



Journées mathématiques X-UPS

Année 2013

Aléatoire

Christophe GIRAUD

Fondements mathématiques de l'apprentissage statistique

Journées mathématiques X-UPS (2013), p. 59-92.

<https://doi.org/10.5802/xups.2013-02>

© Les auteurs, 2013.



Cet article est mis à disposition selon les termes de la licence

LICENCE INTERNATIONALE D'ATTRIBUTION CREATIVE COMMONS BY 4.0.

<https://creativecommons.org/licenses/by/4.0/>

Les Éditions de l'École polytechnique
Route de Saclay
F-91128 PALAISEAU CEDEX
<https://www.editions.polytechnique.fr>

Centre de mathématiques Laurent Schwartz
CMLS, École polytechnique, CNRS,
Institut polytechnique de Paris
F-91128 PALAISEAU CEDEX
<https://portail.polytechnique.edu/cmls/>



Publication membre du

Centre Mersenne pour l'édition scientifique ouverte

www.centre-mersenne.org

FONDEMENTS MATHÉMATIQUES DE L'APPRENTISSAGE STATISTIQUE

par

Christophe Giraud

Résumé. L'objectif d'un algorithme de classification est de prédire au mieux la classe d'un objet à partir d'observations de cet objet. Un exemple typique est le filtre à spam des messageries électroniques qui prédisent (plus ou moins bien) si un courriel est un spam ou non. Nous introduisons dans ces notes les principaux concepts fondamentaux de la théorie de la classification statistique supervisée et quelques uns des algorithmes de classification les plus populaires. Nous soulignons chemin faisant l'importance de certains concepts mathématiques, parmi lesquels la symétrisation, la convexification, les inégalités de concentration, le principe de contraction et les espaces de Hilbert à noyau reproduisant.

Table des matières

1. Introduction.....	60
2. Modélisation mathématique.....	61
3. Minimisation du risque empirique.....	62
3.1. Probabilité de mauvaise classification avec $\hat{h}_{\mathcal{H}}$	63
3.2. Sélection de dictionnaire.....	68
3.3. Dimension combinatoire de Vapnik-Chervonenkis .	70
4. De la théorie vers la pratique.....	73
4.1. Convexification du risque empirique.....	73
4.2. Propriétés statistiques.....	77
4.3. Support Vector Machine.....	81
4.4. AdaBoost.....	86
5. Pour aller au-delà de cette introduction.....	87
Appendice A. Espace de Hilbert à noyaux reproduisants (RKHS).....	87
Appendice B. Inégalités probabilistes.....	91
Références.....	92

1. Introduction

En ce début de siècle, nous observons un accroissement phénoménal de l'utilisation des mathématiques, à la fois dans l'industrie et dans les laboratoires scientifiques. Cet essor de l'importance des mathématiques va de pair avec l'explosion des volumes de données collectées et l'accroissement de la puissance de calculs des ordinateurs. Dans l'industrie, la modélisation mathématique apparaît à tous les stades de la vie d'un produit. Depuis la conception technique, avec force de simulations numériques, via la production, avec l'optimisation des ressources et des flux, jusqu'au marketing et la distribution avec des prévisions basées sur l'analyse de grandes bases de données. Dans les laboratoires scientifiques, la modélisation mathématique devient de plus en plus cruciale, en particulier en biologie et médecine où les scientifiques doivent extraire des informations pertinentes des données massives qu'ils produisent grâce aux récents développements biotechnologiques.

La classification automatique est peut-être l'un des usages quotidiens les plus invasifs des mathématiques. L'objectif de la classification automatique est de prédire au mieux la classe y d'un objet x à partir d'observations de ce dernier. Un exemple typique est le filtre à spam de notre messagerie électronique qui prédit (plus ou moins bien) si un courriel est un spam ou non. La classification automatique est omniprésente dans notre quotidien, filtrant nos courriels, lisant automatiquement les codes postaux sur nos lettres ou reconnaissant les visages sur les photos que nous publions sur les réseaux sociaux. Elle est aussi extrêmement importante en sciences, par exemple en médecine pour effectuer des diagnostics précoces à partir de données à hauts débits, ou pour la recherche *in silico* de médicaments efficaces.

Nous introduisons dans ces notes les principaux concepts fondamentaux de l'apprentissage statistique (supervisé). Nous décrivons dans la partie 2 la modélisation mathématique d'un problème générique de classification. Dans la partie 3, nous analysons la précision prédictive d'un algorithme universel de classification et dans la partie 4 nous dérivons de cet algorithme théorique des algorithmes numériquement implémentables et populaires. Les appendices rassemblent

des résultats techniques nécessaires pour la définition et l'analyse des algorithmes de classification.

2. Modélisation mathématique

Par soucis de simplicité, nous allons nous restreindre au cas où il y a seulement deux classes (comme pour le filtre à spam). Le problème de la classification automatique peut alors être modélisé de la manière suivante. Soit \mathcal{X} un espace mesuré. On observe conjointement un point $X \in \mathcal{X}$ et une étiquette $Y \in \{-1, +1\}$. Notre objectif est de construire une fonction $h : \mathcal{X} \rightarrow \{-1, +1\}$, appelée *classifieur*, telle que $h(X)$ prédit au mieux l'étiquette Y .

Supposons que le couple $(X, Y) \in \mathcal{X} \times \{-1, +1\}$ est issu d'un tirage selon une loi \mathbb{P} . Pour un classifieur $h : \mathcal{X} \rightarrow \{-1, +1\}$ donné, la probabilité de mauvaise classification est

$$L(h) = \mathbb{P}(Y \neq h(X)).$$

Comme $|Y - h(X)| \in \{0, 2\}$, nous avons par le théorème de Pythagore

$$\begin{aligned} L(h) &= \frac{1}{4} \mathbb{E}[(Y - h(X))^2] \\ &= \frac{1}{4} \mathbb{E}[(Y - \mathbb{E}[Y|X])^2] + \frac{1}{4} \mathbb{E}[(\mathbb{E}[Y|X] - h(X))^2]. \end{aligned}$$

En conséquence, $L(h)$ est minimal pour le classifieur de Bayes

$$h_*(X) = \text{sign}(\mathbb{E}[Y|X]) \quad \text{où} \quad \text{sign}(x) = \mathbf{1}_{x>0} - \mathbf{1}_{x\leq 0} \quad \text{pour } x \in \mathbb{R}.$$

Lorsque la loi \mathbb{P} est connue, il suffit donc d'utiliser le classifieur de Bayes h_* pour avoir la plus petite probabilité possible de mauvaise classification.

Malheureusement, la loi \mathbb{P} est en général inconnue et on ne peut donc pas calculer le classifieur de Bayes h_* . En pratique, nous avons uniquement accès à des données dites d'apprentissage $(X_i, Y_i)_{i=1, \dots, n}$ i.i.d. de loi \mathbb{P} et notre objectif est de construire à partir de ces données d'apprentissage un classifieur $\hat{h} : \mathcal{X} \rightarrow \{-1, +1\}$ tel que $L(\hat{h}) - L(h_*)$ est aussi petit que possible.

3. Minimisation du risque empirique

Comme la loi \mathbb{P} est inconnue, nous ne pouvons pas calculer la probabilité de mauvaise classification $L(h)$. On peut cependant calculer à la place la probabilité empirique de mauvaise classification

$$\widehat{L}_n(h) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{Y_i \neq h(X_i)} = \widehat{\mathbb{P}}_n(Y \neq h(X)),$$

où $\widehat{\mathbb{P}}_n = \frac{1}{n} \sum_{i=1}^n \delta_{(X_i, Y_i)}$. À un ensemble \mathcal{H} de classifieurs, appelé *dictionnaire*, on peut associer le classifieur de minimisation du risque empirique

$$(3.1) \quad \widehat{h}_{\mathcal{H}} \in \operatorname{argmin}_{h \in \mathcal{H}} \widehat{L}_n(h).$$

La définition de ce classifieur est très naturelle, cependant deux questions fondamentales se posent immédiatement : quel dictionnaire \mathcal{H} choisir ? et comment $\widehat{h}_{\mathcal{H}}$ se comporte-t-il comparativement à h_* ? Ces deux questions sont bien évidemment fortement reliées. En décomposant la différence entre les probabilités de mauvaise classification $L(\widehat{h}_{\mathcal{H}})$ et $L(h_*)$, on obtient

$$0 \leq L(\widehat{h}_{\mathcal{H}}) - L(h_*) = \underbrace{\min_{h \in \mathcal{H}} L(h) - L(h_*)}_{\text{erreur d'approximation}} + \underbrace{L(\widehat{h}_{\mathcal{H}}) - \min_{h \in \mathcal{H}} L(h)}_{\text{erreur stochastique}}.$$

Le premier terme mesure la qualité de l'approximation de h_* par un classifieur $h \in \mathcal{H}$. Cette erreur d'approximation est purement déterministe et élargir le dictionnaire \mathcal{H} ne peut que la réduire. Le second terme mesure l'erreur résultant de la minimisation sur $h \in \mathcal{H}$ de la probabilité empirique de mauvaise classification $\widehat{L}_n(h)$ au lieu de la vraie probabilité de mauvaise classification $L(h)$. Ce terme est stochastique et il tend à augmenter lorsque \mathcal{H} grossit. Ce phénomène est illustré figure 1. Dans cette illustration dans $\mathcal{X} = \mathbb{R}^2$, les classifieurs du dictionnaire $\mathcal{H}_{\text{lin}} = \{h(x) = \text{sign}(\langle w, x \rangle) : \|w\| = 1\}$ ne sont pas suffisamment flexibles et ils produisent une classification de mauvaise qualité. Dans ce cas, l'erreur d'approximation est grande. À l'opposé, les classifieurs du dictionnaire

$$\mathcal{H}_{\text{poly}} = \{h(x) = 2 \mathbf{1}_A(x) - 1 : A \text{ polygone dans } \mathcal{X}\}$$

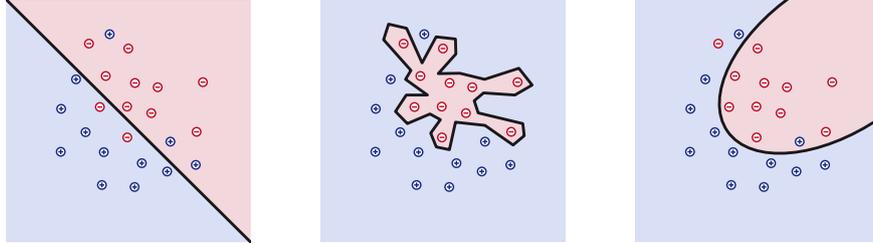


FIGURE 1. Exemples de classifications produites par différents dictionnaires. Gauche : avec des classifieurs linéaires \mathcal{H}_{lin} . Centre : avec des classifieurs polygonaux $\mathcal{H}_{\text{poly}}$. Droite : avec des classifieurs basés sur des formes quadratiques.

sont très flexibles et ils peuvent toujours reproduire exactement la classification des données $(X_i, Y_i)_{i=1, \dots, n}$. L'erreur empirique $\widehat{L}_n(\widehat{h}_{\mathcal{H}_{\text{poly}}})$ est donc toujours nulle, mais $\widehat{h}_{\mathcal{H}_{\text{poly}}}$ tend à mal classifier une nouvelle donnée (X, Y) et l'erreur stochastique $L(\widehat{h}_{\mathcal{H}}) - \min_{h \in \mathcal{H}} L(h)$ est grande. Le dernier exemple, basé sur un ensemble moins flexible de classifieurs quadratiques produit de meilleurs résultats.

Pour choisir un bon dictionnaire \mathcal{H} , il faut donc trouver un bon équilibre entre les propriétés d'approximation de \mathcal{H} et sa taille. La première étape pour développer un principe de sélection de dictionnaire \mathcal{H} est de quantifier la probabilité de mauvaise classification du minimiseur du risque empirique $\widehat{h}_{\mathcal{H}}$.

3.1. Probabilité de mauvaise classification avec $\widehat{h}_{\mathcal{H}}$

Comme mentionné ci-dessus, augmenter la taille de \mathcal{H} tend à augmenter l'erreur stochastique $L(\widehat{h}_{\mathcal{H}}) - \min_{h \in \mathcal{H}} L(h)$. En fait, ce n'est pas exactement la taille du dictionnaire qui compte, mais plutôt sa flexibilité en terme de classification. Par exemple, nous ne pouvons pas classer correctement les trois points étiquetés

$$\{((0, 1), +1), ((1, 1), -1), ((1, 0), +1)\}$$

avec les classifieurs dans \mathcal{H}_{lin} . A l'inverse, pour tout ensemble $(x_i, y_i)_{i=1, \dots, n}$ de points étiquetés, il existe $h \in \mathcal{H}_{\text{poly}}$ tel que $h(x_i) = y_i$.

Nous quantifions cette flexibilité de classification par le coefficient d'éclatement

$$(3.2) \quad \mathfrak{S}_n(\mathcal{H}) = \max_{(x_1, \dots, x_n) \in \mathcal{X}^n} \text{card}\{(h(x_1), \dots, h(x_n)) : h \in \mathcal{H}\},$$

qui donne le nombre maximum d'étiquetages différents de n points qui peuvent être produits par les classifieurs dans \mathcal{H} . Par exemple, comme on peut obtenir à partir de $\mathcal{H}_{\text{poly}}$ tous les étiquetages possibles de n points (distincts), nous avons $\mathbb{S}_n(\mathcal{H}_{\text{poly}}) = 2^n$. À l'inverse, le nombre d'étiquetages possibles de n points avec les classifieurs dans \mathcal{H}_{lin} est plus limité. En effet, la proposition 1 partie 3.3 implique que $\mathbb{S}_n(\mathcal{H}_{\text{lin}}) \leq (n+1)^2$. Le prochain théorème nous fournit une borne supérieure de l'erreur stochastique et un intervalle de confiance pour la probabilité de mauvaise classification $L(\hat{h}_{\mathcal{H}})$ en terme du coefficient d'éclatement.

Théorème 1 (Contrôle de l'erreur stochastique). *Pour tout $t > 0$, nous avons avec probabilité au moins $1 - e^{-t}$*

$$(3.3) \quad L(\hat{h}_{\mathcal{H}}) - \min_{h \in \mathcal{H}} L(h) \leq 4 \sqrt{\frac{2 \log(2 \mathbb{S}_{\mathcal{H}}(n))}{n}} + \sqrt{\frac{2t}{n}}$$

et

$$(3.4) \quad |L(\hat{h}_{\mathcal{H}}) - \hat{L}_n(\hat{h}_{\mathcal{H}})| \leq 2 \sqrt{\frac{2 \log(2 \mathbb{S}_{\mathcal{H}}(n))}{n}} + \sqrt{\frac{t}{2n}}$$

Démonstration. Nous segmentons la preuve en trois lemmes. Le premier lemme montre que les termes de gauche dans (3.3) et (3.4) peuvent être bornés à l'aide du maximum sur \mathcal{H} de la différence entre l'erreur empirique de mauvaise classification et la vraie erreur de mauvaise classification

$$(3.5) \quad \hat{\Delta}_n(\mathcal{H}) = \sup_{h \in \mathcal{H}} |\hat{L}_n(h) - L(h)|.$$

Lemme 1.1. *Pour tout dictionnaire \mathcal{H} et entier n , nous avons les bornes supérieures*

$$L(\hat{h}_{\mathcal{H}}) - \min_{h \in \mathcal{H}} L(h) \leq 2 \hat{\Delta}_n(\mathcal{H}) \quad \text{et} \quad |L(\hat{h}_{\mathcal{H}}) - \hat{L}_n(\hat{h}_{\mathcal{H}})| \leq \hat{\Delta}_n(\mathcal{H}).$$

Démonstration. Pour tout $h \in \mathcal{H}$, nous avons $\hat{L}_n(\hat{h}_{\mathcal{H}}) \leq \hat{L}_n(h)$ et donc

$$\begin{aligned} L(\hat{h}_{\mathcal{H}}) - L(h) &= L(\hat{h}_{\mathcal{H}}) - \hat{L}_n(\hat{h}_{\mathcal{H}}) + \hat{L}_n(\hat{h}_{\mathcal{H}}) - L(h) \\ &\leq L(\hat{h}_{\mathcal{H}}) - \hat{L}_n(\hat{h}_{\mathcal{H}}) + \hat{L}_n(h) - L(h) \\ &\leq 2 \hat{\Delta}_n(\mathcal{H}). \end{aligned}$$

Comme cette inégalité est vraie pour tout $h \in \mathcal{H}$, nous en déduisons directement la première borne du lemme 1.1. La seconde borne est évidente. \square

Pour démontrer le théorème 1, il reste à prouver qu'on a

$$\widehat{\Delta}_n(\mathcal{H}) \leq 2\sqrt{\frac{2\log(2\mathbb{S}_{\mathcal{H}}(n))}{n}} + \sqrt{\frac{t}{2n}}$$

avec probabilité au moins $1 - e^{-t}$. Nous segmentons la preuve de cette borne en deux lemmes.

Lemme 1.2. *Avec probabilité $1 - e^{-t}$, nous avons*

$$\widehat{\Delta}_n(\mathcal{H}) \leq \mathbb{E}[\widehat{\Delta}_n(\mathcal{H})] + \sqrt{\frac{t}{2n}}.$$

Démonstration. Nous avons $\widehat{\Delta}_n(\mathcal{H}) = F((X_1, Y_1), \dots, (X_n, Y_n))$ avec

$$F : (\mathcal{X} \times \{-1, +1\})^n \longrightarrow \mathbb{R}$$

$$((x_1, y_1), \dots, (x_n, y_n)) \longmapsto \frac{1}{n} \sup_{h \in \mathcal{H}} \left| \sum_{i=1}^n \mathbf{1}_{y_i \neq h(x_i)} - L(h) \right|.$$

Pour tout $(x_1, y_1), \dots, (x_n, y_n), (x'_i, y'_i) \in \mathcal{X} \times \{-1, +1\}$, nous avons

$$\begin{aligned} & |F((x_1, y_1), \dots, (x'_i, y'_i), \dots, (x_n, y_n)) \\ & \quad - F((x_1, y_1), \dots, (x_i, y_i), \dots, (x_n, y_n))| \leq \frac{1}{n}. \end{aligned}$$

En conséquence, l'inégalité de concentration de McDiarmid (voir théorème 4 de l'appendice B), nous assure qu'avec probabilité au moins $1 - e^{-2ns^2}$, nous avons $\widehat{\Delta}_n(\mathcal{H}) \leq \mathbb{E}[\widehat{\Delta}_n(\mathcal{H})] + s$. Le lemme 1.2 en découle (avec le changement de variable $s = \sqrt{t/(2n)}$). \square

Il reste à borner l'espérance de $\widehat{\Delta}_n(\mathcal{H})$ à l'aide de $\mathbb{S}_{\mathcal{H}}(n)$.

Lemme 1.3. *Pour tout dictionnaire \mathcal{H} , nous avons l'inégalité*

$$\mathbb{E}[\widehat{\Delta}_n(\mathcal{H})] \leq 2\sqrt{\frac{2\log(2\mathbb{S}_{\mathcal{H}}(n))}{n}}.$$

Démonstration. La preuve du lemme 1.3 est basée sur un argument classique et élégant de symétrisation.

La première étape de la symétrisation est de représenter la probabilité de mauvaise classification $L(h)$ comme l'espérance de la probabilité empirique de mauvaise classification

$$L(h) = \mathbb{P}(Y \neq h(X)) = \tilde{\mathbb{E}} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\tilde{Y}_i \neq h(\tilde{X}_i)} \right],$$

où $(\tilde{X}_i, \tilde{Y}_i)_{i=1, \dots, n}$ est indépendant de $(X_i, Y_i)_{i=1, \dots, n}$ et est identiquement distribué. Dans la suite, nous noterons $\tilde{\mathbb{E}}$ pour l'espérance par rapport aux variables $(\tilde{X}_i, \tilde{Y}_i)_{i=1, \dots, n}$ et \mathbb{E} pour l'espérance par rapport aux variables $(X_i, Y_i)_{i=1, \dots, n}$. Les inégalités de Jensen et Fatou nous donnent

$$\begin{aligned} \mathbb{E}[\widehat{\Delta}_n(\mathcal{H})] &= \mathbb{E} \left[\sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{Y_i \neq h(X_i)} - \tilde{\mathbb{E}} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\tilde{Y}_i \neq h(\tilde{X}_i)} \right] \right| \right] \\ &\leq \mathbb{E} \tilde{\mathbb{E}} \left[\sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{1}_{Y_i \neq h(X_i)} - \mathbf{1}_{\tilde{Y}_i \neq h(\tilde{X}_i)}) \right| \right]. \end{aligned}$$

La seconde étape consiste à capitaliser sur la symétrie des variables $\mathbf{1}_{Y_i \neq h(X_i)} - \mathbf{1}_{\tilde{Y}_i \neq h(\tilde{X}_i)}$. Nous introduisons n variables aléatoires i.i.d. $(\sigma_i)_{i=1, \dots, n}$ indépendantes de $(X_i, Y_i, \tilde{X}_i, \tilde{Y}_i)_{i=1, \dots, n}$ et de loi

$$\mathbb{P}_\sigma(\sigma_i = 1) = \mathbb{P}_\sigma(\sigma_i = -1) = 1/2.$$

Par symétrie, on remarque que $(\sigma_i (\mathbf{1}_{Y_i \neq h(X_i)} - \mathbf{1}_{\tilde{Y}_i \neq h(\tilde{X}_i)}))_{i=1, \dots, n}$ a la même loi que $(\mathbf{1}_{Y_i \neq h(X_i)} - \mathbf{1}_{\tilde{Y}_i \neq h(\tilde{X}_i)})_{i=1, \dots, n}$, donc nous obtenons

$$\begin{aligned} \mathbb{E}[\widehat{\Delta}_n(\mathcal{H})] &\leq \mathbb{E} \tilde{\mathbb{E}} \mathbb{E}_\sigma \left[\sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i (\mathbf{1}_{Y_i \neq h(X_i)} - \mathbf{1}_{\tilde{Y}_i \neq h(\tilde{X}_i)}) \right| \right] \\ &\leq 2 \mathbb{E} \mathbb{E}_\sigma \left[\sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbf{1}_{Y_i \neq h(X_i)} \right| \right] \\ &\leq 2 \max_{y \in \{-1, +1\}^n} \max_{x \in \mathcal{X}^n} \mathbb{E}_\sigma \left[\sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbf{1}_{y_i \neq h(x_i)} \right| \right], \end{aligned}$$

où la seconde inégalité découle de l'inégalité triangulaire. Pour tout $(x, y) \in \mathcal{X}^n \times \{-1, +1\}^n$, nous introduisons l'ensemble

$$\mathcal{V}_{\mathcal{H}}(x, y) = \{(\mathbf{1}_{y_1 \neq h(x_1)}, \dots, \mathbf{1}_{y_n \neq h(x_n)}) : h \in \mathcal{H}\}.$$

La dernière majoration de $\mathbb{E}[\widehat{\Delta}_n(\mathcal{H})]$ peut s'écrire comme

$$\mathbb{E}[\widehat{\Delta}_n(\mathcal{H})] \leq \frac{2}{n} \times \max_{y \in \{-1, +1\}^n} \max_{x \in \mathcal{X}^n} \mathbb{E}_{\sigma} \left[\sup_{v \in \mathcal{V}_{\mathcal{H}}(x, y)} |\langle \sigma, v \rangle| \right],$$

où $\langle x, y \rangle$ est le produit scalaire canonique sur \mathbb{R}^n . Remarquons au passage qu'entre le membre de gauche et le membre de droite nous avons remplacé l'espérance sous la probabilité inconnue \mathbb{P} par une espérance sous la probabilité connue \mathbb{P}_{σ} . Pour tout $y \in \{-1, +1\}^n$ remarquons aussi qu'il y a une bijection entre $\mathcal{V}_{\mathcal{H}}(x, y)$ et l'ensemble $\{(h(x_1), \dots, h(x_n)) : h \in \mathcal{H}\}$. En conséquence, nous avons la majoration

$$\max_{y \in \{-1, +1\}^n} \max_{x \in \mathcal{X}^n} \text{card}(\mathcal{V}_{\mathcal{H}}(x, y)) \leq \mathbb{S}_n(\mathcal{H}).$$

Au regard des deux dernières inégalités, il suffit de démontrer que

$$(3.6) \quad \mathbb{E}_{\sigma}[\sup_{v \in \mathcal{V}} |\langle \sigma, v \rangle|] \leq \sqrt{2n \log(2 \text{card}(\mathcal{V}))},$$

pour tout ensemble $\mathcal{V} \subset \{-1, 0, +1\}^n$

pour conclure la preuve du lemme 1.3. Démontrons (3.6). Avec la notation $\mathcal{V}^{\#} = \mathcal{V} \cup -\mathcal{V}$, l'inégalité de Jensen nous donne pour tout $s > 0$

$$(3.7) \quad \begin{aligned} \mathbb{E}_{\sigma} \left[\sup_{v \in \mathcal{V}} |\langle \sigma, v \rangle| \right] &= \mathbb{E}_{\sigma} \left[\sup_{v \in \mathcal{V}^{\#}} \langle \sigma, v \rangle \right] \leq \frac{1}{s} \log \mathbb{E}_{\sigma} \left[\sup_{v \in \mathcal{V}^{\#}} e^{s \langle \sigma, v \rangle} \right] \\ &\leq \frac{1}{s} \log \left(\sum_{v \in \mathcal{V}^{\#}} \mathbb{E}_{\sigma} [e^{s \langle \sigma, v \rangle}] \right). \end{aligned}$$

En combinant le fait que les σ_i sont indépendants, que $(e^x + e^{-x}) \leq 2e^{x^2/2}$ pour tout $x \in \mathbb{R}$ et que $v_i^2 \in \{0, 1\}$ pour tout $v \in \mathcal{V}^{\#}$, nous obtenons

$$\begin{aligned} \mathbb{E}_{\sigma} [e^{s \langle \sigma, v \rangle}] &= \prod_{i=1}^n \mathbb{E}_{\sigma} [e^{s v_i \sigma_i}] = \prod_{i=1}^n \frac{1}{2} (e^{s v_i} + e^{-s v_i}) \\ &\leq \prod_{i=1}^n e^{s^2 v_i^2 / 2} \leq e^{n s^2 / 2}. \end{aligned}$$

En injectant cette majoration dans (3.7), nous obtenons

$$\mathbb{E}_\sigma[\sup_{v \in \mathcal{V}} |\langle \sigma, v \rangle|] \leq \frac{\log(\text{card}(\mathcal{V}^\#))}{s} + \frac{ns}{2} \quad \text{pour tout } s > 0.$$

Le terme de droite est minimal pour $s = \sqrt{2 \log(\text{card}(\mathcal{V}^\#))}/n$, ce qui nous donne la majoration

$$\mathbb{E}_\sigma[\sup_{v \in \mathcal{V}} |\langle \sigma, v \rangle|] \leq \sqrt{2n \log(\text{card}(\mathcal{V}^\#))}.$$

Nous obtenons finalement la borne (3.6) en remarquant que $\text{card}(\mathcal{V}^\#) \leq 2 \text{card}(\mathcal{V})$. La preuve du lemme 1.3 est complète. \square

Les bornes (3.3) et (3.4) s'obtiennent directement en combinant les trois lemmes. \square

3.2. Sélection de dictionnaire

Soit $\{\mathcal{H}_1, \dots, \mathcal{H}_M\}$ une collection de dictionnaires de classifieurs. Nous aimerions sélectionner parmi cette collection, le dictionnaire \mathcal{H}_* avec la plus petite probabilité de mauvaise classification $L(\hat{h}_{\mathcal{H}_*})$. Ce dictionnaire \mathcal{H}_* , dit *dictionnaire oracle*, dépend de la probabilité inconnue \mathbb{P} , il n'est donc malheureusement pas connu du statisticien. Dans cette partie, nous allons développer à partir du théorème 1 une procédure reposant sur les données pour sélectionner parmi la collection $\{\mathcal{H}_1, \dots, \mathcal{H}_M\}$ un dictionnaire $\mathcal{H}_{\hat{m}}$ ayant des performances similaires à celles de \mathcal{H}_* .

Le dictionnaire oracle \mathcal{H}_* est obtenu en minimisant la probabilité de mauvaise classification $L(\hat{h}_{\mathcal{H}})$ sur $\mathcal{H} \in \{\mathcal{H}_1, \dots, \mathcal{H}_M\}$. Une première idée est de sélectionner $\mathcal{H}_{\hat{m}}$ en minimisant sur la collection $\{\mathcal{H}_1, \dots, \mathcal{H}_M\}$ la probabilité empirique de mauvaise classification $\hat{L}_n(\hat{h}_{\mathcal{H}})$. Cette procédure de sélection ne donne pas de bons résultats car pour tout $\mathcal{H} \subset \mathcal{H}'$ on a toujours $\hat{L}_n(\hat{h}_{\mathcal{H}'}) \leq \hat{L}_n(\hat{h}_{\mathcal{H}})$, donc la procédure tend à sélectionner le dictionnaire le plus gros possible. Pour bâtir une bonne procédure de sélection, nous devons prendre en compte les fluctuations de $\hat{L}_n(\hat{h}_{\mathcal{H}})$ autour de $L(\hat{h}_{\mathcal{H}})$. La borne (3.4) du théorème 1 nous offre un contrôle de ces fluctuations. En élaborant à partir de cette borne, nous obtenons le résultat suivant.

Théorème 2 (Sélection de dictionnaire). *Considérons la procédure de sélection de dictionnaire*

$$\hat{m} = \operatorname{argmin}_{m=1,\dots,M} \{ \widehat{L}_n(\widehat{h}_{\mathcal{H}_m}) + \operatorname{pen}(\mathcal{H}_m) \},$$

$$\text{avec } \operatorname{pen}(\mathcal{H}) = 2 \sqrt{\frac{2 \log(2 \mathbb{S}_n(\mathcal{H}))}{n}}.$$

Alors, pour tout $t > 0$, avec probabilité au moins $1 - e^{-t}$ nous avons

$$(3.8) \quad L(\widehat{h}_{\mathcal{H}_{\hat{m}}}) \leq \min_{m=1,\dots,M} \{ \inf_{h \in \mathcal{H}_m} L(h) + 2 \operatorname{pen}(\mathcal{H}_m) \} + \sqrt{\frac{2 \log(M) + 2t}{n}}.$$

Avant de prouver le théorème 2, commentons la borne (3.8). Comme $\min_{h \in \mathcal{H}} L(h) \leq L(\widehat{h}_{\mathcal{H}})$, nous obtenons avec probabilité $1 - e^{-t}$ que

$$L(\widehat{h}_{\mathcal{H}_{\hat{m}}}) \leq L(\widehat{h}_{\mathcal{H}_*}) + 2 \operatorname{pen}(\mathcal{H}_*) + \sqrt{\frac{2 \log(M) + 2t}{n}}.$$

En particulier, cela nous permet de comparer la probabilité de mauvaise classification du classifieur sélectionné à celle du meilleur classifieur parmi la collection $\{\widehat{h}_{\mathcal{H}_1}, \dots, \widehat{h}_{\mathcal{H}_M}\}$.

Remarquons aussi que le second terme de la borne (3.8) augmente en $\sqrt{2 \log(M)/n}$ avec le nombre M de dictionnaires candidats. Enfin, nous noterons que les résultats restent valides si on prend $\operatorname{pen}(\mathcal{H})$ plus grand que $2 \sqrt{2 \log(2 \mathbb{S}_n(\mathcal{H}))/n}$.

Démonstration du théorème 2

Rappelons la notation $\widehat{\Delta}_n(\mathcal{H}) = \sup_{h \in \mathcal{H}} |\widehat{L}_n(h) - L(h)|$. Les lemmes 1.2 et 1.3 nous assurent qu'avec probabilité au moins $1 - e^{-t}$ nous avons

$$(3.9) \quad \widehat{\Delta}_n(\mathcal{H}_m) \leq \operatorname{pen}(\mathcal{H}_m) + \sqrt{\frac{\log(M) + t}{2n}},$$

simultanément pour tout $m = 1, \dots, M$.

En conséquence, le lemme 1.1 et la définition du critère de sélection nous donnent qu'avec probabilité au moins $1 - e^{-t}$

$$(3.10) \quad L(\widehat{h}_{\mathcal{H}_{\hat{m}}}) \leq \widehat{L}_n(\widehat{h}_{\mathcal{H}_{\hat{m}}}) + \operatorname{pen}(\mathcal{H}_{\hat{m}}) + \sqrt{\frac{\log(M) + t}{2n}}$$

$$\leq \min_{m=1,\dots,M} \{ \widehat{L}_n(\widehat{h}_{\mathcal{H}_m}) + \operatorname{pen}(\mathcal{H}_m) \} + \sqrt{\frac{\log(M) + t}{2n}}.$$

Pour conclure, il suffit de contrôler $\widehat{L}_n(\widehat{h}_{\mathcal{H}_m})$ en fonction de $\inf_{h \in \mathcal{H}_m} L(h)$. Cela peut être fait directement en combinant les inégalités (3.3) et (3.4), mais la majoration résultante peut être améliorée par le raisonnement ci-dessous.

Pour comparer $\widehat{L}_n(\widehat{h}_{\mathcal{H}_m})$ à $\inf_{h \in \mathcal{H}_m} L(h)$, commençons par remarquer que pour tout $h \in \mathcal{H}_m$ nous avons

$$\widehat{L}_n(\widehat{h}_{\mathcal{H}_m}) \leq \widehat{L}_n(h) \leq L(h) + \widehat{\Delta}_n(\mathcal{H}_m),$$

donc en prenant l'infimum sur $h \in \mathcal{H}_m$ nous obtenons pour tout $m = 1, \dots, M$

$$\widehat{L}_n(\widehat{h}_{\mathcal{H}_m}) \leq \inf_{h \in \mathcal{H}_m} L(h) + \widehat{\Delta}_n(\mathcal{H}_m).$$

En combinant cette majoration avec (3.9) et (3.10), on trouve

$$L(\widehat{h}_{\mathcal{H}_{\widehat{m}}}) \leq \min_{m=1, \dots, M} \left\{ \inf_{h \in \mathcal{H}_m} L(h) + 2 \operatorname{pen}(\mathcal{H}_m) \right\} + 2 \sqrt{\frac{\log(M) + t}{2n}}.$$

La preuve du théorème 2 est complète. \square

Remarque. En combinant (3.9) avec le lemme 1.1, on obtient un intervalle de confiance pour la probabilité de mauvaise classification

$$\mathbb{P}(L(\widehat{h}_{\mathcal{H}_{\widehat{m}}}) \in [\widehat{L}_n(\widehat{h}_{\mathcal{H}_{\widehat{m}}}) - \delta(\widehat{m}, t), \widehat{L}_n(\widehat{h}_{\mathcal{H}_{\widehat{m}}}) + \delta(\widehat{m}, t)]) \geq 1 - e^{-t}$$

avec $\delta(\widehat{m}, t) = \operatorname{pen}(\mathcal{H}_{\widehat{m}}) + \sqrt{\frac{\log(M) + t}{2n}}.$

3.3. Dimension combinatoire de Vapnik-Chervonenkis

Le calcul des coefficients d'éclatement $\mathbb{S}_n(\mathcal{H})$ peut être complexe en pratique. Cependant, une propriété combinatoire remarquable de ces coefficients d'éclatement permet d'obtenir une majoration simple de $\mathbb{S}_n(\mathcal{H})$, ne dépendant de \mathcal{H} qu'à travers une unique quantité, la dimension de Vapnik-Chervonenkis de \mathcal{H} définie ci-dessous.

Par convention on posera $\mathbb{S}_0(\mathcal{H}) = 1$. On appelle VC-dimension de \mathcal{H} l'entier $d_{\mathcal{H}}$ défini par

$$d_{\mathcal{H}} = \sup\{d \in \mathbb{N} : \mathbb{S}_d(\mathcal{H}) = 2^d\} \in \mathbb{N} \cup \{+\infty\}.$$

Elle correspond au nombre maximum de points de \mathcal{X} qui peuvent être classifiés de toutes les façons possibles par les classifieurs de \mathcal{H} . La proposition suivante donne une majoration du coefficient d'éclatement $\mathbb{S}_n(\mathcal{H})$ en fonction de la VC-dimension $d_{\mathcal{H}}$.

Proposition 1 (Lemme de Sauer). *Soit \mathcal{H} un ensemble de classifieurs de VC-dimension $d_{\mathcal{H}}$ finie. Pour tout $n \in \mathbb{N}$, on a*

$$\mathbb{S}_n(\mathcal{H}) \leq \sum_{i=0}^{d_{\mathcal{H}}} C_n^i \leq (n+1)^{d_{\mathcal{H}}} \quad \text{avec} \quad C_n^i = \begin{cases} \frac{n!}{i!(n-i)!} & \text{pour } n \geq i \\ 0 & \text{pour } n < i. \end{cases}$$

Démonstration. Nous commençons par prouver par récurrence sur k l'inégalité

$$(3.11) \quad \mathbb{S}_k(\mathcal{H}) \leq \sum_{i=0}^{d_{\mathcal{H}}} C_k^i$$

pour tout \mathcal{H} de VC-dimension finie $d_{\mathcal{H}}$.

Considérons dans un premier temps le cas $k = 1$. Si $d_{\mathcal{H}} = 0$, aucun point ne peut être étiqueté de deux manières différentes, donc tous les points ne peuvent être étiquetés que d'une seule manière. Il en découle que $\mathbb{S}_1(\mathcal{H}) = 1$, qui est égal à C_1^0 . Si $d_{\mathcal{H}} \geq 1$, nous avons $\mathbb{S}_1(\mathcal{H}) = 2$ qui est égal à $C_1^0 + C_1^1$.

Supposons maintenant que (3.11) est vraie pour tout $k \leq n - 1$. Considérons \mathcal{H} de VC-dimension $d_{\mathcal{H}}$ finie. Comme on vient de le voir, lorsque $d_{\mathcal{H}} = 0$ tous les points peuvent être étiquetés d'une seule manière donc $\mathbb{S}_k(\mathcal{H}) = 1$ et (3.11) est vraie pour tout k . Considérons maintenant le cas $d_{\mathcal{H}} \geq 1$. Soient $x_1, \dots, x_n \in \mathcal{X}$ et définissons

$$\mathcal{H}(x_1, \dots, x_n) = \{(h(x_1), \dots, h(x_n)) : h \in \mathcal{H}\}.$$

L'ensemble $\mathcal{H}(x_1, \dots, x_n)$ ne dépend que des valeurs de h sur $\{x_1, \dots, x_n\}$, donc on peut remplacer \mathcal{H} par $\mathcal{F} = \{h|_{\{x_1, \dots, x_n\}} : h \in \mathcal{H}\}$ dans la définition de $\mathcal{H}(x_1, \dots, x_n)$. Comme $d_{\mathcal{F}}$ est inférieur à $d_{\mathcal{H}}$, on peut supposer sans perte de généralité que $\mathcal{X} = \{x_1, \dots, x_n\}$ et $\mathcal{H} = \mathcal{F}$, ce que nous faisons par la suite. Considérons l'ensemble

$$\mathcal{H}' = \{h \in \mathcal{H} : h(x_n) = 1 \text{ et } h' = h - 2 \times \mathbf{1}_{\{x_n\}} \in \mathcal{H}\}.$$

Comme $\mathcal{H}(x_1, \dots, x_n) = \mathcal{H}'(x_1, \dots, x_n) \cup (\mathcal{H} \setminus \mathcal{H}')(x_1, \dots, x_n)$ nous avons

$$(3.12) \quad \text{card}(\mathcal{H}(x_1, \dots, x_n)) \leq \text{card}(\mathcal{H}'(x_1, \dots, x_n)) + \text{card}((\mathcal{H} \setminus \mathcal{H}')(x_1, \dots, x_n)).$$

Nous allons majorer séparément le cardinal de $\mathcal{H}'(x_1, \dots, x_n)$ et celui de $(\mathcal{H} \setminus \mathcal{H}')(x_1, \dots, x_n)$.

(1) Pour commencer, remarquons que

$$\text{card}(\mathcal{H}'(x_1, \dots, x_n)) = \text{card}(\mathcal{H}'(x_1, \dots, x_{n-1}))$$

car on a $h(x_n) = 1$ pour tout $h \in \mathcal{H}'$. Ensuite, on observe que la VC-dimension $d_{\mathcal{H}'}$ de \mathcal{H}' est au plus $d_{\mathcal{H}} - 1$. En effet, si d points x_{i_1}, \dots, x_{i_d} de $\mathcal{X} = \{x_1, \dots, x_n\}$ peuvent être étiquetés de toutes les manières possibles par \mathcal{H}' , alors $x_n \notin \{x_{i_1}, \dots, x_{i_d}\}$ vu que $h(x_n) = 1$ pour tout $h \in \mathcal{H}'$. De plus, l'ensemble $\{x_{i_1}, \dots, x_{i_d}, x_n\}$ peut être étiqueté de toutes les manières possibles par \mathcal{H} par définition de \mathcal{H}' , donc $d + 1 \leq d_{\mathcal{H}}$, ce qui implique $d_{\mathcal{H}'} \leq d_{\mathcal{H}} - 1$. En appliquant (3.11) avec $k = n - 1$ on obtient

$$(3.13) \quad \begin{aligned} \text{card}(\mathcal{H}'(x_1, \dots, x_n)) &= \text{card}(\mathcal{H}'(x_1, \dots, x_{n-1})) \\ &\leq \sum_{i=0}^{d_{\mathcal{H}'}-1} C_{n-1}^i. \end{aligned}$$

(2) Si $h, h' \in \mathcal{H} \setminus \mathcal{H}'$ vérifient $h(x_i) = h'(x_i)$ pour $i = 1, \dots, n - 1$, alors ils vérifient aussi $h(x_n) = h'(x_n)$ par définition de \mathcal{H}' . En conséquence, on a comme précédemment l'égalité

$$\text{card}((\mathcal{H} \setminus \mathcal{H}')(x_1, \dots, x_n)) = \text{card}((\mathcal{H} \setminus \mathcal{H}')(x_1, \dots, x_{n-1})).$$

De plus $d_{\mathcal{H} \setminus \mathcal{H}'}$ est inférieur à $d_{\mathcal{H}}$ car $\mathcal{H} \setminus \mathcal{H}' \subset \mathcal{H}$ donc l'équation (3.11) avec $k = n - 1$ donne

$$(3.14) \quad \begin{aligned} \text{card}((\mathcal{H} \setminus \mathcal{H}')(x_1, \dots, x_n)) \\ = \text{card}((\mathcal{H} \setminus \mathcal{H}')(x_1, \dots, x_{n-1})) &\leq \sum_{i=0}^{d_{\mathcal{H}}-1} C_{n-1}^i. \end{aligned}$$

En combinant (3.12), (3.13) et (3.14), on obtient

$$\text{card}(\mathcal{H}(x_1, \dots, x_n)) \leq \sum_{i=1}^{d_{\mathcal{H}}} C_{n-1}^{i-1} + \sum_{i=0}^{d_{\mathcal{H}}} C_{n-1}^i = \sum_{i=0}^{d_{\mathcal{H}}} C_n^i,$$

car $C_{n-1}^i + C_{n-1}^{i-1} = C_n^i$ pour $i \geq 1$. Il en découle que (3.11) est vrai pour $k = n$ et la récurrence est établie.

La seconde borne de la proposition est obtenue par

$$\sum_{i=0}^d C_n^i \leq \sum_{i=0}^d \frac{n^i}{i!} \leq \sum_{i=0}^d C_d^i n^i = (1+n)^d.$$

La preuve de la proposition 1 est complète. \square

Donnons quelques exemples de VC-dimension pour quelques dictionnaires simples de $\mathcal{X} = \mathbb{R}^d$. Les preuves sont laissées en exercice.

Exemple 1 : classifieurs linéaires. La VC-dimension de l'ensemble

$$\mathcal{H} = \{h(x) = \text{sign}(\langle w, x \rangle) : \|w\| = 1\}$$

de classifieurs linéaires est d .

Exemple 2 : classifieurs affines. La VC-dimension de l'ensemble

$$\mathcal{H} = \{h(x) = \text{sign}(\langle w, x \rangle + b) : \|w\| = 1, b \in \mathbb{R}\}$$

de classifieurs affines est $d + 1$.

Exemple 3 : classifieurs hyper-rectangulaires. La VC-dimension de l'ensemble

$$\mathcal{H} = \{h(x) = 2\mathbf{1}_A(x) - 1 : A \text{ hyper-rectangle de } \mathbb{R}^d\}$$

de classifieurs hyper-rectangulaires est $2d$.

Exemple 4 : classifieurs à polygones convexes. La VC-dimension de l'ensemble

$$\mathcal{H} = \{h(x) = 2\mathbf{1}_A(x) - 1 : A \text{ polygone convexe de } \mathbb{R}^d\}$$

de classifieurs par polygones convexes est $+\infty$ (prendre n points sur la sphère unité : pour tout sous-ensemble de ces points, on peut prendre leur enveloppe convexe comme polygone convexe).

4. De la théorie vers la pratique

4.1. Convexification du risque empirique

Les classifieurs par minimisation du risque empirique analysés précédemment possèdent de bonnes qualités statistiques, mais ils ne peuvent pas être mis en œuvre en pratique à cause de leur coût en

temps de calcul. En effet, il n'existe pas de façon efficace de minimiser (3.1) car ni \mathcal{H} ni \widehat{L}_n ne sont convexes. Un certain nombre des algorithmes de classification les plus populaires sont obtenus par une simple relaxation convexe du problème de minimisation (3.1). La probabilité empirique de mauvaise classification \widehat{L}_n est remplacée par un substitut convexe et l'ensemble des classifieurs \mathcal{H} est remplacé par un ensemble fonctionnel convexe $\mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$.

Considérons un ensemble convexe \mathcal{F} de fonctions de \mathcal{X} dans \mathbb{R} . Une fonction $f \in \mathcal{F}$ n'est pas un classifieur, mais on peut l'utiliser pour classifier en étiquetant les points en fonction du signe de f . En d'autres mots, on peut associer à f le classifieur $\text{sign}(f)$. La probabilité empirique de mauvaise classification de ce classifieur peut être réécrite sous la forme

$$\widehat{L}_n(\text{sign}(f)) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{Y_i \text{sign}(f)(X_i) < 0\}} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{Y_i f(X_i) < 0\}}.$$

Remplaçons maintenant la probabilité empirique de mauvaise classification \widehat{L}_n par un substitut convexe se prêtant bien au calcul numérique. Un moyen simple et efficace d'obtenir un critère convexe est de remplacer la fonction de perte $z \rightarrow \mathbf{1}_{z < 0}$ par une fonction convexe $z \rightarrow \ell(z)$. Nous allons développer cette idée et considérer dans la suite les classifieurs obtenus par la procédure

$$(4.1) \quad \widehat{h}_{\mathcal{F}} = \text{sign}(\widehat{f}_{\mathcal{F}})$$

$$\text{où } \widehat{f}_{\mathcal{F}} = \underset{f \in \mathcal{F}}{\text{argmin}} \widehat{L}_n^{\ell}(f) \quad \text{avec} \quad \widehat{L}_n^{\ell}(f) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i f(X_i)).$$

Ce classifieur peut être calculé numériquement de façon efficace car à la fois \mathcal{F} et \widehat{L}_n^{ℓ} sont convexes. De nombreux classifieurs populaires sont obtenus en résolvant (4.1) pour des choix spécifiques de \mathcal{F} et ℓ , voir par exemple les parties 4.3 et 4.4 pour quelques exemples.

Quelques fonctions de perte convexes ℓ classiques. Il est naturel de considérer une fonction de perte ℓ qui est décroissante et positive. Habituellement, on demande aussi que $\ell(z) \geq \mathbf{1}_{z < 0}$ pour tout $z \in \mathbb{R}$ car cela nous permet de donner une majoration de la probabilité

de mauvaise classification, voir le théorème 3. Quelques fonctions de pertes classiques sont

- la perte exponentielle $\ell(z) = e^{-z}$
- la perte logit $\ell(z) = \log_2(1 + e^{-z})$
- la perte hinge $\ell(z) = (1 - z)_+$ (avec $(x)_+ = \max(0, x)$)

voir la figure 2 pour un tracé de ces trois fonctions.

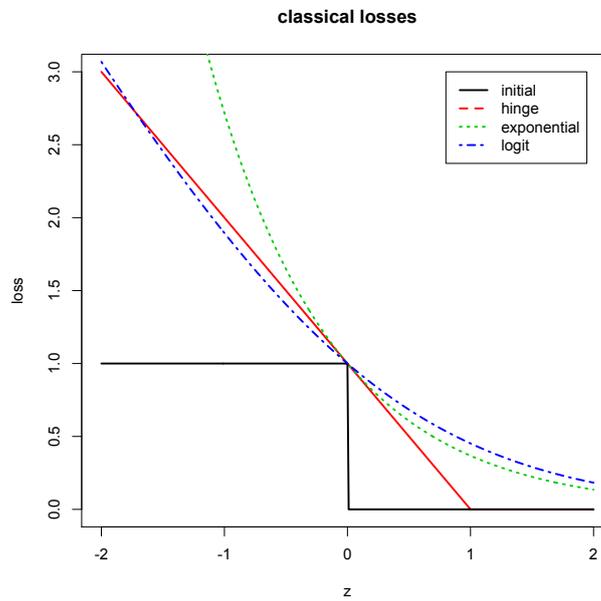


FIGURE 2. Tracé des pertes exponentielle, hinge et logit

Quelques exemples classiques d'ensembles fonctionnels \mathcal{F} . Les principaux exemples classiques d'ensembles fonctionnels \mathcal{F} peuvent être regroupés en deux classes.

Une première classe d'ensembles \mathcal{F} est obtenue en prenant des combinaisons linéaires d'une famille finie $\mathcal{H} = \{h_1, \dots, h_p\}$ de classifieurs

$$(4.2) \quad \mathcal{F} = \left\{ f : f(x) = \sum_{j=1}^p \beta_j h_j(x) \text{ avec } \beta_j \in \mathcal{C} \right\},$$

où \mathcal{C} est un sous-ensemble convexe de \mathbb{R}^p . Des choix typiques pour \mathcal{C} sont

- le simplexe $\{\beta \in \mathbb{R}^p : \beta_j \geq 0, \sum_{j=1}^p \beta_j \leq 1\}$,
- la boule $\{\beta \in \mathbb{R}^p : |\beta|_1 \leq R\}$
- ou tout l'ensemble \mathbb{R}^p .

Ces choix apparaissent notamment dans les méthodes de « boosting », voir partie 4.4. Les classifieurs de base $\{h_1, \dots, h_p\}$ sont souvent appelés *weak learners*. Un choix populaire pour les weak learners est $h_j(x) = \text{sign}(x_j - t_j)$ avec $t_j \in \mathbb{R}$.

Une seconde classe classique d'ensembles \mathcal{F} est obtenue en prenant une boule d'une espace de Hilbert à noyaux reproduisant (qu'on appelle RKHS, pour *Reproducing Kernel Hilbert Space* en anglais). On renvoie à l'appendice A pour une brève introduction aux RKHS. Plus précisément, considérons un RKHS \mathcal{F}_k de noyau k et notons par $\|f\|_{\mathcal{F}}$ la norme hilbertienne de $f \in \mathcal{F}_k$. Pour alléger les notations, dans la suite on écrira simplement \mathcal{F} pour \mathcal{F}_k . Minimiser \tilde{L}_n^ℓ sur la boule $\{f \in \mathcal{F} : \|f\|_{\mathcal{F}} \leq R\}$ est équivalent à minimiser sur \mathcal{F} le problème lagrangien dual

$$(4.3) \quad \hat{f}_{\mathcal{F}} = \underset{f \in \mathcal{F}}{\text{argmin}} \tilde{L}_n^\ell(f) \quad \text{avec} \quad \tilde{L}_n^\ell(f) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i f(X_i)) + \lambda \|f\|_{\mathcal{F}}^2,$$

pour un certain $\lambda > 0$. Ce type de classifieur apparaît notamment dans les algorithmes de type *Support Vector Machine*, présentés partie 4.3. Il s'avère que le minimiseur $\hat{f}_{\mathcal{F}}$ de (4.3) est de la forme

$$(4.4) \quad \hat{f}_{\mathcal{F}} = \sum_{i=1}^n \hat{\beta}_i k(X_i, \cdot).$$

En effet, notons V l'espace linéaire engendré par $k(X_1, \cdot), \dots, k(X_n, \cdot)$ et décomposons $f = f_V + f_{V^\perp}$ sur $V \oplus V^\perp$. Par la propriété de reproduction on a $f(X_i) = \langle f, k(X_i, \cdot) \rangle_{\mathcal{F}} = \langle f_V, k(X_i, \cdot) \rangle_{\mathcal{F}} = f_V(X_i)$, donc la formule de Pythagore nous donne

$$\tilde{L}_n^\ell(f_V + f_{V^\perp}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i f_V(X_i)) + \lambda \|f_V\|_{\mathcal{F}}^2 + \lambda \|f_{V^\perp}\|_{\mathcal{F}}^2.$$

Comme λ est strictement positif, tout minimiseur de \widehat{f} de \widetilde{L}_n^ℓ doit satisfaire $\widehat{f}_{V^\perp} = 0$, donc il est de la forme (4.4). De plus, la propriété de reproduction nous assure de nouveau que $\langle k(X_i, \cdot), k(X_j, \cdot) \rangle_{\mathcal{F}} = k(X_i, X_j)$ donc

$$\left\| \sum_{j=1}^n \beta_j k(X_j, \cdot) \right\|_{\mathcal{F}}^2 = \sum_{i,j=1}^n \beta_i \beta_j k(X_i, X_j).$$

Le problème de minimisation (4.3) est donc équivalent à

$$(4.5) \quad \widehat{f}_{\mathcal{F}} = \sum_{j=1}^n \widehat{\beta}_j k(X_j, \cdot) \quad \text{avec}$$

$$\widehat{\beta} = \operatorname{argmin}_{\beta \in \mathbb{R}^n} \left\{ \frac{1}{n} \sum_{i=1}^n \ell \left(\sum_{j=1}^n \beta_j Y_i k(X_j, X_i) \right) + \lambda \sum_{i,j=1}^n \beta_i \beta_j k(X_i, X_j) \right\}.$$

Cette formulation est cruciale en pratique, car elle réduit le problème de minimisation infini-dimensionnel (4.3) à un problème de minimisation n -dimensionnel qui peut être résolu efficacement. Dans la partie 4.3 sur les Support Vector Machines, nous donnerons une description plus précise des solutions de ce problème lorsque ℓ est la perte Hinge.

4.2. Propriétés statistiques

Le classifieur $\widehat{h}_{\mathcal{F}}$ donné par (4.1) avec \mathcal{F} et ℓ convexes a l'avantage de se calculer facilement, mais a-t-il de bonnes propriétés statistiques ?

Lien avec le classifieur de Bayes. Le minimiseur du risque empirique $\widehat{h}_{\mathcal{H}}$ de la partie 3 était un minimiseur de la version empirique de la probabilité de mauvaise classification $\mathbb{P}(Y \neq h(X))$ sur un ensemble \mathcal{H} de classifieurs. La fonction $\widehat{f}_{\mathcal{F}}$ minimise à la place la version empirique de $\mathbb{E}[\ell(Yf(X))]$ sur un ensemble fonctionnel \mathcal{F} . Le classifieur $\widehat{h}_{\mathcal{H}}$ peut donc être vu comme une version empirique du classifieur de Bayes h_* qui minimise $\mathbb{P}(Y \neq h(X))$ sur l'ensemble des fonctions mesurables $h : \mathcal{X} \rightarrow \{-1, +1\}$, alors que la fonction $\widehat{f}_{\mathcal{F}}$ est une version empirique de la fonction f_*^ℓ qui minimise $\mathbb{E}[\ell(Yf(X))]$ sur l'ensemble des fonctions mesurables $f : \mathcal{X} \rightarrow \mathbb{R}$. Une première

étape est de comprendre le lien entre le classifieur de Bayes h_* et le signe de la fonction f_*^ℓ . Il s'avère que sous des hypothèses très faibles sur ℓ , le signe de f_*^ℓ coïncide exactement avec le classifieur de Bayes h_* , donc $\text{sign}(f_*^\ell)$ minimise la probabilité de mauvaise classification $\mathbb{P}(Y \neq h(X))$. Démontrons maintenant ce résultat.

En conditionnant par X on obtient

$$\begin{aligned} \mathbb{E}[\ell(Yf(X))] &= \mathbb{E}[\mathbb{E}[\ell(Yf(X)|X)]] \\ &= \mathbb{E}\left[\ell(f(X))\mathbb{P}(Y=1|X) + \ell(-f(X))(1 - \mathbb{P}(Y=1|X))\right]. \end{aligned}$$

Supposons que ℓ est décroissante, dérivable et strictement convexe (par exemple perte exponentielle ou logit). En minimisant l'expression ci-dessus on obtient que $f_*^\ell(X)$ est solution de

$$\frac{\ell'(-f(X))}{\ell'(f(X))} = \frac{\mathbb{P}(Y=1|X)}{1 - \mathbb{P}(Y=1|X)}.$$

Comme ℓ est strictement convexe, on a $f(X) > 0$ si et seulement si $\ell'(-f(X))/\ell'(f(X)) > 1$, donc

$$\begin{aligned} f_*^\ell(X) > 0 &\iff \mathbb{P}(Y=1|X) > 1/2 \\ &\iff \mathbb{E}[Y|X] = 2\mathbb{P}(Y=1|X) - 1 > 0. \end{aligned}$$

Comme $h_*(X) = \text{sign}(\mathbb{E}[Y|X])$ (voir partie 2), on obtient $\text{sign}(f_*^\ell) = h_*$. Cette égalité reste vraie pour la perte Hinge (exercice!).

Pour résumer la discussion ci-dessus, la fonction cible f_*^ℓ approchée par $\widehat{f}_{\mathcal{F}}$ est intéressante pour le problème de classification car son signe coïncide avec le meilleur classifieur possible h_* .

Majoration de la probabilité de mauvaise classification. Nous nous concentrons maintenant sur la probabilité de mauvaise classification $L(\widehat{h}_{\mathcal{F}})$ du classifieur $\widehat{h}_{\mathcal{F}} = \text{sign}(\widehat{f}_{\mathcal{F}})$ défini par (4.1). En pratique, il est intéressant d'avoir une majoration de la probabilité de mauvaise classification $L(\widehat{h}_{\mathcal{F}})$ qui puisse être évaluée à partir des données. Le prochain théorème fournit une telle majoration pour des exemples typiques d'ensembles \mathcal{F} .

Théorème 3 (Borne de confiance pour $L(\widehat{h}_{\mathcal{F}})$). *Pour tout $R > 0$, on définit $\Delta\ell(R) = |\ell(R) - \ell(-R)|$. On suppose que la fonction de perte ℓ est convexe, décroissante, positive, α -Lipschitz sur $[-R, R]$ et vérifie*

$\ell(z) \geq \mathbf{1}_{z < 0}$ pour tout z dans \mathbb{R} . On s'intéresse au classifieur $\widehat{h}_{\mathcal{F}}$ défini par (4.1).

(a) Lorsque \mathcal{F} est de la forme (4.2) avec $\mathcal{C} = \{\beta \in \mathbb{R}^p : |\beta|_1 \leq R\}$, nous avons avec probabilité au moins $1 - e^{-t}$

$$(4.6) \quad L(\widehat{h}_{\mathcal{F}}) \leq \widehat{L}_n^{\ell}(\widehat{f}_{\mathcal{F}}) + 2\alpha R \sqrt{\frac{2 \log(2p)}{n}} + \Delta \ell(R) \sqrt{\frac{t}{2n}}.$$

(b) Soit \mathcal{F} une boule de rayon R d'un RKHS de noyau k vérifiant $k(x, x) \leq 1$ pour tout $x \in \mathcal{X}$. Alors, avec probabilité au moins $1 - e^{-t}$, on a

$$(4.7) \quad L(\widehat{h}_{\mathcal{F}}) \leq \widehat{L}_n^{\ell}(\widehat{f}_{\mathcal{F}}) + \frac{2\alpha R}{\sqrt{n}} + \Delta \ell(R) \sqrt{\frac{t}{2n}}.$$

Démonstration. Commençons par prouver une majoration générale de $L(\widehat{h}_{\mathcal{F}})$ très similaire au théorème 1.

Lemme 3.1. Supposons que $\sup_{f \in \mathcal{F}} |f(x)| \leq R < +\infty$. Alors pour toute fonction de perte ℓ vérifiant les hypothèses du théorème 3, on a avec probabilité au moins $1 - e^{-t}$

$$(4.8) \quad L(\widehat{h}_{\mathcal{F}}) \leq \widehat{L}_n^{\ell}(\widehat{f}_{\mathcal{F}}) + \frac{2\alpha}{n} \max_{x \in \mathcal{X}^n} \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i f(x_i) \right| \right] + \Delta \ell(R) \sqrt{\frac{t}{2n}}$$

où $\sigma_1, \dots, \sigma_n$ sont des variables aléatoires i.i.d. de loi $\mathbb{P}_{\sigma}(\sigma_i = 1) = \mathbb{P}_{\sigma}(\sigma_i = -1) = 1/2$.

Démonstration du lemme 3.1. La preuve de ce lemme repose sur exactement les mêmes arguments que la preuve du théorème 1. La première étape est de remarquer que l'inégalité $\ell(z) \geq \mathbf{1}_{z < 0}$ pour tout réel z nous donne

$$\begin{aligned} L(\widehat{h}_{\mathcal{F}}) &= \mathbb{P}(Y \widehat{f}_{\mathcal{F}}(X) < 0) \leq L^{\ell}(f) \quad \text{avec } L^{\ell}(f) = \mathbb{E}[\ell(Yf(X))] \\ &\leq \widehat{L}_n^{\ell}(f) + \widehat{\Delta}_n^{\ell}(\mathcal{F}) \\ &\quad \text{où } \widehat{\Delta}_n^{\ell}(\mathcal{F}) = \sup_{f \in \mathcal{F}} |\widehat{L}_n^{\ell}(f) - L^{\ell}(f)|. \end{aligned}$$

Comme dans le lemme 1.2, l'inégalité de concentration de McDiarmid (théorème 4 de l'appendice B) assure qu'avec probabilité au moins

$1 - e^{-t}$ on a

$$\widehat{\Delta}_n^\ell(\mathcal{F}) \leq \mathbb{E}[\widehat{\Delta}_n^\ell(\mathcal{F})] + \Delta\ell(R)\sqrt{\frac{t}{2n}}.$$

Pour conclure la preuve du lemme, il suffit de montrer que

$$(4.9) \quad \mathbb{E}[\widehat{\Delta}_n^\ell(\mathcal{F})] \leq \frac{2\alpha}{n} \max_{x \in \mathcal{X}^n} \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i f(x_i) \right| \right].$$

En suivant exactement les mêmes arguments que dans la preuve du lemme 1.3 (en remplaçant $\mathbf{1}_{Y_i \neq h(X_i)}$ par $\ell(Y_i f(X_i))$) on obtient

$$\mathbb{E}[\widehat{\Delta}_n^\ell(\mathcal{F})] \leq \frac{2}{n} \max_{y \in \{-1, +1\}^n} \max_{x \in \mathcal{X}^n} \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i \ell(y_i f(x_i)) \right| \right].$$

Finalement, on utilise que ℓ est α -Lipschitz pour conclure : le principe de Contraction (proposition 5 de l'appendice B) nous donne

$$\begin{aligned} \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i \ell(y_i f(x_i)) \right| \right] &\leq \alpha \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i y_i f(x_i) \right| \right] \\ &= \alpha \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i f(x_i) \right| \right]. \end{aligned}$$

En combinant les deux dernières majorations on obtient (4.9) et le lemme 3.1 est démontré. \square

(a) Prouvons maintenant la borne (4.6). L'application $\beta \rightarrow \sum_{i=1}^n \sigma_i \sum_{j=1}^p \beta_j h_j(x_i)$ est linéaire, donc elle atteint à l'un des sommets de \mathcal{C} son maximum et son minimum sur la boule \mathcal{C} . En conséquence, on a

$$\mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i f(x_i) \right| \right] = R \mathbb{E}_\sigma \left[\max_{j=1, \dots, p} \left| \sum_{i=1}^n \sigma_i h_j(x_i) \right| \right].$$

En appliquant l'inégalité (3.6) avec

$$\mathcal{V} = \{(h_j(x_1), \dots, h_j(x_n)) : j = 1, \dots, p\},$$

dont le cardinal est inférieur à p , on obtient

$$\mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \sigma_i f(x_i) \right| \right] \leq R \sqrt{2n \log(2p)}.$$

La borne (4.6) découle du lemme 3.1.

(b) Démontrons maintenant la borne (4.7). On note $\|\cdot\|_{\mathcal{F}}$ la norme dans le RKHS. La formule de reproduction et l'inégalité de Cauchy-Schwarz nous donne

$$\left| \sum_{i=1}^n \sigma_i f(x_i) \right| = \left| \left\langle f, \sum_{i=1}^n \sigma_i k(x_i, \cdot) \right\rangle_{\mathcal{F}} \right| \leq \|f\|_{\mathcal{F}} \left\| \sum_{i=1}^n \sigma_i k(x_i, \cdot) \right\|_{\mathcal{F}}.$$

En appliquant de nouveau l'inégalité de Cauchy-Schwarz, on obtient

$$\begin{aligned} \mathbb{E}_{\sigma} \left[\sup_{\|f\|_{\mathcal{F}} \leq R} \left| \sum_{i=1}^n \sigma_i f(x_i) \right| \right] &\leq R \mathbb{E}_{\sigma} \left[\left\| \sum_{i=1}^n \sigma_i k(x_i, \cdot) \right\|_{\mathcal{F}} \right] \\ &\leq R \sqrt{\mathbb{E}_{\sigma} \left[\left\| \sum_{i=1}^n \sigma_i k(x_i, \cdot) \right\|_{\mathcal{F}}^2 \right]} \\ &\leq R \sqrt{\sum_{i=1}^n k(x_i, x_i) \mathbb{E}_{\sigma}[\sigma_i^2]} \leq R\sqrt{n}, \end{aligned}$$

où on a utilisé $k(x, x) \leq 1$ dans la dernière inégalité et $\mathbb{E}[\sigma_i \sigma_j] = 0$ pour $i \neq j$ dans l'avant-dernière. En combinant de nouveau la propriété de reproduction avec l'inégalité de Cauchy-Schwarz, on obtient

$$|f(x)| = |\langle f, k(x, \cdot) \rangle_{\mathcal{F}}| \leq R\sqrt{k(x, x)} \leq R.$$

Le lemme 3.1 donne alors

$$L(\widehat{h}_{\mathcal{F}}) \leq \widehat{L}_n^{\ell}(\widehat{f}_{\mathcal{F}}) + \frac{2\alpha R}{\sqrt{n}} + \Delta\ell(R)\sqrt{\frac{t}{2n}},$$

ce qui achève la preuve du théorème 3. \square

Il est possible de démontrer des bornes de risque similaire à (3.3) pour $L(\widehat{h}_{\mathcal{H}_C})$, nous renvoyons à Bousquet, Boucheron et Lugosi [1] pour un panorama de tels résultats. Dans la dernière partie de ces notes, nous allons décrire deux algorithmes de classification très populaires : les Support Vector Machines et AdaBoost.

4.3. Support Vector Machine

Les Support Vector Machines (SVM) correspondent à l'estimateur (4.3) avec la perte Hinge $\ell(z) = (1 - z)_+$. La classification finale s'effectue avec $\widehat{h}_{\mathcal{F}}(x) = \text{sign}(\widehat{f}_{\mathcal{F}}(x))$. Nous allons donner une interprétation géométrique de la solution $\widehat{f}_{\mathcal{F}}$, interprétation qui est à l'origine du nom « Support Vector Machine ».

Proposition 2 (Vecteurs supports). *La solution de (4.3) avec $\ell(z) = (1 - z)_+$ est de la forme $\widehat{f}_{\mathcal{F}}(x) = \sum_{i=1}^n \widehat{\beta}_i k(X_i, x)$ avec*

$$\begin{cases} \widehat{\beta}_i = 0 & \text{si } Y_i \widehat{f}_{\mathcal{F}}(X_i) > 1 \\ \widehat{\beta}_i = Y_i / (2\lambda n) & \text{si } Y_i \widehat{f}_{\mathcal{F}}(X_i) < 1 \\ 0 \leq Y_i \widehat{\beta}_i \leq 1 / (2\lambda n) & \text{si } Y_i \widehat{f}_{\mathcal{F}}(X_i) = 1. \end{cases}$$

Les vecteurs X_i d'indice i tel que $\widehat{\beta}_i \neq 0$ sont appelés vecteurs supports.

Démonstration. Notons K la matrice $[k(X_i, X_j)]_{i,j=1,\dots,n}$. Nous savons par (4.5) que la solution de (4.3) est de la forme

$$\widehat{f}_{\mathcal{F}} = \sum_{j=1}^n \widehat{\beta}_j k(X_j, \cdot)$$

avec
$$\widehat{\beta} = \operatorname{argmin}_{\beta \in \mathbb{R}^n} \left\{ \frac{1}{n} \sum_{i=1}^n (1 - Y_i [K\beta]_i)_+ + \lambda \beta^T K \beta \right\}.$$

Le problème de minimisation ci-dessus n'est pas régulier, donc on introduit des variables de relâchement $\widehat{\xi}_i = (1 - Y_i [K\widehat{\beta}]_i)_+$ et on reformule le problème de minimisation comme suit

$$(4.10) \quad (\widehat{\beta}, \widehat{\xi}) = \operatorname{argmin}_{\substack{\beta, \xi \in \mathbb{R}^n \text{ tels que} \\ \xi_i \geq 1 - Y_i [K\beta]_i \\ \xi_i \geq 0}} \left\{ \frac{1}{n} \sum_{i=1}^n \xi_i + \lambda \beta^T K \beta \right\}.$$

Ce problème est maintenant régulier et convexe et les conditions de Karush-Kuhn-Tucker pour le problème lagrangien dual

$$(\widehat{\beta}, \widehat{\xi}) = \operatorname{argmin}_{\beta, \xi \in \mathbb{R}^n} \left\{ \frac{1}{n} \sum_{i=1}^n \xi_i + \lambda \beta^T K \beta - \sum_{i=1}^n (\alpha_i (\xi_i - 1 + Y_i [K\beta]_i) + \gamma_i \xi_i) \right\}.$$

donnent les formules

conditions de 1^{er} ordre :

$$2\lambda [K\widehat{\beta}]_j = \sum_{i=1}^n K_{ij} \alpha_i Y_i \quad \text{et} \quad \alpha_j + \gamma_j = \frac{1}{n},$$

conditions de relâchement :

$$\min(\alpha_i, \widehat{\xi}_i - 1 + Y_i [K\widehat{\beta}]_i) = 0 \quad \text{et} \quad \min(\gamma_i, \widehat{\xi}_i) = 0.$$

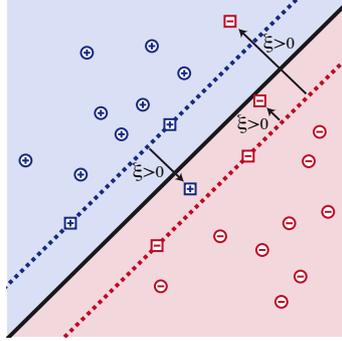


FIGURE 3. Classification par un SVM linéaire : l'hyperplan frontière $\{x \in \mathbb{R}^d : \langle \hat{w}, x \rangle = 0\}$ est représenté en noir, les hyperplans de marge $\{x \in \mathbb{R}^d : \langle \hat{w}, x \rangle = +1\}$ et $\{x \in \mathbb{R}^d : \langle \hat{w}, x \rangle = -1\}$ sont représentés en pointillés bleus et rouges respectivement. Les vecteurs supports sont représentés par des carrés.

On déduit de la première condition de premier ordre que $\hat{\beta}_i = \alpha_i Y_i / (2\lambda)$. Comme $\hat{f}_{\mathcal{F}}(X_i) = [K\hat{\beta}]_i$, la première condition de relâchement donne $\hat{\beta}_i = 0$ si $Y_i \hat{f}_{\mathcal{F}}(X_i) > 1$. La seconde condition de relâchement couplée à la seconde condition de premier ordre implique que $\hat{\beta}_i = Y_i / (2\lambda n)$ si $\hat{\xi}_i > 0$ et $0 \leq Y_i \hat{\beta}_i \leq 1 / (2\lambda n)$ sinon. Pour conclure la preuve de la proposition, on remarque que lorsque que $\hat{\xi}_i > 0$ on a $\hat{\beta}_i$ et α_i non nuls, et donc $Y_i \hat{f}_{\mathcal{F}}(X_i) = 1 - \hat{\xi}_i < 1$ vu la première condition de relâchement. \square

Interprétons géométriquement la proposition 2.

Interprétation géométrique : noyau linéaire. Considérons le noyau le plus simple $k(x, y) = \langle x, y \rangle$ pour tout $x, y \in \mathbb{R}^d$. Le RKHS associé est l'espace des formes linéaires $\mathcal{F} = \{\langle w, \cdot \rangle : w \in \mathbb{R}^d\}$. Dans ce cas

$$\hat{f}_{\mathcal{F}}(x) = \sum_{i=1}^n \hat{\beta}_i \langle X_i, x \rangle = \langle \hat{w}, x \rangle \quad \text{avec} \quad \hat{w} = \sum_{i=1}^n \hat{\beta}_i X_i,$$

donc le classifieur $\hat{h}_{\mathcal{F}}(x) = \text{sign}(\langle \hat{w}, x \rangle)$ classe les points en fonction de leur position par rapport à l'hyperplan $\{x \in \mathbb{R}^d : \langle \hat{w}, x \rangle = 0\}$. La normale à l'hyperplan \hat{w} est une combinaison linéaire des vecteurs supports, qui sont les points X_i tels que $Y_i \langle \hat{w}, X_i \rangle \leq 1$. Ils sont

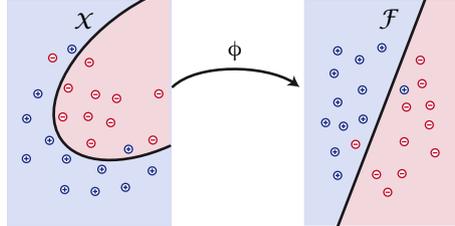


FIGURE 4. Classification avec un noyau non linéaire : la classification linéaire dans \mathcal{F} produit une classification non linéaire dans \mathcal{X} via l'image réciproque de ϕ .

représentés par des carrés sur la figure 3. Les hyperplans

$$\{x \in \mathbb{R}^d : \langle \hat{w}, x \rangle = +1\} \quad \text{et} \quad \{x \in \mathbb{R}^d : \langle \hat{w}, x \rangle = -1\}$$

sont habituellement appelés hyperplans de marge. Pour finir, notons la propriété suivante des SVM. Si on ajoute à l'ensemble d'apprentissage $(X_i, Y_i)_{i=1, \dots, n}$ un point X_{n+1} vérifiant $Y_{n+1} \langle \hat{w}, X_{n+1} \rangle > 1$, alors le vecteur \hat{w} et le classifieur $\hat{h}_{\mathcal{F}}$ ne changent pas. En d'autres mots, seules les données mal classifiées ou classifiées avec une marge insuffisante (i.e. $Y_i \langle \hat{w}, X_i \rangle \leq 1$) influencent l'hyperplan frontière $\{x \in \mathbb{R}^d : \langle \hat{w}, x \rangle = 0\}$.

Interprétation géométrique : noyau défini positif quelconque. Notons par $\phi : \mathcal{X} \rightarrow \mathcal{F}$ l'application $\phi(x) = k(x, \cdot)$. La propriété de reproduction combinée à la proposition 2 nous donne

$$\hat{f}_{\mathcal{F}}(x) = \langle \hat{f}_{\mathcal{F}}, \phi(x) \rangle_{\mathcal{F}} = \left\langle \sum_{i=1}^n \hat{\beta}_i \phi(X_i), \phi(x) \right\rangle_{\mathcal{F}}.$$

Un point $x \in \mathcal{X}$ est classifié selon le signe du produit scalaire ci-dessus. En conséquence, les points $\phi(x) \in \mathcal{F}$ sont classifiés selon le classifieur linéaire sur \mathcal{F}

$$f \mapsto \text{sign}(\langle \hat{w}_{\phi}, f \rangle_{\mathcal{F}}) \quad \text{où} \quad \hat{w}_{\phi} = \sum_{i=1}^n \hat{\beta}_i \phi(X_i).$$

La frontière $\{x \in \mathcal{X} : \hat{f}_{\mathcal{F}}(x) = 0\}$ du classifieur $\hat{h}_{\mathcal{F}}$ est donc l'image réciproque par ϕ de l'hyperplan $\{f \in \mathcal{F} : \langle \hat{w}_{\phi}, f \rangle_{\mathcal{F}} = 0\}$, comme représenté figure 4. Nous observons que le noyau k délinéarise le SVM, dans

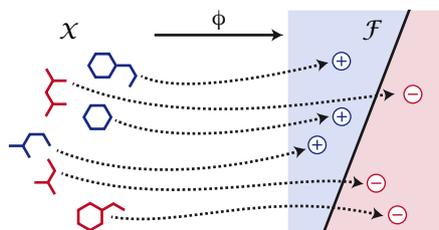


FIGURE 5. Classification de molécules à l'aide d'un SVM.

le sens où il produit un classifieur non linéaire $\hat{h}_{\mathcal{F}}$ pour le même coût numérique qu'un classifieur linéaire dans \mathbb{R}^n .

Vous pouvez observer les SVM à l'œuvre avec la petite démo :
<http://svm.cs.rhul.ac.uk/pagesnew/GPat.shtml>

Pourquoi les RKHS sont-ils utiles ? Il y a principalement deux raisons pour utiliser un RKHS. La première raison est de permettre de délinéariser un algorithme en envoyant \mathcal{X} dans \mathcal{F} avec $\phi : x \rightarrow k(x, \cdot)$, comme représenté dans la figure 4. Cela permet d'obtenir un algorithme non linéaire pour un coût similaire à celui d'un algorithme linéaire.

La seconde raison pour utiliser un RKHS est de pouvoir employer sur n'importe quel ensemble \mathcal{X} des algorithmes définis pour des vecteurs. Imaginons par exemple qu'on veuille classifier des protéines ou des molécules en fonction de leur propriétés thérapeutiques. Notons par \mathcal{X} notre ensemble de molécules. Pour tout $x, y \in \mathcal{X}$, représentons par $k(x, y)$ leur similarité (selon des critères à définir). Si le noyau résultant $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ est défini positif, on peut alors directement appliquer un SVM pour les classifier, comme illustré figure 5. Bien sûr, le point clef dans ce cas est de construire le noyau k . En général, le noyau $k(x, y)$ est défini en fonction de certaines propriétés de x, y qui sont connues pour avoir de l'importance pour le problème de classification. Par exemple, le nombre de courtes séquences communes est un indice utile pour quantifier la similarité entre deux protéines. La complexité du calcul de $k(x, y)$ est aussi un critère crucial pour les applications sur des données complexes. Nous renvoyons aux excellents slides de Jean-Philippe Vert pour la

description d'applications prometteuses en biologie et médecine :
<https://members.cbio.mines-paristech.fr/~jvert/talks/120302ensae/ensae.pdf>

4.4. AdaBoost

AdaBoost est un algorithme qui tend à calculer une solution approximative de l'estimateur (4.1) avec la perte exponentielle $\ell(z) = e^{-z}$ et l'espace fonctionnel $\mathcal{F} = \text{span}\{h_1, \dots, h_p\}$ où h_1, \dots, h_p sont p classifieurs donnés.

Le principe de l'algorithme AdaBoost est de réaliser une minimisation dite agressive de (4.1). Plus précisément, AdaBoost produit une suite de fonctions \widehat{f}_m pour $m = 0, \dots, M$ en partant de $\widehat{f}_0 = 0$ puis en résolvant récursivement pour $m = 1, \dots, M$

$$\widehat{f}_m = \widehat{f}_{m-1} + \beta_m h_{j_m}$$

$$\text{où } (\beta_m, j_m) = \underset{\substack{j=1, \dots, p \\ \beta \in \mathbb{R}}}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n \exp\left(-Y_i (\widehat{f}_{m-1}(X_i) + \beta h_j(X_i))\right).$$

La classification finale est effectuée à l'aide de $\widehat{h}_M(x) = \text{sign}(\widehat{f}_M(x))$ qui est une approximation de $\widehat{h}_{\mathcal{H}}$ défini par (4.1).

La perte exponentielle permet de calculer (β_m, j_m) très efficacement. En effet, en posant $w_i^{(m)} = n^{-1} \exp(-Y_i \widehat{f}_{m-1}(X_i))$, on peut écrire

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \exp\left(-Y_i (\widehat{f}_{m-1}(X_i) + \beta h_j(X_i))\right) \\ = (e^\beta - e^{-\beta}) \sum_{i=1}^n w_i^{(m)} \mathbf{1}_{h_j(X_i) \neq Y_i} + e^{-\beta} \sum_{i=1}^n w_i^{(m)}. \end{aligned}$$

Lorsque la condition suivante

$$\text{err}_m(j) = \frac{\sum_{i=1}^n w_i^{(m)} \mathbf{1}_{h_j(X_i) \neq Y_i}}{\sum_{i=1}^n w_i^{(m)}} \leq \frac{1}{2} \quad \text{pour tout } j = 1, \dots, p,$$

est satisfaite, les minimiseurs (β_m, j_m) sont donnés par

$$j_m = \underset{j=1, \dots, p}{\text{argmin}} \text{err}_m(j) \quad \text{et} \quad \beta_m = \frac{1}{2} \log\left(\frac{1 - \text{err}_m(j_m)}{\text{err}_m(j_m)}\right).$$

En remarquant que $-Y_i h(X_i) = 2\mathbf{1}_{Y_i \neq h(X_i)} - 1$ on obtient la formulation standard de l'algorithme AdaBoost.

AdaBoost
Init : $w_i^{(1)} = 1/n$, for $i = 1, \dots, n$ Iterate : For $m = 1, \dots, M$ do $j_m = \underset{j=1, \dots, p}{\operatorname{argmin}} \operatorname{err}_m(j)$ $2\beta_m = \log(1 - \operatorname{err}_m(j_m)) - \log(\operatorname{err}_m(j_m))$ $w_i^{(m+1)} = w_i^{(m)} \exp(2\beta_m \mathbf{1}_{h_{j_m}(X_i) \neq Y_i}), \quad \text{for } i = 1, \dots, n$ STOP if $\min_{j=1, \dots, p} \operatorname{err}_{m+1}(j) > 1/2$ Output : $\hat{f}_M(x) = \sum_{m=1}^M \beta_m h_{j_m}(x)$.

On remarque que l'algorithme AdaBoost donne de plus en plus de poids dans $\operatorname{err}_m(j)$ aux points X_i qui sont mal classifiés à l'étape m .

Vous pouvez observer AdaBoost à l'œuvre (pour des classifieurs h_j par demi-espaces) avec l'application récréative :
<http://cseweb.ucsd.edu/~yfreund/adaboost/>

5. Pour aller au-delà de cette introduction

Nous renvoyons le lecteur intéressé à aller au-delà des concepts présentés dans ces notes à l'article de synthèse de Boucheron, Bousquet et Lugosi [1]. Cet article comprend notamment une bibliographie assez exhaustive. D'un point de vue plus pratique et appliqué, l'incontournable livre de Hastie, Tibshirani et Friedman [5] décrit et discute de très nombreux algorithmes opérationnels. Enfin, nous soulignons que les concepts introduits ici apparaissent aussi pour le problème de *ranking* (ordonner au mieux des données, l'exemple le plus célèbre étant les moteurs de recherche sur internet), voir Cléménçon, Lugosi et Vayatis [3].

Appendice A. Espace de Hilbert à noyaux reproduisants (RKHS)

Les espaces de Hilbert à noyaux reproduisants, dit *Reproducing Kernel Hilbert Spaces* (RKHS) en anglais, sont des espaces de Hilbert

fonctionnels dont la norme décrit la régularité des fonctions. Les RKHS vérifient aussi une propriété de reproduction particulière, qui s'avère cruciale en pratique car elle permet une mise œuvre numérique efficace.

Une fonction $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ est dite définie positive si elle est symétrique ($k(x, y) = k(y, x)$ pour tout $x, y \in \mathcal{X}$) et si pour tout $N \in \mathbb{N}$, $x_1, \dots, x_N \in \mathcal{X}$ et $a_1, \dots, a_N \in \mathbb{R}$ on a

$$(A.1) \quad \sum_{i,j=1}^N a_i a_j k(x_i, x_j) \geq 0.$$

Exemples de noyaux définis positifs :

- noyau linéaire : $k(x, y) = \langle x, y \rangle$
- noyau gaussien : $k(x, y) = e^{-\|x-y\|^2/2\sigma^2}$
- noyau histogramme ($d = 1$) : $k(x, y) = \min(x, y)$
- noyau exponentiel : $k(x, y) = e^{-\|x-y\|/\sigma}$

On peut associer à un noyau défini positif k un sous-espace de Hilbert \mathcal{F} de $\mathbb{R}^{\mathcal{X}}$ appelé espace de Hilbert de noyau reproduisant k .

Proposition 3 (Espace de Hilbert à noyau reproduisant (RKHS))

À tout noyau défini positif k sur \mathcal{X} , on peut associer un (unique) espace de Hilbert $\mathcal{F} \subset \mathbb{R}^{\mathcal{X}}$ vérifiant

- $k(x, \cdot) \in \mathcal{F}$ pour tout $x \in \mathcal{X}$
- propriété de reproduction : $f(x) = \langle f, k(x, \cdot) \rangle_{\mathcal{F}}$ pour tout $x \in \mathcal{X}$ et $f \in \mathcal{F}$.

L'espace \mathcal{F} est appelé espace de Hilbert de noyau reproduisant k .

Démonstration. Considérons l'espace linéaire \mathcal{F}_0 engendré par la famille $\{k(x, \cdot) : x \in \mathcal{X}\}$

$$\mathcal{F}_0 = \left\{ f : \mathcal{X} \rightarrow \mathbb{R} : f(x) = \sum_{i=1}^N a_i k(x_i, x), \quad N \in \mathbb{N}, \right. \\ \left. x_1, \dots, x_N \in \mathcal{X}, \quad a_1, \dots, a_N \in \mathbb{R} \right\}.$$

À tout $f = \sum_{i=1}^N a_i k(x_i, \cdot)$ et $g = \sum_{j=1}^M b_j k(y_j, \cdot)$ on associe

$$\langle f, g \rangle_{\mathcal{F}_0} := \sum_{i=1}^N \sum_{j=1}^M a_i b_j k(x_i, y_j) = \sum_{i=1}^N a_i g(x_i) = \sum_{j=1}^M b_j f(y_j).$$

Les deux dernières égalités indiquent que $\langle f, g \rangle_{\mathcal{F}_0}$ ne dépend pas du choix de la décomposition de f et g , donc cette quantité est bien définie. De plus, l'application $(f, g) \rightarrow \langle f, g \rangle_{\mathcal{F}_0}$ est bilinéaire, symétrique, positive (vu (A.1)) et on a la propriété de reproduction

$$f(x) = \langle f, k(x, \cdot) \rangle_{\mathcal{F}_0} \quad \text{pour tout } x \in \mathcal{X} \quad \text{et } f \in \mathcal{F}_0.$$

L'inégalité de Cauchy-Schwarz $\langle f, g \rangle_{\mathcal{F}_0} \leq \|f\|_{\mathcal{F}_0} \|g\|_{\mathcal{F}_0}$ combinée à la formule de reproduction donne

$$(A.2) \quad |f(x)| \leq \sqrt{k(x, x)} \|f\|_{\mathcal{F}_0}.$$

En conséquence $\|f\|_{\mathcal{F}_0} = 0$ implique $f = 0$ donc $\langle f, g \rangle_{\mathcal{F}_0}$ est un produit scalaire. L'espace \mathcal{F}_0 est donc un espace pré-hilbertien et pour obtenir \mathcal{F} il nous suffit de le compléter. Considérons deux suites $(x_i) \in \mathcal{X}^{\mathbb{N}}$ et $(a_i) \in \mathbb{R}^{\mathbb{N}}$ vérifiant $\sum_{i,j \geq 1} a_i a_j k(x_i, x_j) < +\infty$. L'inégalité (A.2) nous assure que pour tout $M < N$ et $x \in \mathcal{X}$ on a

$$\left| \sum_{i=M+1}^N a_i k(x_i, x) \right| \leq \sqrt{k(x, x)} \sum_{i,j=M+1}^N a_i a_j k(x_i, x_j).$$

Lorsque $\sum_{i,j \geq 1} a_i a_j k(x_i, x_j)$ est fini, le membre de droite tend vers 0 lorsque M, N tendent vers l'infini, donc les séries partielles $\sum_{i=1}^N a_i k(x_i, x)$ sont de Cauchy et convergent lorsque $N \rightarrow \infty$. On peut donc définir l'espace

$$\mathcal{F} = \left\{ f : \mathcal{X} \rightarrow \mathbb{R} : f(x) = \sum_{i=1}^{\infty} a_i k(x_i, x), \quad (x_i) \in \mathcal{X}^{\mathbb{N}}, \right. \\ \left. (a_i) \in \mathbb{R}^{\mathbb{N}}, \quad \sum_{i,j \geq 1} a_i a_j k(x_i, x_j) < +\infty \right\}$$

et la forme bilinéaire

$$\langle f, g \rangle_{\mathcal{F}} := \sum_{i,j=1}^{\infty} a_i b_j k(x_i, y_j) = \sum_{i=1}^{\infty} a_i g(x_i) = \sum_{j=1}^{\infty} b_j f(y_j)$$

pour $f = \sum_{i=1}^{\infty} a_i k(x_i, \cdot)$ et $g = \sum_{j=1}^{\infty} b_j k(y_j, \cdot)$. Exactement comme précédemment, l'application $(f, g) \rightarrow \langle f, g \rangle_{\mathcal{F}}$ est un produit scalaire vérifiant la propriété de reproduction

$$f(x) = \langle f, k(x, \cdot) \rangle_{\mathcal{F}} \quad \text{pour tout } x \in \mathcal{X} \text{ et } f \in \mathcal{F}.$$

Finalement, l'espace \mathcal{F} muni du produit scalaire $\langle \cdot, \cdot \rangle_{\mathcal{F}}$ est la complétion dans $\mathbb{R}^{\mathcal{X}}$ de \mathcal{F}_0 muni de $\langle \cdot, \cdot \rangle_{\mathcal{F}_0}$, donc c'est un espace de Hilbert à noyau reproduisant. \square

La norme d'une fonction f dans un RKHS \mathcal{F} est fortement liée à la régularité de f . Cette propriété apparaît clairement dans l'inégalité $|f(x) - f(x')| = |\langle f, k(x, \cdot) - k(x', \cdot) \rangle_{\mathcal{F}}| \leq \|f\|_{\mathcal{F}} \|k(x, \cdot) - k(x', \cdot)\|_{\mathcal{F}}$. Nous illustrons ce point en décrivant le RKHS associé aux noyaux histogrammes et gaussiens.

Exemple 1 : RKHS associé au noyau histogramme

L'espace de Sobolev

$$\mathcal{F} = \left\{ f \in C([0, 1], \mathbb{R}) : f \text{ est p.p. différentiable} \right. \\ \left. \text{avec } f' \in L^2([0, 1]) \text{ et } f(0) = 0 \right\}$$

muni du produit scalaire $\langle f, g \rangle_{\mathcal{F}} = \int_0^1 f' g'$ est un RKHS de noyau reproduisant $k(x, y) = \min(x, y)$ sur $[0, 1]$. En effet, $k(x, \cdot) \in \mathcal{F}$ pour tout $x \in [0, 1]$ et

$$f(x) = \int_0^1 f'(y) \mathbf{1}_{y \leq x} dy = \langle f, k(x, \cdot) \rangle_{\mathcal{F}}, \quad \text{pour tout } f \in \mathcal{F} \text{ et } x \in [0, 1].$$

Dans ce cas, la norme $\|f\|_{\mathcal{F}}$ correspond simplement à la norme L^2 de la dérivée de f . Plus cette norme est petite, plus f est régulière.

Exemple 2 : RKHS associé au noyau gaussien. Notons par $\mathbf{F}[f]$ la transformée de Fourier dans \mathbb{R}^d de normalisation

$$\mathbf{F}[f](\omega) = \frac{1}{(2\pi)^{d/2}} \int_{\mathbb{R}^d} f(t) e^{-i\langle \omega, t \rangle}, \\ \text{pour } f \in L^1(\mathbb{R}^d) \cap L^2(\mathbb{R}^d) \text{ et } \omega \in \mathbb{R}^d.$$

Pour tout $\sigma > 0$, l'espace fonctionnel

$$\mathcal{F}_{\sigma} = \left\{ f \in C_0(\mathbb{R}^d) \cap L^1(\mathbb{R}^d) \text{ tel que } \int_{\mathbb{R}^d} |\mathbf{F}[f](\omega)|^2 e^{\sigma|\omega|^2/2} d\omega < +\infty \right\},$$

muni du produit scalaire

$$\langle f, g \rangle_{\mathcal{F}_\sigma} = (2\pi\sigma^2)^{-d/2} \int_{\mathbb{R}^d} \overline{\mathbf{F}[f](\omega)} \mathbf{F}[g](\omega) e^{\sigma|\omega|^2/2} d\omega$$

est un RKHS associé au noyau gaussien

$$k(x, y) = \exp(-\|y - x\|^2/2\sigma^2).$$

En effet, pour tout $x \in \mathbb{R}^d$ la fonction $k(x, \cdot)$ appartient à \mathcal{F}_σ et un calcul direct donne

$$\langle k(x, \cdot), f \rangle_{\mathcal{F}_\sigma} = \mathbf{F}^{-1}[\mathbf{F}[f]](x) = f(x) \quad \text{pour tout } f \in \mathcal{F} \text{ et } x \in \mathbb{R}^d.$$

L'espace \mathcal{F}_σ est composé de fonctions très régulières et la norme $\|f\|_{\mathcal{F}_\sigma}$ contrôle directement la régularité de f . Remarquons que lorsque σ augmente l'espace \mathcal{F}_σ rétrécit pour ne contenir que des fonctions de plus en plus régulières.

Appendice B. Inégalités probabilistes

Les analyses non asymptotiques en apprentissage statistique reposent très souvent sur l'inégalité de concentration de McDiarmid [6].

Théorème 4 (McDiarmid (1989)). *Soit \mathcal{X} un ensemble mesurable et considérons $F : \mathcal{X}^n \rightarrow \mathbb{R}$ telle qu'il existe $\delta_1, \dots, \delta_n$ vérifiant*

$$|F(x_1, \dots, x'_i, \dots, x_n) - F(x_1, \dots, x_i, \dots, x_n)| \leq \delta_i,$$

pour tout $x_1, \dots, x_n, x'_i \in \mathcal{X}$, pour $i = 1, \dots, n$. Alors, pour tout $t > 0$ et pour toutes variables indépendantes X_1, \dots, X_n , on a

$$\mathbb{P}(F(X_1, \dots, X_n) > \mathbb{E}[F(X_1, \dots, X_n)] + t) \leq \exp\left(-\frac{2t^2}{\delta_1^2 + \dots + \delta_n^2}\right).$$

Nous renvoyons le lecteur au Chapitre 9 du livre de Devroye, Györfi et Lugosi [4] pour la preuve classique de ce résultat, preuve qui combine simplement l'inégalité de Markov avec un argument martingale (par conditionnements télescopiques). Une preuve plus conceptuelle basée sur la méthode d'entropie est donnée dans le livre de Boucheron, Lugosi et Massart [2], Chapitre 6.

Un second concept important en apprentissage statistique est le principe de contraction.

Théorème 5 (Principe de contraction). Soit \mathcal{Z} un sous-ensemble borné de \mathbb{R}^n et $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ une fonction α -Lipschitz vérifiant $\varphi(0) = 0$. Pour des variables aléatoires $\sigma_1, \dots, \sigma_n$ i.i.d. de loi

$$\mathbb{P}_\sigma(\sigma_i = 1) = \mathbb{P}_\sigma(\sigma_i = -1) = 1/2,$$

nous avons

$$\mathbb{E}_\sigma \left[\sup_{z \in \mathcal{Z}} \left| \sum_{i=1}^n \sigma_i \varphi(z_i) \right| \right] \leq \alpha \mathbb{E}_\sigma \left[\sup_{z \in \mathcal{Z}} \left| \sum_{i=1}^n \sigma_i z_i \right| \right].$$

Une preuve concise est donnée dans le Chapitre 11 du livre de Boucheron, Lugosi et Massart [2].

Références

- [1] S. BOUCHERON, O. BOUSQUET & G. LUGOSI – « Theory of classification : a survey of some recent advances », *ESAIM Probab. Stat.* **9** (2005), p. 323–375.
- [2] S. BOUCHERON, G. LUGOSI & P. MASSART – *Concentration inequalities. A nonasymptotic theory of independence*, Oxford University Press, Oxford, 2013.
- [3] S. CLÉMENÇON, G. LUGOSI & N. VAYATIS – « Ranking and empirical minimization of U-statistics », *The Annals of Statistics* **36** (2008), no. 2, p. 844–874.
- [4] L. DEVROYE, L. GYÖRFI & G. LUGOSI – *A probabilistic theory of pattern recognition*, Applications of Math., vol. 31, Springer-Verlag, New York, 1996.
- [5] T. HASTIE, R. TIBSHIRANI & J. FRIEDMAN – *The elements of statistical learning. Data mining, inference, and prediction*, 2^e éd., Springer Series in Statistics, Springer, New York, 2009.
- [6] C. MCDIARMID – « On the method of bounded differences », in *Surveys in combinatorics, 1989 (Norwich, 1989)*, London Math. Soc. Lecture Note Ser., vol. 141, Cambridge University Press, Cambridge, 1989, p. 148–188.

Christophe Giraud, CMAP, UMR CNRS 7641, École Polytechnique, 91128 Palaiseau Cedex
 E-mail : christophe.giraud@math.u-psud.fr