

TP3 : CONSISTANCE UNIVERSELLE, NO FREE LUNCH ET CLASSIFIEUR KNN

1 Consistance universelle uniforme : résultats positifs et négatifs

On considère le problème de classification binaire avec

$$(X_i, Y_i) \stackrel{\text{iid}}{\sim} P, \quad i = 1, \dots, n$$

où P est une loi de probabilité sur $\mathcal{X} \times \mathcal{Y}$ avec $\mathcal{Y} = \{0, 1\}$. La qualité de prédiction de la valeur y par la valeur y' est mesurée

$$\ell(y, y') = \mathbb{1}(y \neq y').$$

Le risque d'une fonction de prédiction $g : \mathcal{X} \rightarrow \mathcal{Y}$ est alors calculé par

$$R_P(g) = \mathbb{E}_P[\ell(Y, g(X))] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, g(x)) dP(x, y).$$

Rappelons que ce risque est minimisé par le classifieur de Bayes défini par $g_P^*(x) = \mathbb{1}(\eta^*(x) > 1/2)$ où $\eta^*(x) = \mathbb{E}_P[Y|X = x]$. Soit $\hat{g}_n : (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{F}(\mathcal{X}, \mathcal{Y})$ un classifieur. On dit qu'il est uniformément universellement consistant en probabilité, si $\forall \delta > 0$,

$$\sup_P P\left(|R_P(\hat{g}_n) - R_P(g_P^*)| > \delta\right) \xrightarrow[n \rightarrow \infty]{P} 0. \quad (1)$$

Dans toute la suite, on utilisera la notation $R_P^* = R_P(g_P^*)$.

1. On suppose d'abord que \mathcal{X} est fini : $\text{Card}(\mathcal{X}) = K$.
 - (a) Quel est le cardinal de $\mathcal{F}(\mathcal{X}, \mathcal{Y})$?
 - (b) Rappeler la borne de risque obtenue en cours pour le minimiseur du risque empirique $\hat{g}_{n, \mathcal{G}}$ pour $\mathcal{G} = \mathcal{F}(\mathcal{X}, \mathcal{Y})$. Peut-on en déduire que $\hat{g}_{n, \mathcal{G}}$ est uniformément universellement consistant ?
 - (c) On suppose maintenant que $K = K_n$ dépend de la taille de l'échantillon. Montrer que si K_n est sous-linéaire en n , alors $\hat{g}_{n, \mathcal{G}}$ est uniformément universellement consistant.
2. On montre maintenant que la consistance uniforme universelle est impossible pour les \mathcal{X} de cardinal infini. Pour simplifier et sans perte de généralité, on suppose que $\mathcal{X} = \mathbb{N}^*$. On va montrer que pour tout $K \in \mathbb{N}^*$, on a

$$\sup_P P\left(|R_P(\hat{g}_n) - R_P(g_P^*)| > \delta\right) \geq \frac{1}{2} \left(1 - \frac{1}{K}\right)^n - \delta. \quad (2)$$

- (a) Soit S un ensemble fini, $\{P_s : s \in S\}$ une famille de probabilités sur $\mathcal{X} \times \mathcal{Y}$ et $w : S \rightarrow \mathbb{R}_+$ tel que $\sum_{s \in S} w(s) = 1$. Montrer que

$$\sup_P P\left(|R_P(\hat{g}_n) - R_P(g_P^*)| > \delta\right) \geq \sum_{s \in S} \mathbb{E}_{P_s}[R_{P_s}(\hat{g}_n) - R_{P_s}^*] w(s) - \delta. \quad (3)$$

- (b) On choisit $S = \{0;1\}^K$, l'ensemble des vecteurs de longueur K dont toutes les coordonnées sont soit 0 soit 1. Pour chaque $s \in S$, on définit la loi P_s de la manière suivante :

$$P_s(X = k, Y = y) = \frac{1}{K} \mathbb{1}(k \leq K) \mathbb{1}(y = s_k).$$

Vérifier que P_s est une probabilité et montrer que $R_{P_s}^* = 0$.

- (c) Montrer que

$$R_{P_s}(g) = \frac{1}{K} \sum_{k=1}^K \mathbb{1}(g(k) \neq s_k), \quad \forall s \in S; \quad \forall g : \mathbb{N}^* \rightarrow \mathcal{Y}.$$

- (d) Toujours pour $S = \{0;1\}^K$, on définit $w : S \rightarrow \mathbb{R}_+$ par $w(s) = 1/2^K$. Soit

$$\hat{g}_n(\cdot) = \hat{g}_n(\cdot, X_1, \dots, X_n, Y_1, \dots, Y_n)$$

un prédicteur quelconque. Montrer que

$$P_s(\hat{g}_n(1, \mathcal{D}_n) \neq s_1) \geq P_X(\hat{g}_n(1, X_1, \dots, X_n, s_{X_1}, \dots, s_{X_n}) \neq s_1; 1 \notin \{X_1, \dots, X_n\}).$$

En déduire que

$$\sum_{s_1 \in \{0;1\}} P_s(\hat{g}_n(1, \mathcal{D}_n) \neq s_1) \geq P_X(1 \notin \{X_1, \dots, X_n\}) = \left(1 - \frac{1}{K}\right)^n.$$

- (e) En combinant les résultats précédents, montrer que

$$\sup_P P\left(|R_P(\hat{g}_n) - R_P(g_P^*)| > \delta\right) \geq \frac{1}{2} \left(1 - \frac{1}{K}\right)^n - \delta, \quad \forall K \in \mathbb{N}^*; \quad \forall \delta > 0.$$

En déduire qu'il n'existe pas de prédicteur universellement uniformément consistant.

2 Consistance de l'algorithme kNN

Le but de cet exercice est de montrer que l'algorithme kNN employé avec $k = 1$ n'est pas consistant. Pour cela, nous considérons le problème de classification binaire avec $\mathcal{X} = [0, 1]$ et $\mathcal{Y} = \{0; 1\}$. On note P_X la loi marginale des X_i et suppose que

$$\eta^*(x) = P(Y_1 = 1 | X_1 = x) \equiv \frac{3}{4}, \quad \forall x \in \mathcal{X}.$$

L'objectif des questions suivantes est de calculer le risque du classifieur oracle g_P^* ainsi que celui du classifieur kNN $\hat{g}_{n,k}$ avec $k = 1$. On verra que ce dernier ne dépend pas de la taille de l'échantillon et est strictement plus grand que le risque de l'oracle.

1. Montrer que pour tout application déterministe $g : \mathcal{X} \rightarrow \{0; 1\}$, on a

$$R_P(g) = \mathbf{E}_{P_X}[\eta^*(X)] + \mathbf{E}_{P_X}[g(X)(1 - 2\eta^*(X))].$$

2. En déduire que si $\eta^* \equiv 3/4$, alors le classifieur oracle (appelé aussi classifieur de Bayes) est donné par $g_P^* \equiv 1$ et son risque vaut $R_P(g_P^*) = 1/4$.
3. Montrer que pour tout $g : \mathcal{X} \rightarrow \{0; 1\}$, on a

$$R_P(g) = \frac{3}{4} - \frac{1}{2} \int_{\mathcal{X}} g(x) P_X(dx).$$

4. Soit $\mathcal{D}_n = \{(X_i, Y_i); i = 1, \dots, n\}$ et $\hat{g}_{n,1}(x) = \hat{g}_{\text{PPV}}(x, \mathcal{D}_n)$ le classifieur du plus proche voisin (PPV). Fixons $x \in \mathcal{X}$ et cherchons à calculer $\mathbf{E}_P[\hat{g}_{\text{PPV}}(x, \mathcal{D}_n)]$, où l'espérance est par rapport à l'échantillon \mathcal{D}_n . Pour tout $i = 1, \dots, n$, posons

$$Z_i = \begin{cases} 1, & \text{si } X_i \text{ est le PPV de } x \\ 0, & \text{sinon.} \end{cases}$$

Montrer que

$$\mathbf{E}_P[\hat{g}_{\text{PPV}}(x, \mathcal{D}_n)] = \sum_{i=1}^n \mathbf{E}_P[Y_i Z_i]. \quad (4)$$

5. Vérifier que pour tout i , Y_i est indépendant de (X_1, \dots, X_n) . En déduire que Y_i et Z_i sont indépendantes.
6. En utilisant la question précédente et la relation évidente $\sum_{i=1}^n Z_i = 1$ montrer que

$$\mathbf{E}_P[R_P(\hat{g}_{\text{PPV}})] = \frac{3}{8}.$$

Conclure.

7. Considérer le cas des 3 plus proches voisins $\hat{g}_{3-\text{PPV}}$. Montrer que son risque moyen $\mathbf{E}_P[R_P(\hat{g}_{3-\text{PPV}})]$ est égal à 21/64.
8. Passons maintenant au cas général d'un prédicteur kNN $\hat{g}_{k-\text{PPV}}$. Soient V_1, \dots, V_k des variables aléatoires i.i.d. de loi de Bernoulli de paramètre 3/4. Montrer que

$$\mathbf{E}_P[\hat{g}_{k-\text{PPV}}(x, \mathcal{D}_n)] = \mathbf{P}(\bar{V}_k > 1/2).$$

En déduire que cette espérance tend vers 1 lorsque $k \rightarrow \infty$ et, par conséquent, le risque espéré $\mathbf{E}_P[R_P(\hat{g}_{k-\text{PPV}})]$ tend vers le risque de l'oracle, c'est à dire vers 1/4.