

Apprentissage: cours 2

Méthodes par moyennage local -

Consistance des méthodes par partition

Simon Lacoste-Julien

25 septembre 2015

Référence : chap. 6 of [Hastie et al., 2009] and chap. 6 of [Devroye et al., 1996].

Résumé

On obtient notre première analyse théorique d'apprentissage (on va faire des preuves!). On considère des algorithmes très simples pour la régression et la classification (le moyennage local), et on prouve déjà la consistance universelle d'un algorithme de classification. On montre aussi comment obtenir une règle de classification à partir de régression, grâce à l'astuce du 'plug-in'.

1 Algorithmes par moyennage local

On considère la régression au sens des moindres carrés avec des entrées dans $\mathcal{X} = \mathbb{R}^d$ et des sorties réelles bornées : $\mathcal{Y} = [-B, B]$ pour $B > 0$ et $\ell(y, y') = \frac{1}{2}(y - y')^2$. Une fonction cible est donc $f^*(x) = \mathbb{E}[Y|X = x]$. On considère un ensemble d'entraînement $D_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$.

Principe des méthodes par moyennage local : Prédire par la moyenne pondérée des Y_i pour des X_i voisins de x . On considère les prédicteurs de la forme

$$\hat{\eta} : x \mapsto \sum_{i=1}^n W_i(x) Y_i,$$

où $\sum_{i=1}^n W_i(x) = 1$ (les poids sont normalisés).

1) Algorithme par partition : méthode d'histogrammes

On fixe en avance une partition $\{A_1, A_2, \dots\}$ finie ou dénombrable de \mathcal{X} . Soit $A(x)$ l'élément de la partition contenant x (cette notation reviendra durant toute cette note – à retenir!). On choisit les poids :

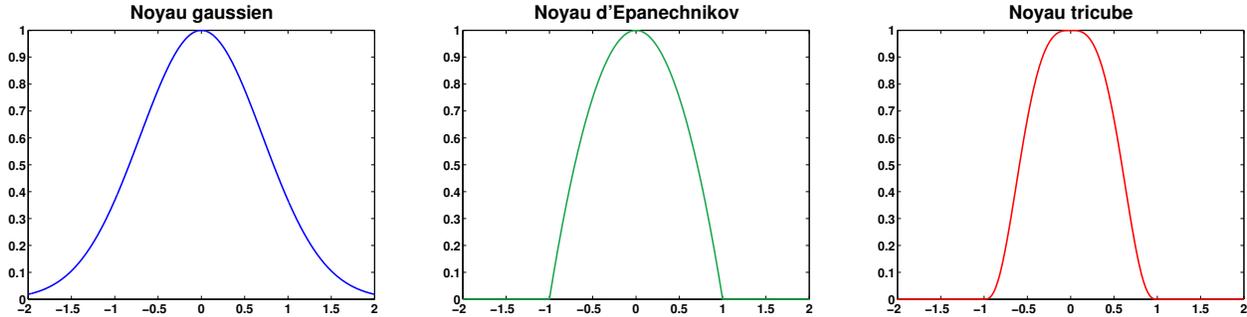
$$W_i(x) = \frac{\mathbb{1}_{\{X_i \in A(x)\}}}{\sum_{l=1}^n \mathbb{1}_{\{X_l \in A(x)\}}} = \frac{\mathbb{1}_{\{X_i \in A(x)\}}}{N_A(x)}, \quad \text{où } N_A(x) := \sum_{l=1}^n \mathbb{1}_{\{X_l \in A(x)\}},$$

avec la convention $\frac{0}{0} = 0$. Ici, $\hat{\eta}$ fait une prédiction constante sur chacun des A_j qui est tout simplement la moyenne des Y_i dans cette cellule. La fonction de prédiction ressemble donc en effet à un 'histogramme'.

2) Algorithme des k plus proches voisins (k -p.p.v.)

Ici, c'est comme la méthode d'histogrammes, mais avec une partition *définie par les données*. On définit les k plus proches voisins de x comme un ensemble $V_k(x)$ de k éléments de $\mathcal{X}_n = \{X_1, \dots, X_n\}$ tel que $\forall (X_i, X_j) \in V_k(x) \times \mathcal{X}_n \setminus V_k(x)$, $\|X_i - x\| \leq \|X_j - x\|$. Ce sont exactement les k -p.p.v. s'il n'y a pas d'ex aequo. On définit alors les poids :

$$W_i(x) = \frac{\mathbb{1}_{\{X_i \in V_k(x)\}}}{k}.$$



3) Algorithme par noyau : méthode de Nadaraya-Watson

On considère une fonction $K : \mathbb{R}^p \rightarrow \mathbb{R}_+$ appelé *noyau de convolution*. Les *noyaux de convolution* ont d'abord été utilisés par Parzen et Rosenblatt pour faire de l'estimation de densité (méthode des fenêtres de Parzen). On appelle h la *largeur de bande* du noyau.

Les poids de la méthode sont alors définis comme :

$$W_i(x) = \frac{K\left(\frac{x-X_i}{h}\right)}{\sum_{l=1}^n K\left(\frac{x-X_l}{h}\right)}$$

Quelques noyaux classiques sur \mathbb{R} :

- le noyau gaussien : $t \mapsto \exp(-t^2)$
- le noyau quadratique d'Epanechnikov : $t \mapsto (1-t^2)_+$
- le noyau "tricube" : $t \mapsto (1-|t|^3)_+^3$

où on a noté la fonction *partie positive* par $(x)_+ = \max(0, x)$.

Pour une liste d'autres noyaux on pourra consulter [http://en.wikipedia.org/wiki/Kernel_\(statistics\)](http://en.wikipedia.org/wiki/Kernel_(statistics)).

En dimension p on utilisera typiquement des noyaux de la forme $K : x \mapsto K_1(\|x\|)$ pour K_1 l'un des noyaux définis pour la dimension 1.

Exercice 1. A quoi correspond la méthode de Nadaraya-Watson pour un noyau gaussien lorsque $h \rightarrow 0$?

2 Analyse de l'algorithme par partition (histogrammes)

Définition 1 (Algorithme par partition). Soit $A = \{A_1, \dots, A_k, \dots\}$ une partition finie ou dénombrable de \mathcal{X} et $D_n = (X_i, Y_i)_{1 \leq i \leq n}$ un échantillon. Pour tout $x \in \mathcal{X}$, on note $A(x)$ l'élément de la partition qui contient x et $N_A(x) := \text{card} \{i \in \{1, \dots, n\} \text{ t.q. } X_i \in A(x)\}$. L'algorithme d'apprentissage par partition associé est défini comme suit.

— en régression :

$$\forall x \in \mathcal{X}, \quad \hat{\eta}_A(x; D_n) := \frac{1}{N_A(x)} \sum_{i=1}^n \mathbb{1}_{\{X_i \in A(x)\}} Y_i \quad (1)$$

— en classification binaire ($\mathcal{Y} = \{0, 1\}$) :

$$\forall x \in \mathcal{X}, \quad \hat{f}_A(x; D_n) := \mathbb{1}_{\{\hat{\eta}_A(x; D_n) \geq \frac{1}{2}\}} \quad (2)$$

Exemple classique dans $[0, 1]^d$ ou \mathbb{R}^d : partition régulière de pas $h > 0$.

2.1 Lien entre régression et classification pour la méthode plug-in

L'algorithme de classification par partition est un algorithme de type "plug-in" : on apprend d'abord la fonction de régression par $\hat{\eta}$, puis on considère le prédicteur pour la classification donné par $\hat{f} = \mathbb{1}_{\{\hat{\eta} \geq 1/2\}}$.

La motivation de cette approche est la suivante. Considérons le problème de la classification binaire : $\mathcal{X} = \mathbb{R}^d$ et $\mathcal{Y} = \{0, 1\}$. Comme $\{0, 1\} \subset \mathbb{R}$, on peut aussi considérer le problème de régression au sens des moindres carrés de Y sur X . Pour la régression la fonction cible est $\eta^*(x) = \mathbb{E}[Y|X=x] = \mathbb{P}(Y=1|X=x)$. Pour la classification la fonction cible est $g^*(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \mathbb{P}(Y=y|X=x) = \mathbb{1}_{\{\eta^*(x) \geq \frac{1}{2}\}}$.

Comme l'objectif de la régression dans ce cas-ci est d'apprendre à prédire la probabilité conditionnelle de Y sachant X il paraît naturel de se servir de ce prédicteur pour construire un prédicteur de classification.

Définition 2 (Estimateur plug-in). On suppose $\mathcal{Y} = \{0, 1\}$. Soit $\hat{\eta} : \mathcal{X} \mapsto \mathbb{R}$ un estimateur de la fonction de régression $\eta^* : x \mapsto \mathbb{E}[Y|X=x] = \mathbb{P}(Y=1|X=x)$. On lui associe la règle de classification "plug-in" \hat{f} défini par

$$\forall x \in \mathcal{X}, \quad \hat{f}(x) := \mathbb{1}_{\{\hat{\eta}(x) \geq \frac{1}{2}\}}. \quad (3)$$

Théorème 1. Soit $\hat{\eta}_n : \mathbb{R}^d \rightarrow \mathbb{R}$ un prédicteur pour la régression des moindres carrés et $\hat{f}_n : \mathbb{R}^d \rightarrow \{0, 1\}$ défini par $\hat{f}_n(x) = \mathbb{1}_{\{\hat{\eta}_n(x) \geq \frac{1}{2}\}}$. Alors le risque pour la régression de η noté $\mathcal{R}^{\text{reg}}(\hat{\eta}_n) := \frac{1}{2} \mathbb{E}[(\hat{\eta}_n(X) - Y)^2 | D_n]$ et le risque pour la classification de \hat{f}_n noté $\mathcal{R}^{0-1}(\hat{f}_n) := \mathbb{P}(\hat{f}_n(X) \neq Y | D_n)$ sont liés par la relation

$$\mathcal{R}^{0-1}(\hat{f}_n) - \mathcal{R}^{0-1}(f^*) \leq \sqrt{8} \sqrt{\mathcal{R}^{\text{reg}}(\hat{\eta}_n) - \mathcal{R}^{\text{reg}}(\eta^*)}.$$

Démonstration.

$$\begin{aligned} \mathcal{R}^{0-1}(\hat{f}_n) - \mathcal{R}^{0-1}(f^*) &= \mathbb{E} \left[|2\eta^*(X) - 1| \mathbb{1}_{\{\hat{f}_n(X) \neq f^*(X)\}} \mid D_n \right] \quad \dots \text{détails au tableau!} \\ &\leq 2 \mathbb{E} [|\hat{\eta}_n(X) - \eta^*(X)| \mid D_n] \quad \dots \text{détails au tableau!} \\ &\leq 2 \sqrt{\mathbb{E} [(\hat{\eta}_n(X) - \eta^*(X))^2 \mid D_n]} \quad \text{par inégalité de Jensen} \\ &= \sqrt{8} \sqrt{\mathcal{R}^{\text{reg}}(\hat{\eta}_n) - \mathcal{R}^{\text{reg}}(\eta^*)} \quad \dots \text{détails au tableau!} \end{aligned}$$

□

En intégrant par rapport à D_n , on obtient

$$\begin{aligned} \mathbb{E} \left[\mathcal{R}^{0-1}(\hat{f}_n) - \mathcal{R}^{0-1}(f^*) \right] &\leq 2 \mathbb{E} [|\hat{\eta}_n(X) - \eta^*(X)|] \\ &\leq \sqrt{8} \sqrt{\mathbb{E} [\mathcal{R}^{\text{reg}}(\hat{\eta}_n) - \mathcal{R}^{\text{reg}}(\eta^*)]} \end{aligned} \quad (4)$$

en utilisant la preuve du Théorème 1. En particulier, si $\hat{\eta}_n$ est un estimateur consistant de η^* au sens du risque des moindres carrés, alors l'estimateur plug-in associé est consistant au sens du risque 0-1.

2.2 Condition suffisante de consistance en classification pour les estimateur par partition

Théorème 2. On se place en classification binaire ($\mathcal{Y} = \{0, 1\}$) et l'on note $\mathcal{R}^{0-1}(\cdot)$ le risque associé à la perte 0-1. Soit $A_n = \{A_{1,n}, \dots, A_{k,n}, \dots\}_{n \in \mathbb{N}}$ une suite (déterministe) de partitions finies ou dénombrables de $\mathcal{X} = \mathbb{R}^d$. Soit \hat{f}_{A_n} l'algorithme par partition associé, donné par la Définition 1. Pour tout $E \subset \mathbb{R}^d$, on définit son diamètre (éventuellement infini) $\operatorname{diam}(E) = \sup_{x,y \in E} \{ \|x - y\| \}$ où $\|\cdot\|$ désigne la norme euclidienne dans \mathbb{R}^d . Si

1. $\operatorname{diam}(A_n(X)) \rightarrow 0$ en probabilité, et
2. $N_{A_n}(X) \rightarrow +\infty$ en probabilité,

alors, l' algorithme de classification $(\widehat{f}_{A_n})_{n \in \mathbb{N}}$ est consistant :

$$\mathbb{E} \left[\mathcal{R}^{0-1} \left(\widehat{f}_{A_n}(\cdot; D_n) \right) \right] - \mathcal{R}^{0-1}(f^*) \xrightarrow{n \rightarrow +\infty} 0 . \quad (5)$$

La démonstration du Théorème 2 repose notamment sur le lemme suivant (qui peut être prouvé par l'inégalité de Jensen et la variance d'une binomiale).

Lemme 1. Soit $N \geq 1$, $p \in [0, 1]$ et Z une variable de loi binomiale de paramètres N et p . Alors,

$$\mathbb{E} \left| \frac{Z - Np}{N} \right| \leq \frac{1}{2\sqrt{N}} . \quad (6)$$

Démonstration du Théorème 2. Notons $\widehat{\eta}_n = \widehat{\eta}_{A_n}$. Au vu de (4), il suffit de montrer que

$$\mathbb{E} |\widehat{\eta}_n(X; D_n) - \eta^*(X)| \xrightarrow{n \rightarrow +\infty} 0 .$$

Introduisons la fonction $\bar{\eta}_n : \mathcal{X} \mapsto \mathbb{R}$ définie par

$$\forall x \in \mathcal{X}, \quad \bar{\eta}_n(x) := \mathbb{E} [\eta^*(X) \mid X \in A_n(x)] = \mathbb{E}[Y \mid X \in A_n(x)] .$$

($\bar{\eta}_n$ est la meilleure fonction histogramme sur la partition A_n , c'est à dire, constante par morceau sur chacune des cellules). Alors, d'après l'inégalité triangulaire,

$$\mathbb{E} |\widehat{\eta}_n(X; D_n) - \eta^*(X)| \leq \underbrace{\mathbb{E} |\widehat{\eta}_n(X; D_n) - \bar{\eta}_n(X)|}_{\text{'variance'}} + \underbrace{\mathbb{E} |\bar{\eta}_n(X) - \eta^*(X)|}_{\text{'biais'}} ,$$

et il suffit de montrer que chacun de ces deux termes tend vers zéro.

Contrôle du premier terme $\mathbb{E} |\widehat{\eta}_n(X; D_n) - \bar{\eta}_n(X)|$ Pour tout $x \in \mathcal{X}$, conditionnellement à $N_{A_n}(x)$,

$$N_{A_n}(x) \widehat{\eta}_n(x; D_n) = \sum_{i / X_i \in A(x)} Y_i$$

est la somme de $N_{A_n}(x)$ variables de Bernoulli de paramètre $\mathbb{P}(Y_i = 1 \mid X_i \in A(x)) = \bar{\eta}_n(x)$, donc c'est une variable binomiale de paramètres $(N_{A_n}(x), \bar{\eta}_n(x))$. Ainsi,

$$\begin{aligned} & \mathbb{E} [|\widehat{\eta}_n(X; D_n) - \bar{\eta}_n(X)| \mid X, N_{A_n}(X)] \\ & \leq \mathbb{1}_{\{N_{A_n}(X)=0\}} + \frac{\mathbb{1}_{\{N_{A_n}(X)>0\}}}{N_{A_n}(X)} \mathbb{E} [|N_{A_n}(X) \widehat{\eta}_n(X; D_n) - N_{A_n}(X) \bar{\eta}_n(X)| \mid X, N_{A_n}(X)] \\ & \leq \mathbb{1}_{\{N_{A_n}(X)=0\}} + \frac{\mathbb{1}_{\{N_{A_n}(X)>0\}}}{2\sqrt{N_{A_n}(X)}} \quad \text{d'après le Lemme 1.} \end{aligned}$$

En intégrant, on en déduit

$$\begin{aligned} \mathbb{E} |\widehat{\eta}_n(X; D_n) - \bar{\eta}_n(X)| & \leq \mathbb{P}(N_{A_n}(X) = 0) + \mathbb{E} \left[\frac{\mathbb{1}_{\{N_{A_n}(X)>0\}}}{2\sqrt{N_{A_n}(X)}} \right] \\ & \leq \mathbb{P}(N_{A_n}(X) = 0) + \frac{\mathbb{P}(N_{A_n}(X) \leq k)}{2} + \frac{1}{2\sqrt{k}} \end{aligned}$$

pour tout $k > 0$. Comme $N_{A_n}(X) \rightarrow +\infty$ en probabilité, ceci implique

$$\limsup_{n \rightarrow +\infty} \mathbb{E} |\widehat{\eta}_n(X; D_n) - \bar{\eta}_n(X)| \leq \frac{1}{2\sqrt{k}}$$

pour tout $k > 0$, et donc

$$\lim_{n \rightarrow +\infty} \mathbb{E} |\widehat{\eta}_n(X; D_n) - \bar{\eta}_n(X)| = 0 .$$

Contrôle du deuxième terme $\mathbb{E} |\bar{\eta}_n(X) - \eta^*(X)|$ Soit $\varepsilon > 0$ quelconque. Comme l'ensemble des fonctions continues à support compact est dense dans $L_1(P_X)$ (où P_X désigne la loi commune des X_i), on peut trouver une fonction g continue à support compact (donc uniformément continue) telle que $\mathbb{E} |\eta^*(X) - g(X)| \leq \varepsilon$. Quitte à remplacer g par $\min\{1, \max\{0, g\}\}$, on peut supposer que g est à valeurs dans $[0, 1]$. On définit alors $\bar{g}_n : \mathcal{X} \mapsto \mathbb{R}$ par

$$\forall x \in \mathcal{X}, \quad \bar{g}_n(x) = \mathbb{E}[g(X) \mid X \in A_n(x)]$$

et l'inégalité triangulaire donne

$$\mathbb{E} |\bar{\eta}_n(X) - \eta^*(X)| \leq \mathbb{E} |\bar{\eta}_n(X) - \bar{g}_n(X)| + \mathbb{E} |\bar{g}_n(X) - g(X)| + \mathbb{E} |g(X) - \eta^*(X)| .$$

Le troisième terme est majoré par ε par définition de g . Le premier terme est inférieur au troisième (donc à ε) car le *conditionnement réduit la distance L1 entre les variables*. Plus spécifiquement, soit Z_1, Z_2 and B des variables aléatoires. Considérons $\bar{Z}_1 = \mathbb{E}[Z_1|B]$ et $\bar{Z}_2 = \mathbb{E}[Z_2|B]$. Alors nous avons $\mathbb{E}[|\bar{Z}_1 - \bar{Z}_2|] \leq \mathbb{E}[|Z_1 - Z_2|]$. Pour voir cela :

$$\mathbb{E}[|\bar{Z}_1 - \bar{Z}_2|] = \mathbb{E}[|\mathbb{E}[Z_1 - Z_2|B]|] \leq \mathbb{E}[\mathbb{E}[|Z_1 - Z_2| \mid B]] = \mathbb{E}[|Z_1 - Z_2|].$$

On applique ce résultat pour le premier terme en identifiant $X \in A_n(x)$ comme étant B ; et $Z_1 = g(X)$, $Z_2 = \eta^*(X)$.

Enfin, pour le deuxième terme, remarquons que g étant uniformément continue, il existe $\delta > 0$ tel que g varie d'au plus ε sur tout ensemble de diamètre inférieur à δ . Comme $g(X) \in [0, 1]$ p.s., on en déduit que

$$\mathbb{E} |\bar{g}_n(X) - g(X)| \leq \varepsilon + \mathbb{P}(\text{diam}(A_n(X)) > \delta) \leq 2\varepsilon$$

pour n assez grand, puisque $\text{diam}(A_n(X)) \rightarrow 0$ en probabilité. Comme $\varepsilon > 0$ peut être choisi arbitrairement petit, on a bien montré que $\mathbb{E} |\bar{\eta}_n(X) - \eta^*(X)| \rightarrow 0$ quand $n \rightarrow +\infty$, d'où le résultat. \square

2.3 Consistance universelle pour les partitions régulières en classification

Pour l'instant, le Théorème 2 nous donne une condition suffisante de consistance, mais sans préciser si cette condition peut être satisfaite sans connaître a priori la loi P des données. L'exemple suivant montre qu'un bon choix de partition de \mathbb{R}^d permet d'avoir la consistance universelle.

Définition 3. Pour tout $h > 0$, on définit $A^r(h)$ la *partition régulière de pas $h > 0$* de \mathbb{R}^d en cubes de taille h définis par la grille régulière $h\mathbb{Z}^d$. Autrement dit,

$$A^r(h) = \left\{ \prod_{i=1}^d [k_i h, (k_i + 1)h[\mid k_1, \dots, k_d \in \mathbb{Z} \right\} .$$

Théorème 3. Soit $(h_n)_{n \in \mathbb{N}}$ une suite de réels strictement positifs, $(A^r(h_n))_{n \in \mathbb{N}}$ la suite de partitions régulières de \mathbb{R}^d associée (donnée par la Définition 3) et $(\hat{f}_{A^r(h_n)})_{n \in \mathbb{N}}$ le classifieur par partition associé. Alors, si

$$h_n \rightarrow 0 \quad \text{et} \quad nh_n^d \rightarrow +\infty$$

quand $n \rightarrow +\infty$, l'algorithme $(\hat{f}_{A^r(h_n)})_{n \in \mathbb{N}}$ est universellement consistant : pour toute loi P des données,

$$\mathbb{E} \left[\mathcal{R}^{0-1} \left(\hat{f}_{A^r(h_n)}(D_n) \right) \right] - \mathcal{R}^{0-1}(f^*) \xrightarrow{n \rightarrow +\infty} 0 . \quad (7)$$

La démonstration du Théorème 3 repose notamment sur le lemme suivant.

Lemme 2. Soit $N \geq 1$, $p \in]0, 1]$ et Z une variable de loi binomiale de paramètres N et p . Alors,

$$\mathbb{P} \left(Z \leq \frac{Np}{2} \right) \leq \frac{4}{Np} . \quad (8)$$

Démonstration du Théorème 3. Il suffit de vérifier que les deux hypothèses du Théorème 2 sont satisfaites quelle que soit la loi P . Comme $\text{diam}(A^r(h_n)) = \sqrt{d}h_n$, la condition sur le diamètre est satisfaite dès lors que $h_n \rightarrow 0$. Pour la deuxième condition, fixons $M \in]0, +\infty[$ et majorons $\mathbb{P}(N_{A^r(h_n)}(X) \geq M)$.

Soit S une boule centrée à l'origine.¹ Elle a une intersection non-vide avec au plus $c_1 + c_2 h_n^{-d}$ cellules de la forme $\prod_{i=1}^d [k_i h_n, (k_i + 1)h_n[$ avec $k_1, \dots, k_d \in \mathbb{Z}$, où $c_1, c_2 > 0$ sont des constantes. On a alors

$$\mathbb{P}(N_{A^r(h_n)}(X) \leq M) \leq \sum_{B \in A^r(h_n) / B \cap S \neq \emptyset} \mathbb{P}(X \in B, N_B \leq M) + \mathbb{P}(X \notin S) .$$

Considérons une cellule $B \in A^r(h_n)$ quelconque telle que $\mathbb{P}(X \in B) > 0$ et majorons $\mathbb{P}(X \in B, N_B \leq M)$. Deux cas sont possibles :

(i) Soit $\mathbb{P}(X \in B) \leq \frac{2M}{n}$, et alors

$$\mathbb{P}(X \in B, N_B \leq M) \leq \frac{2M}{n} .$$

(ii) Soit $\mathbb{P}(X \in B) \geq \frac{2M}{n}$, et alors on va appliquer le Lemme 2, en remarquant que sachant $X \in B$, $Z = N_B$ suit une loi binomiale de paramètres n et $\mathbb{P}(X \in B)$. Ainsi,

$$\begin{aligned} \mathbb{P}(X \in B, N_B \leq M) &= \mathbb{P}(X \in B) \mathbb{P}(N_B \leq M \mid X \in B) \\ &\leq \mathbb{P}(X \in B) \mathbb{P}\left(N_B \leq \frac{n\mathbb{P}(X \in B)}{2} \mid X \in B\right) \\ &\leq \mathbb{P}(X \in B) \frac{4}{n\mathbb{P}(X \in B)} = \frac{4}{n} . \end{aligned}$$

Au final, on a donc

$$\begin{aligned} \mathbb{P}(N_{A^r(h_n)}(X) \leq M) &\leq \text{card}\{B \in A^r(h_n) / B \cap S \neq \emptyset\} \frac{\max\{4, 2M\}}{n} + \mathbb{P}(X \notin S) \\ &\leq (c_1 + c_2 h_n^{-d}) \frac{\max\{4, 2M\}}{n} + \mathbb{P}(X \notin S) \end{aligned}$$

et lorsque n tend vers l'infini, par l'hypothèse $nh_n^d \rightarrow +\infty$, ce majorant tend vers $\mathbb{P}(X \notin S)$ qui peut être rendu arbitrairement petit. Ainsi, pour tout $M > 0$, $\mathbb{P}(N_{A^r(h_n)}(X) \leq M) \rightarrow 0$ lorsque $n \rightarrow +\infty$, c'est-à-dire, $N_{A^r(h_n)}(X) \rightarrow +\infty$ en probabilité. \square

2.4 Partition et minimisation du risque empirique

Il est intéressant de noter que les prédicteur par partition peuvent aussi s'interpréter comme des prédicteurs minimisant le risque empirique (définies au premier cours).

Proposition 3. *Soit A une partition finie ou dénombrable de \mathcal{X} . En régression avec le risque quadratique, le prédicteur par partition associée à A minimise le risque empirique sur l'ensemble de prédicteurs*

$$S_A^r = \{f : \mathcal{X} \mapsto \mathbb{R} \text{ mesurable} / \forall k \geq 1, f \text{ est constante sur } A_k\}$$

des fonctions constantes sur chaque élément de la partition.

De même, en classification binaire avec le risque 0-1, le prédicteur par partition minimise le risque empirique sur l'ensemble de prédicteurs

$$S_A^c = \{f : \mathcal{X} \mapsto \{0, 1\} \text{ mesurable} / \forall k \geq 1, f \text{ est constante sur } A_k\} .$$

1. Cette boule est introduite seulement pour gérer les distributions sur \mathcal{X} avec un support infini.

3 Théorème de Stone (optionnel)

Le théorème de Stone que nous ne démontrerons pas généralise l'argument que nous avons vu pour les algorithmes par partition à un ensemble assez général de méthodes par moyennage local. Il s'applique aux plus proches voisins, aux méthodes d'histogrammes et au prédicteurs de Nadaraya-Watson (avec des conditions sur leurs hyperparamètres).

Théorème 4. (Stone) *Supposons que les poids W_i et la loi des données d'entraînement satisfont*

$$(i) \exists c > 0, \quad \forall f : \mathcal{X} \rightarrow \mathbb{R}_+, \quad \forall n \in \mathbb{N}, \quad \mathbb{E} \left[\sum_{i=1}^n |W_i(X)| f(X_i) \right] \leq c \mathbb{E}[f(X)],$$

$$(ii) \exists D > 0, \quad \forall n \in \mathbb{N}, \quad \sum_{i=1}^n |W_i(X)| \leq D \quad \mathbb{P}\text{-p.s.},$$

$$(iii) \forall a > 0, \quad \mathbb{E} \left[\sum_{i=1}^n |W_i(X)| \mathbf{1}_{\{\|X_i - X\| > a\}} \right] \xrightarrow{n \rightarrow \infty} 0,$$

$$(iv) \sum_{i=1}^n W_i(X) \xrightarrow{\mathbb{P}} 1,$$

$$(v) \mathbb{E} \left[\sum_{i=1}^n |W_i(X)|^2 \right] \xrightarrow{n \rightarrow \infty} 0.$$

Alors $\hat{\eta} : x \mapsto \sum_{i=1}^n W_i(x) Y_i$ donne une loi de classification consistante pour la loi des données d'entraînement avec la méthode de plug-in.

Note : pour que la loi des k plus proches voisins satisfasse les conditions du théorème, on a besoin de faire varier k avec n . Les conditions : $k(n) \xrightarrow{n \rightarrow \infty} \infty$ et $\frac{k(n)}{n} \xrightarrow{n \rightarrow \infty} 0$ sont suffisantes pour avoir la consistance universelle.

Références

- [Devroye et al., 1996] Devroye, L., Györfi, L., and Lugosi, G. (1996). *A probabilistic theory of pattern recognition*, volume 31. Springer Verlag.
- [Hastie et al., 2009] Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The elements of statistical learning*. Springer.