

Bornes PAC-Bayésiennes, classification et clustering en grande dimension

Olivier Catoni

CNRS, INRIA (projet CLASSIC)

Département de Mathématiques et Applications,

ENS, 45 rue d'Ulm, 75 230 Paris Cedex 05,

`Olivier.Catoni@ens.fr`

JOURNÉE THÉMATIQUE IMAGE ET APPRENTISSAGE

Séminaire Méthodes Mathématiques du Traitement d'Images

Laboratoire Jacques Louis Lions,

Université Pierre et Marie Curie & CNRS

jeudi 5 juillet 2012

On considère :

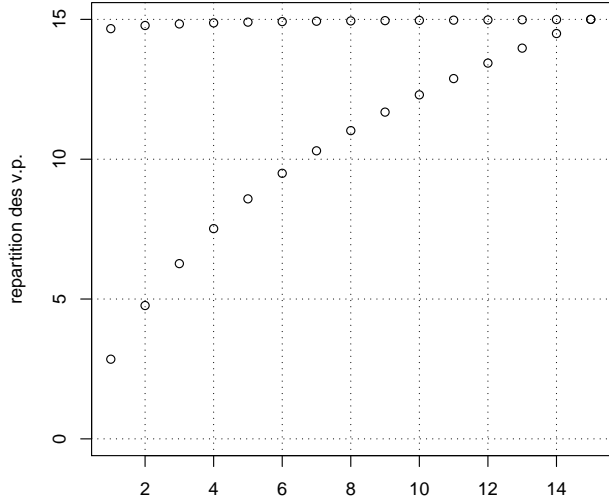
- des données (par exemple des images) $(X_1, \dots, X_n) \in \mathcal{F}$, i.i.d. de loi \mathbb{P} , à valeurs dans un *espace de formes* \mathcal{F} ;
- un noyau symétrique positif $k : \mathcal{F} \times \mathcal{F} \rightarrow \mathbb{R}_+$, normalisé, $\left(k(x, y) = k(y, x), \quad \sum_{i,j} \alpha_i k(x_i, x_j) \alpha_j \geq 0, \quad k(x, x) = 1 \right)$;
- pour tout $x \in \mathcal{F}$ la fonction $K_x = y \mapsto k(x, y) : \mathcal{F} \rightarrow \mathbb{R}$;
- le produit scalaire $\langle K_x, K_y \rangle = k(x, y)$;
- l'espace de Hilbert \mathcal{H} complété de $\mathbf{Vect}\{K_x, x \in \mathcal{F}\}$ pour ce produit scalaire ;
- le plongement $\Psi = x \mapsto K_x : \mathcal{F} \rightarrow \mathcal{H}$ (construction de Moore-Aronszajn).

Exemple type : $\mathcal{F} \subset \mathbb{R}^d$, $k(x, y) = \exp(-\lambda \|x - y\|^2)$.

Effet du choix du noyau :

- pour la classification : La règle linéaire dans \mathcal{H} (introduite par Vapnik sous le nom de Support Vector Machine) $\{x \in \mathcal{F} : \langle x, \theta \rangle > 0\}$, où $\theta \in \mathcal{H}$, a une frontière dont la régularité dépend du choix de k ;
- pour le clustering : le choix de k influence la décroissance des valeurs propres de l'opérateur $y \mapsto \int \langle y, x \rangle x d\mathbb{P}(x) : \mathcal{H} \rightarrow \mathcal{H}$, dont les vecteurs propres v_k définissent une partition de \mathcal{F} suivant la valeur de $\arg \max_k |\langle x, v_k \rangle|$.





Enjeux statistiques

Classification supervisée : on observe un échantillon étiqueté, (X_i, Y_i) , où $Y_i \in \{-1, +1\}$. On souhaite relier l'erreur empirique

$$\int \mathbf{1}(\langle x, \theta \rangle y \leq 0) d\bar{\mathbb{P}}(x, y)$$

où $\bar{\mathbb{P}} = \frac{1}{n} \sum_{i=1}^n \delta_{(X_i, Y_i)}$, à son espérance $\int \mathbf{1}(\langle x, \theta \rangle y \leq 0) d\mathbb{P}(x, y)$.

Clustering : on souhaite relier l'énergie empirique dans la direction $\theta \in \mathcal{H}$

$$e(\theta) = \int \langle \theta, x \rangle^2 d\bar{\mathbb{P}}(x), \text{ et son espérance } \mathcal{E}(\theta) = \int \langle \theta, x \rangle^2 d\mathbb{P}(x).$$

Inégalités PAC-Bayésiennes pour la classification

d'après Langford, Shawe-Taylor, McAllester, C.

Technique de perturbation du paramètre θ : on suppose pour le moment que $\mathcal{H} = \mathbb{R}^d$. On introduit la perturbation Gaussienne

$$\frac{d\mu_\theta}{d\theta'}(\theta') = \left(\frac{\beta}{2\pi}\right)^{d/2} \exp\left(-\frac{\beta\|\theta' - \theta\|^2}{2}\right), \quad \theta \in \mathbb{R}^d.$$

Elle vérifie $\mathcal{K}(\mu_\theta, \mu_0) = \int \log(d\mu_\theta/d\mu_0)d\mu_\theta = \frac{\beta}{2}\|\theta\|^2$.

Fonction d'erreur perturbée :

$$\int \mathbf{1}(\langle x, \theta' \rangle y \leq 0) d\mu_{\theta}(\theta') = \varphi\left(\frac{\sqrt{\beta} \langle x, \theta \rangle y}{\|x\|}\right),$$

où $\varphi(a) = \int_a^{+\infty} (2\pi)^{-1/2} \exp(-t^2/2) dt$

$(\leq \frac{1}{2} \exp(-a^2/2), \text{ pour tout } a \geq 0).$

On prend un peu de marge :

$$\begin{aligned} & \int \mathbf{1} \left(\|x\|^{-1} \langle x, \theta' \rangle y \leq 1 \right) d\mu_{\theta}(\theta') \\ &= \varphi \left[\sqrt{\beta} \left(\|x\|^{-1} \langle x, \theta \rangle y - 1 \right) \right] \geq \varphi(-\sqrt{\beta}) \mathbf{1} \left(\langle x, \theta \rangle y \leq 0 \right) \\ & \qquad \qquad \qquad \left[1 - \varphi(\sqrt{\beta}) \right] \mathbf{1} \left(\langle x, \theta \rangle y \leq 0 \right), \end{aligned}$$

pour aboutir à

$$\begin{aligned} & \int \mathbf{1} \left(\langle x, \theta \rangle y \leq 0 \right) d\mathbb{P}(x, y) \leq \left[1 - \varphi(\sqrt{\beta}) \right]^{-1} \\ & \qquad \qquad \qquad \times \underbrace{\int \mathbf{1} \left(\|x\|^{-1} \langle x, \theta' \rangle y \leq 1 \right) d\mu_{\theta}(\theta')}_{\text{erreur avec marge}} \underbrace{d\mathbb{P}(x, y)}_{\text{perturbation}}. \end{aligned}$$

Inégalité PAC-Bayésienne :

$$\exp\left(\int h(\theta')\mu_\theta(\theta') - \mathcal{K}(\mu_\theta, \mu_0)\right) \leq \int \exp[h(\theta')]d\mu_0(\theta').$$

Donne

$$\int \exp\left\{ \sup_{\theta \in \mathbb{R}^d} n\lambda \int \left[\Phi_\lambda\left(\int \mathbf{1}[\|x\|^{-1}\langle x, \theta' \rangle y \leq 1]d\mathbb{P}(x, y)\right) - \int \mathbf{1}(\|x\|^{-1}\langle x, \theta' \rangle y \leq 1)d\bar{\mathbb{P}}(x, y)\right] d\mu_\theta(\theta') - \mathcal{K}(\mu_\theta, \mu_0) \right\} d\mathbb{P}^{\otimes n} \leq 1,$$

où $\Phi_\lambda(p) = -\lambda^{-1} \log[1 - p + p \exp(-\lambda)]$.

Avec probabilité au moins $1 - \epsilon$, pour tout $\theta \in \mathbb{R}^d$,

$$\begin{aligned}
 & \int \mathbf{1}(\langle x, \theta \rangle y \leq 0) d\mathbb{P}(x, y) \leq [1 - \varphi(\sqrt{\beta})]^{-1} \\
 & \quad \times \inf_{\lambda \in \Lambda} \Phi_{\lambda}^{-1} \left\{ \int \varphi[\sqrt{\beta}(\|x\|^{-1} \langle x, \theta \rangle y - 1)] d\bar{\mathbb{P}}(x, y) \right. \\
 & \quad \quad \quad \left. + \frac{\beta \|\theta\|^2 / 2 + \log(|\Lambda|/\epsilon)}{n\lambda} \right\} \\
 & \quad \leq [1 - \varphi(\sqrt{\beta})]^{-1} \\
 & \quad \times \inf_{\lambda \in \Lambda} \Phi_{\lambda}^{-1} \left\{ \int \left(2 - \|x\|^{-1} \langle \theta, x \rangle y \right)_+ d\bar{\mathbb{P}}(x, y) + \frac{\beta \|\theta\|^2}{2n\lambda} \right. \\
 & \quad \quad \quad \left. + \frac{\log(|\Lambda|/\epsilon)}{n\lambda} + \varphi(\sqrt{\beta}) \right\},
 \end{aligned}$$

indépendamment de la dimension !

Estimation de l'énergie

en cours ...

Résultat dépendant de la dimension :

Considérons l'approximation suivante de l'identité

$$\psi(a) = \begin{cases} -\log(1 - a + a^2/2), & 0 \leq a \leq 1, \\ \log(2), & a \geq 1, \\ -\psi(-a), & a \leq 0. \end{cases}$$

Soit $e(\theta)$ la solution de

$$\int \psi\left(\lambda \left[\langle x, \theta \rangle^2 e(\theta)^{-1} - 1\right]\right) d\bar{\mathbb{P}}(x) = 0.$$

En perturbant θ comme précédemment et en utilisant le changement de variable $G^{1/2}\theta, G^{-1/2}x$ (qui laisse le produit scalaire invariant), on obtient avec probabilité au moins $1 - \epsilon$

$$\sup_{\theta \in \mathbb{R}^d \setminus \{0\}} |\mathcal{E}(\theta)/e(\theta) - 1| \leq 2(\eta + \gamma),$$

où

$$\eta = \sqrt{\frac{2(\kappa^2 - 1)}{n} \left[\log(2/\epsilon) + \frac{10s_4^2}{9\kappa} \right]},$$

$$\gamma = 6\sqrt{\frac{5s_4^2\kappa}{n}},$$

$$s_4^2 = \sqrt{\int \|G^{-1/2}x\|^4 d\mathbb{P}(x)} \geq d = \int \|G^{-1/2}x\|^2 d\mathbb{P}(x),$$

$$\kappa = \sup \left\{ \sqrt{\int \langle x, \theta \rangle^4 d\mathbb{P}(x)}, \theta \in \mathbb{R}^d, \int \langle x, \theta \rangle^2 d\mathbb{P}(x) \leq 1 \right\}.$$