

Apprentissage statistique et grande dimension

M. Hoffmann. M1 Paris Dauphine 2014-2015

Chapitre 1 : Expérience statistique et grande dimension (Séances 1-2)

1. Exemples de modèles statistiques non-paramétriques
 - Estimation d'une densité de probabilité (cadre paramétrique, non-paramétrique)
 - Régression non-paramétrique (cadre paramétrique, cadre grande dimension)
 - Exemples de problèmes statistiques en grande dimension.
2. L'heuristique du modèle de suite gaussienne, notion d'expériences statistiques (asymptotiquement) équivalentes.
 - Le cas du modèle de densité
 - La cas du modèle de régression
 - Définition du modèle de suite gaussienne

Chapitre 2 : Estimateurs linéaires et régression non-paramétrique (Séances 2-3-4)

1. Situation : modèle de régression, définition : on observe

$$Y_i = f(X_i) + \xi_i, \dots, i = 1, \dots, n$$

où $f : [0, 1] \rightarrow \mathbb{R}$ est la fonction inconnue, les X_i sont le design i.i.d. et les ξ_i sont des bruits i.i.d. de variance finie et centrés.

Notation vectorielle : $\mathbf{y} = f + \boldsymbol{\xi}$. Lorsqu'il n'y a pas (trop) d'ambiguïté, on identifie la fonction f et le vecteur colonne $(f(X_1) \cdots f(X_n))^T$. Notion de dictionnaire de taille M et modèle paramétrique associé. Définition de l'estimateur des moindres carrés non-paramétrique. L'EMC non-paramétrique est un estimateur linéaire si $\mathbf{X}^T \mathbf{X} > 0$. Le terme $\boldsymbol{\xi}$ est un bruit blanc fort de variance σ^2 .

2. Propriétés statistiques des estimateurs linéaires en régression. Matrice de lissage S . Risque quadratique

$$\mathcal{R}(\hat{f}, f) = \mathbb{E}[\|\hat{f} - f\|_n^2]$$

où, pour $g : [0, 1] \rightarrow \mathbb{R}$, on note $\|g\|_n^2 = n^{-1} \sum_{i=1}^n g(X_i)^2$ la semi-norme L^2 -empirique associée au design $\{X_i, i = 1, \dots, n\}$. Décomposition biais-variance du risque quadratique des estimateurs linéaires : si $\hat{f} = S\mathbf{y}$, alors

$$\mathcal{R}(\hat{f}, f) = \mathbb{E}[\|Sf - f\|_n^2] + \frac{\sigma^2}{n} \text{Trace}(S^T S).$$

Si l'on travaille conditionnellement au design $\{X_i = x_i, i = 1, \dots, n\}$ (par exemple pour le design déterministe uniforme $X_i = i/n$) alors la décomposition précédente devient

$$\mathcal{R}(\hat{f}, f) = \|Sf - f\|_n^2 + \frac{\sigma^2}{n} \text{Trace}(S^T S).$$

Corollaire : si \hat{f}^{MC} désigne l'EMC non-paramétrique, alors, conditionnellement au design $\{X_i = x_i, i = 1, \dots, n\}$, on a

$$\mathcal{R}(\hat{f}^{\text{MC}}, f) \leq \min_{\theta \in \mathbb{R}^M} \|f_\theta - f\|_n^2 + \frac{\sigma^2}{n} M \wedge n, \quad (\star)$$

où f_θ est la paramétrisation de f via un dictionnaire fini.

3. Vitesses de convergence de l'EMC non-paramétrique. Définition des espaces (boules) de Sobolev $W(\beta, L)$ périodiques en dimension 1, traduction séquentielle via la base de Fourier. Si $f \in W(\beta, L)$, alors

$$\inf_{\theta \in \mathbb{R}^M} \mathbb{E}[\|f_\theta - f\|_n^2] \leq L^2 M^{-2\beta}$$

pour un design aléatoire uniforme. Extension pour le risque conditionnel au design déterministe uniforme $X_i = i/n$. Bornes supérieures d'estimation : pour un choix de taille de dictionnaire $M_n \approx n^{1/(2\beta+1)}$ dans la base de Fourier, l'EMC non-paramétrique atteint la vitesse $n^{-2\beta/(2\beta+1)}$ pour le risque \mathcal{R} . Optimalité et notion de vitesse minimax.

Chapitre 3 : Représentation parcimonieuse (sparse) en régression (Séances 5-6)

1. Situation. Ecriture séquentielle de (\star) :

$$\mathbb{E}[n^{-1} \|X(\theta^{\text{MC}} - \theta^*)\|^2] \leq \sigma^2 \frac{M \wedge n}{n}$$

si $f = f_{\theta^*}$ est représentés sans erreur via un dictionnaire de taille M . Problématique M "grand". Notion de vecteur ε -sparse (le nombre $M(\theta)$ de composantes non-nulles est majoré par εM). Recherche d'un estimateur $\hat{\theta}$ de sorte que

$$\mathbb{E}[n^{-1} \|X(\tilde{\theta} - \theta^*)\|^2] \leq \sigma^2 \frac{M(\theta^*)}{n} \log M. \quad (\star\star)$$

2. Seuillage dans le modèle de suite gaussienne. Hypothèse restrictive ORT: $n^{-1} \mathbf{X}^T \mathbf{X} = \text{Id}_{\mathbb{R}^n}$. Lemme de concentration gaussienne :

$$\mathbb{P}\left(\max_{1 \leq j \leq M} |\eta_j| > t\sqrt{\log M}\right) \leq M^{1-t^2/2}$$

si $M \geq 2$, $t \geq \sqrt{2}$ et les η_j sont des gaussiennes standard. Notion de “hard-thresholding” et interprétation (à la Hodge-Lehmann). Théorème : l’estimateur par seuillage réalise (dans le modèle de suite gaussienne en ORT) le programme $(\star\star)$ en espérance et en probabilité. Il peut de plus estimer l’ensemble d’indices révélateurs $J(\theta) = \{j, \theta_j \neq 0\}$ sous une condition apparentée à la séparation d’hypothèses en théorie des tests.

3. Remarques et développements : soft-thresholding, formulations variationnelles exactes dans le cadre ORT. Au delà, lien entre “hard-thresholding” et estimateur BIC, et “soft-thresholding” et estimateur LASSO. Extension du cadre ORT : hypothèse de cohérence.

Chapitre 4 : Classification et apprentissage statistique (séances 7-8-9)

1. Exemples de problèmes de classification
2. Formulation non-paramétrique de la classification (binaire): notion de classifieur $h : \mathcal{X} \rightarrow \{0, 1\}$ où \mathcal{X} désigne l’espace des covariables, erreur/risque de classification $R(h)$.
3. Le classifieur bayésien. Le classifieur $h^*(x) = \mathbf{1}_{\{\eta(x) > 1/2\}}$ minimise l’erreur de classification, où $\eta(x) = \mathbb{P}(h(X) \neq Y)$. Pour tout classifieur h , on a

$$R(h) - R(h^*) = \int_{x, h(x) \neq h^*(x)} |2\eta(x) - 1| \mathbb{P}_X(dx)$$

où \mathbb{P}_X désigne la loi du design. Limites de l’estimation directe du classifieur bayésien par estimation de la fonction de régression η .

4. Apprentissage et minimisation du risque empirique. Notion d’excès de risque.

- Minimisation du risque empirique $R_n(h) = n^{-1} \sum_{i=1}^n \mathbf{1}_{\{Y_i \neq h(X_i)\}}$ où

$$\mathcal{D}_n = \{(X_i, Y_i), i = 1, \dots, n\}$$

est un échantillon d’apprentissage. Pour un dictionnaire \mathcal{H} défini *a priori*, le classifieur minimisant le risque empirique est défini comme $\widehat{h}_n^{\text{erm}}$ vérifiant $R_n(\widehat{h}_n^{\text{erm}}) = \min_{h \in \mathcal{H}} R_n(h)$. Notion de sur-apprentissage et de décomposition erreur stochastique + erreur d’approximation.

- Contrôle de l'erreur stochastique

$$R_{\mathcal{D}_n}(\widehat{h}_n^{\text{erm}}) - \inf_{h \in \mathcal{H}} R(h) \leq 2 \sup_{h \in \mathcal{H}} |R_n(h) - R(h)|$$

où $R_{\mathcal{D}_n}(h) = \mathbb{P}(h(X, \mathcal{D}_n) \neq Y | \mathcal{D}_n)$ est l'erreur de classification conditionnelle à l'échantillon d'apprentissage pour un classifieur $h(x) = h(x, \mathcal{D}_n)$ construit à l'aide d'un échantillon d'apprentissage \mathcal{D}_n .

- Classification pour un dictionnaire fini. Si \mathcal{H} est fini et de taille M , alors

$$R_{\mathcal{D}_n}(\widehat{h}_n^{\text{erm}}) \leq \min_{h \in \mathcal{H}} R(h) + \sqrt{\frac{2}{n} \log \frac{2M}{\delta}}$$

avec probabilité plus grande que $1 - \delta$, et

$$\mathbb{E}[R_{\mathcal{D}_n}(\widehat{h}_n^{\text{erm}})] \leq \min_{h \in \mathcal{H}} R(h) + \sqrt{\frac{\log 2M}{2n}}.$$

Preuves via Chernoff.

Chapitre 5. Introduction à la théorie de Vapnik-Chervonenkis (séance 10-11)

1. Si $\mu_n = n^{-1} \sum_{i=1}^n \delta_{X_i}$ et $\mu = \mathbb{P}(X \in \cdot)$, alors, en identifiant un classifieur $h = \mathbf{1}_A$ avec un ensemble mesurable $A \subset \mathcal{X}$, on a

$$\sup_{h \in \mathcal{X}} |R_n(h) - R(h)| = \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)|.$$

Symétrisation du risque :

$$\sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| = \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu'_n(A)|$$

où μ'_n est obtenue via une copie indépendante de (X_1, \dots, X_n) .

2. Complexité de Rademacher : pour $B \subset \mathbb{R}^n$ borné, elle est définie comme

$$\mathcal{R}_n(B) = \mathbb{E} \left[\sup_{(b_1, \dots, b_n) \in B} \left| n^{-1} \sum_{i=1}^n \sigma_i b_i \right| \right],$$

où les σ_i sont des variables de Rademacher indépendantes :

$$\mathbb{P}(\sigma_i = 1) = 1 - \mathbb{P}(\sigma_i = -1) = \frac{1}{2}.$$

Notion d'emprunte binaire d'un vecteur $(x_1, \dots, x_n) \in \mathcal{X}^n$ sur une classe \mathcal{A} d'ensembles mesurables de \mathcal{X} :

$$\mathcal{A}(x_1, \dots, x_n) = \{(b_1, \dots, b_n) \in \{0, 1\}^n, b_i = \mathbf{1}_{\{x_i \in A\}}, A \in \mathcal{A}\}.$$

Alors

$$\begin{aligned} \mathbb{E} \left[\sup_{A \in \mathcal{A}} |\mu_n(A) - \mu'_n(A)| \right] &\leq 2\mathbb{E} \left[\sup_{A \in \mathcal{A}} n^{-1} \left| \sum_{i=1}^n \sigma_i \mathbf{1}_{\{X_i \in A\}} \right| \right] \\ &\leq 2\mathbb{E} [\mathcal{R}_n(\mathcal{A}(X_1, \dots, X_n))]. \end{aligned}$$

3. Contrôle de la complexité de Rademacher. Si B est un sous-ensemble fini de \mathbb{R}^n , alors

$$\mathcal{R}_n(B) \leq \max_{b \in B} \|b\|_{\ell^2} \frac{\sqrt{2 \log(2 \text{Card} B)}}{n}$$

et conclusion :

$$\mathbb{E} \left[\sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| \right] \leq 2\sqrt{\frac{2 \log(2 \mathbb{S}_{\mathcal{A}}(n))}{n}},$$

où $\mathbb{S}_{\mathcal{A}}(n) = \max_{(x_1, \dots, x_n) \in \mathbb{R}^n} \text{Card} \mathcal{A}(x_1, \dots, x_n)$ est le coefficient d'éclatement de la classe \mathcal{A} . Ce résultat ne dépend plus de $\mathcal{L}(X, Y)$.

4. Dimension de Vapnik-Chervonenkis, lemme de Sauer, exemples, interprétations et limites du résultat.

Chapitre 6 : Apprentissage et convexification du risque (séance 11-12)

1. On réécrit $Y \in \{0, 1\}$ sous la forme $Y' \in \{-1, 1\}$ via $Y = 2Y' - 1$. Sans changer de notation dans la suite, si $h : \mathcal{X} \rightarrow \{-1, 1\}$ est un classifieur,

$$\mathbb{P}(h(X) \neq Y) = \mathbb{E}[\varphi_0(-Yh(X))]$$

où $\varphi_0(z) = \mathbf{1}_{\{z \geq 0\}}$. Notion de substitut convexe : toute fonction φ t.q. $\varphi_0 \leq \varphi$. Hinger loss, fonction exponentielle; logistique.

2. Principe de convexification: si \mathcal{F} est un ensemble convexe de fonctions de $\mathcal{X} \rightarrow \mathbb{R}$, on cherche

$$\hat{f}_{n,\varphi} = \operatorname{argmin}_{f \in \mathcal{F}} R_{n,\varphi}(f),$$

où

$$R_{n,\varphi}(f) = n^{-1} \sum_{i=1}^n \varphi(-Y_i f(X_i))$$

et on obtient un classifieur en posant $\hat{h}_{n,\varphi}(x) = \operatorname{sign}(\hat{f}_{n,\varphi}(x))$. Choix de \mathcal{F} : enveloppe convexe ou boule ℓ^1 construite à partir d'un dictionnaire fini de classifieurs.

3. Minimisation du φ -risque empirique. Si φ est strictement convexe et différentiable, alors $f^* = \operatorname{argmin}_f R_\varphi(f)$ est bien défini et $\operatorname{sign}(f^*)$ est le classifieur bayésien.
4. Inégalité de Zhang : pour toute fonction $f : \mathcal{X} \rightarrow \mathbb{R}$ mesurable, il existe deux constantes $c > 0$ et $0 < \gamma \leq 1$ telles que

$$R(\operatorname{sign}(f)) - R^* \leq 2c(R_\varphi(f) - R_\varphi^*)^\gamma,$$

dès lors que $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ est convexe et vérifie $\varphi(0) = 1$, $\varphi(z) \geq \varphi(-z)$ si $z \geq 0$ et $\tau(\eta) = \inf_\alpha (\eta\varphi(-\alpha) + (1 - \eta)\varphi(\alpha))$ vérifie

$$|\tfrac{1}{2} - \eta| \leq c(1 + \tau(\eta))^\gamma.$$

5. Inégalité oracle pour le φ -risque. Si φ est L -Lipischitz , $\varphi(0) = 1$ et \mathcal{F} est une boule ℓ^1 construite sur un dictionnaire de classifieurs de taille M , alors, sous les hypothèses de l'inégalité de Zhang,

$$\mathbb{E}[R_\varphi(\hat{f}_{n,\varphi})] - \inf_{f \in \mathcal{F}} R_\varphi(f) \leq 8L\sqrt{\frac{2\log(2M)}{n}}$$

et

$$\mathbb{E}[R(\hat{h}_{n,\varphi})] - R^* \leq 2c\left(8L\sqrt{\frac{2\log(2M)}{n}}\right)^\gamma + 2c\left(\inf_{f \in \mathcal{F}} R_\varphi(f) - R_\varphi^*\right)^\gamma.$$

Preuve à partir de l'inégalité de Zhang et du principe de contraction dû à Ledoux et Talagrand (admis).

Chapitre 7 : Introduction aux SVM (séance 13)