)$$

avec probabilité $1 - \exp(-x_m)$. Cette stratégie peut mener à une inégalité oracle avec forte probabilité (voir discussion section 2.4). C'est aussi un bon résultat, à m fixé, pour borner l'erreur d'estimation du modèle. (Pour en savoir plus, voir par exemple [3, Chapitre 13].)

- (b) deuxième option : la troisième étape peut nous inciter à utiliser l'entropie combinatoire H dans la pénalité, c'est-à-dire à prendre

$$\text{pen}(m) = C \sqrt{\frac{H_m(X_1, \dots, X_n)}{n}},$$

en utilisant le fait que $H_m(X_1, \dots, X_n)$ se concentre autour de son espérance [6, Section 5.3.2].

- (c) troisième option : Inspirés par la deuxième étape, nous pouvons aussi utiliser la complexité de Rademacher (globale) de S_m :

$$R_n(S_m; D_n) = \frac{1}{n} \mathbb{E} \left[\sup_{t \in S_m} \sum_{i=1}^n \varepsilon_i \gamma(t; (X_i, Y_i)) \mid D_n \right]$$

qui se concentre autour de son espérance, comme H_m [6, Section 5.3.2]. La pénalité correspondante sera alors deux fois cette complexité.

Ces stratégies ont un défaut commun : l'utilisation de $\text{pen}_{\text{id,g}}(m)$, qui est souvent d'un ordre de grandeur plus gros que $\text{pen}_{\text{id}}(m)$ (parce qu'il s'agit d'un supremum sur tout S_m , alors qu' \widehat{s}_m n'est pas « n'importe où » dans S_m), et qui a des déviations en $1/\sqrt{n}$.

L'idée de la localisation est de remettre en cause la première majoration $\text{pen}_{\text{id}}(m) \leq \text{pen}_{\text{id,g}}(m)$.

2.6. Bornes inférieures. Référence : [4] ou [3, Chapitre 14]. Soit \mathcal{P} une famille de lois sur $\mathcal{X} \times \mathcal{Y}$. Nous définissons dans cette section le risque minimax. Ce risque quantifie l'espérance de la perte relative "en pire cas sur $P \in \mathcal{P}$ ", loi de l'échantillon D_n , minimisée en l'estimateur construit à partir de l'échantillon $\tilde{s} : (\mathcal{X}, \mathcal{Y})^n \rightarrow \mathbb{S}$. Formellement,

Définition 2 (Risque minimax).

$$\mathcal{R}_{\min \max}(\mathcal{P}, n) := \inf_{\tilde{s} \text{ estimateur}} \sup_{P \in \mathcal{P}} \{ \mathbb{E}_{D_n \sim P^{\otimes n}} [\ell(s^*, \tilde{s}(D_n))] \}.$$

À titre d'exemple, le « No Free Lunch Theorem » (NFLT) exposé dans le premier cours, nous assure que lorsque \mathcal{X} est infini, $\mathcal{Y} = \{0, 1\}$ et $\gamma = \gamma_{0-1}$, ce risque minimax sur l'ensemble \mathcal{P} des lois de probabilité sur $\mathcal{X} \times \mathcal{Y}$ est supérieur ou égal à $\frac{1}{2}$.

Le théorème qui suit nous donne une minoration de ce risque lorsque l'on sait que s^* appartient à une classe VC connue (à comparer avec la majoration du risque de l'ERM sur S_m que l'on peut déduire de la Section 2.5).

Théorème 4 (Devroye & Lugosi, [4]). *Si S_m est une classe de Vapnik-Chervonenkis de dimension $V_m \geq 2$ et*

$$\mathcal{P}(S_m) = \{P \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y}) \text{ t.q. } s^* \in S_m\} ,$$

alors

$$\mathcal{R}_{\min \max}(\mathcal{P}(S_m), n) \geq \kappa \sqrt{\frac{V_m - 1}{n}} .$$

Aide pour la preuve : choisir x_1, \dots, x_n qui maximisent l'entropie $H_m(x_1, \dots, x_n)$, et prendre la loi uniforme sur ces valeurs pour X , et $\eta(x_i) = \frac{1}{2} + \varepsilon_i h$ pour Y avec $h \propto \frac{1}{\sqrt{n}}$, et $\varepsilon_1, \dots, \varepsilon_n$ v.a. de Rademacher iid.

3. LOCALISATION ET CONDITION DE MARGE

Référence : [2]

3.1. Condition de marge.

Définition 3 (condition de marge; Mammen et Tsybakov, 1999).

Pour $\beta \geq 0$,

$$\exists B > 0, \forall \varepsilon \geq 0, \quad \mathbb{P}(0 < |2\eta(X) - 1| \leq \varepsilon) \leq B\varepsilon^\beta \quad (M_\beta)$$

$$\exists h > 0, \quad \mathbb{P}(0 < |2\eta(X) - 1| < h) = 0 \quad (M_\infty)$$

Ces conditions se comprennent bien en pensant à la perte relative en classification binaire avec un contraste 0-1. En effet cette perte vaut $\ell(s^*, t) = \mathbb{E}[\mathbf{1}_{t(X) \neq s^*(X)} |2\eta(X) - 1|]$, ainsi la condition rend toute erreur de classification plus grave, et ainsi facilite ce problème lorsque $\beta > 0$.

Définition 4. pour $\alpha \in [0; 1]$,

$$\exists C > 0, \forall t \in \mathbb{S}, \quad \mathbb{E}_{(X,Y) \sim P} [(\gamma(t; (X, Y)) - \gamma(s^*; (X, Y)))^2] \leq C(\ell(s^*, t))^\alpha \quad (N_\alpha)$$

En notant $v = \text{var}(\gamma(t, \cdot) - \gamma(s^*, \cdot))$ et en utilisant l'inégalité de Bernstein, nous obtenons

$$(P - P_n)(\gamma(t) - \gamma(s^*)) \leq \sqrt{\frac{2x}{n}v} + \frac{x}{3n}$$

avec probabilité $1 - e^{-x}$. Sous la condition (N_α) nous obtenons la majoration

$$\sqrt{\frac{2x}{n}C\ell(s^*, t)^\alpha} + \frac{x}{3n} .$$

Exercice 1.

- (1) Pour tout $\beta \geq 0$, $(M_\beta) \Rightarrow (N_{\beta/(1+\beta)})$.
- (2) Pour tout $\alpha \in [0, 1[$, $(N_\alpha) \Rightarrow (M_{\alpha/(1-\alpha)})$.
- (3) $(M_\infty) \Leftrightarrow (N_1)$.

3.2. Bornes supérieures sur le risque : vitesses rapides. Référence : [7] ou [6, Section 8.3]. Pour un argument informel, voir [2, Section 5.3].

Dans cette section, l'inégalité de concentration clef s'appelle l'inégalité de Talagrand. [6, Section 5.3.4].

Théorème 5 (Talagrand, puis Bousquet).

Soit $b > 0$ et \mathcal{F} un ensemble de fonctions $\Xi \mapsto \mathbb{R}$ telles que $\forall f \in \mathcal{F}$,

$P(f) - f \leq b$, où $P(f) = \mathbb{E}_{\xi \sim P} [f(\xi)]$.

Soit $Z = \sup_{f \in \mathcal{F}} \{(P - P_n)(f)\}$.

Alors, pour tout $x \geq 0$, avec probabilité au moins $1 - e^{-x}$,

$$Z \leq \mathbb{E}[Z] + \sqrt{\frac{2(\sup_{f \in \mathcal{F}} \{\text{var}_P(f)\} + 2b\mathbb{E}[Z])x}{n}} + \frac{bx}{3n} .$$

Remarque : nous retrouvons l'inégalité de Bernstein lorsque $\text{Card}(\mathcal{F}) = 1$. L'idée est d'appliquer l'inégalité de Talagrand à $\Xi = \mathcal{X} \times \mathcal{Y}$ et \mathcal{F} un sous-ensemble de

$$\mathcal{F}_m = \{\gamma(t; \cdot) - \gamma(s^*; \cdot) \text{ t.q. } t \in S_m\} ,$$

par exemple,

$$\mathcal{F} = \mathcal{F}_m(r) = \{f \in \mathcal{F}_m, \text{var}_P(f) \leq r\} .$$

Principe de la localisation : si (M_∞) est satisfaite et si $\ell(s^*, \hat{s}_m) \leq r_0$, alors l'inégalité de Talagrand appliquée à $\mathcal{F}_m(r_0/h)$ donne une nouvelle borne $r_1 < r_0$ sur $\ell(s^*, \hat{s}_m)$ avec grande probabilité. En itérant le processus, on peut obtenir des bornes meilleures que $n^{-1/2}$ pour $\ell(s^*, \hat{s}_m(D_n))$.

Théorème 6 (Massart et Nédélec, [7]). Si S_m est une classe de Vapnik de dimension $V_m \geq 1$, $s^* \in S_m$ et si (M_∞) est vérifiée avec $h \geq \sqrt{V/n}$, alors tout minimiseur \hat{s}_m du risque empirique sur S_m vérifie l'inégalité

$$\mathbb{E}[\ell(s^*, \hat{s}_m)] \leq C \frac{V}{nh} \left[1 + \ln \left(\frac{nh^2}{V} \right) \right]$$

où C est une constante numérique.

3.3. Bornes inférieures sous condition de marge. Référence : [7].

Théorème 7 (Massart et Nédélec, [7]). Si S_m est une classe de Vapnik de dimension $V_m \geq 2$ et

$$\mathcal{P}(S_m, h) = \{P \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y}) \text{ t.q. } s^* \in S_m \text{ et } (M_\infty) \text{ est vérifiée}\} ,$$

alors

$$\mathcal{R}_{\min \max}(\mathcal{P}(S_m), n) \geq \kappa \min \left\{ \frac{V_m}{nh}, \sqrt{\frac{V_m}{n}} \right\} .$$

\Rightarrow Adaptation du minimiseur du risque empirique \hat{s}_m au paramètre de marge h (au facteur $\ln(n)$ près, parfois nécessaire).

3.4. **Sélection de modèles sous condition de marge.** Références pour en savoir plus : [5] et [1].

- Pour avoir une inégalité-oracle tenant compte de vitesses rapides, on doit prendre $\text{pen}(m)$ plus proche de $\text{pen}_{\text{id}}(m)$ que $\text{pen}_{\text{id,g}}(m)$.
- Il existe des pénalités « locales » (par exemple, les complexités de Rademacher locales), reposant sur une majoration de $(P - P_n)(\gamma(\hat{s}_m) - \gamma(s^*))$ utilisant le principe de localisation (en particulier l'inégalité de Talagrand).
- Défauts de telles pénalités.

RÉFÉRENCES

- [1] Peter L. Bartlett, Olivier Bousquet, and Shahar Mendelson. Local Rademacher complexities. *Ann. Statist.*, 33(4) :1497–1537, 2005.
- [2] Stéphane Boucheron, Olivier Bousquet, and Gábor Lugosi. Theory of classification : a survey of some recent advances. *ESAIM Probab. Stat.*, 9 :323–375 (electronic), 2005.
- [3] Luc P. Devroye, László Györfi, and Gábor Lugosi. *A probabilistic theory of pattern recognition*, volume 31 of *Applications of Mathematics (New York)*. Springer-Verlag, New York, 1996.
- [4] Luc P. Devroye and Gábor Lugosi. Lower bounds in pattern recognition. *Pattern recognition*, 28 :1011–1018, 1995.
- [5] Vladimir Koltchinskii. Local Rademacher complexities and oracle inequalities in risk minimization. *Ann. Statist.*, 34(6) :2593–2656, 2006.
- [6] Pascal Massart. *Concentration Inequalities and Model Selection*, volume 1896 of *Lecture Notes in Mathematics*. Springer, Berlin, 2007. Lectures from the 33rd Summer School on Probability Theory held in Saint-Flour, July 6–23, 2003, With a foreword by Jean Picard.
- [7] Pascal Massart and Élodie Nédélec. Risk bounds for statistical learning. *Ann. Statist.*, 34(5) :2326–2366, 2006.

URL: <http://www.di.ens.fr/~arlot/>