

Statistiques en grande dimension

Cours 5

Previously...

- X_1, \dots, X_n iid sur $[0, 1]$ de densité f
- $x_0 \in]0, 1[$, on cherche à estimer $f(x_0)$
- Pour un $h > 0$, on pose :

$$\hat{f}_{n,h}^{\text{LC}}(x_0) = \frac{1}{2nh} \sum_{i=1}^n \mathbf{1}_{[x_0-h, x_0+h]}(X_i)$$

- Pour choisir h de façon adaptative dans une famille finie, $\{h_1, \dots, h_M\}$, on utilise la méthode «Intersection of Confidence Intervals» $ICI \rightarrow \hat{f}^{\text{ICI}}(x_0)$

Théorème 1. Si f est bornée par L et si $h_1 < h_2 < \dots < h_M$, alors :

$$\mathbb{P}(|\hat{f}^{\text{ICI}}(x_0) - f(x_0)| \leq 6\rho_{j^*}) \geq 1 - \varepsilon$$

...où :

$$\rho_j = \sqrt{\frac{3L}{2nh_j} \ln\left(\frac{2M}{\varepsilon}\right)}$$

Démonstration. On a noté $\xi_j(x_0) = \hat{f}_{n,h_j}^{\text{LC}}(x_0) - \mathbb{E}[\hat{f}_{n,h_j}^{\text{LC}}(x_0)]$ l'erreur stochastique de $\hat{f}_{n,h_j}^{\text{LC}}$.

Soit $\Omega_0 = \{\forall j = 1, \dots, M, |\xi_j(x_0)| < \rho_j\}$, on a prouvé que :

$$\mathbb{P}(|\hat{f}^{\text{ICI}}(x_0) - f(x_0)| \leq 6\rho_{j^*}) \geq \mathbb{P}(\Omega_0)$$

Il reste à prouver que $\mathbb{P}(\Omega_0) \geq 1 - \varepsilon \Leftrightarrow \mathbb{P}(\overline{\Omega_0}) < \varepsilon$. On a :

$$\begin{aligned} \mathbb{P}(\overline{\Omega_0}) &= \mathbb{P}(\exists j \in [1, M] : |\xi_j(x_0)| > \rho_j) \\ &= \mathbb{P}\left(\bigcup_{j=1}^M \{|\xi_j(x_0)| > \rho_j\}\right) \\ &\leq \mathbb{P}(|\xi_0(x_0)| > \rho_0) + \dots + \mathbb{P}(|\xi_M(x_0)| > \rho_M) \\ &\leq M \max_{1 \leq j \leq M} \mathbb{P}(|\xi_j(x_0)| > \rho_j) \end{aligned}$$

De plus, on peut écrire :

$$\xi_j(x_0) = \frac{1}{2nh_j} \sum_{i=1}^n \left[\underbrace{\mathbf{1}_{[x_0-h_j, x_0+h_j]}(X_i)}_{\in [0,1]} - \underbrace{\int_{x_0-h_j}^{x_0+h_j} f}_{\in [0,1]} \right]$$

Posons $Z_i = \frac{1}{2nh} \left(\mathbf{1}_{[x_0-h, x_0+h]}(X_i) - \int_{x_0-h}^{x_0+h} f \right)$: les Z_i sont iid, $\mathbb{E}[Z_i] = 0$ et $|Z_i| \leq \frac{1}{2nh}$.

On veut majorer $\mathbb{P}(|Z_1 + \dots + Z_n| > \rho)$, et on va pouvoir le faire grâce à l'inégalité de Bernstein :

Théorème 2. (*Inégalité de Bernstein*)

Si on dispose de variables indépendantes Z_i , de moyenne nulle, et bornées presque sûrement par une constante $c > 0$, alors :

$$\forall t > 0, \mathbb{P}\left(\left|\sum_{i=1}^n Z_i\right| \geq t\right) \leq 2 \exp\left\{-\frac{t^2}{2 \sum_{i=1}^n \text{Var}(Z_i) + \frac{2}{3} c t}\right\}$$

Cette inégalité est plus fine que celle de Hoeffding (on remarquera que dans Hoeffding la variance n'intervient pas, uniquement la constante c). On utilise donc cette inégalité :

$$\mathbb{P}\left(\left|\sum_{i=1}^n Z_i\right| > \rho\right) \leq 2 \exp\left\{-\frac{\rho^2}{2 n \text{Var}(Z_1) + \frac{2}{3} \rho \frac{1}{2 n h}}\right\}$$

Rappelons par ailleurs que :

$$\rho = \sqrt{\frac{3 L}{2 n h} \ln\left(\frac{2 M}{\varepsilon}\right)}$$

$$\begin{aligned} \text{Var}(Z_1) &= \mathbb{E}[Z_1]^2 \\ &= \frac{1}{4 n^2 h^2} \text{Var}[\mathbf{1}_{[x_0-h, x_0+h]}(X_1)] \\ &\leq \frac{1}{4 n^2 h^2} \mathbb{E}[\mathbf{1}_{[x_0-h, x_0+h]}(X_1)] \\ &= \frac{1}{4 n^2 h^2} \int_{x_0-h}^{x_0+h} f \\ &\leq \frac{2 h L}{4 n^2 h^2} \end{aligned}$$

D'où,

$$\begin{aligned} \frac{\rho^2}{2 n \text{Var}(Z_1) + \frac{\rho}{3 n h}} &\geq \frac{\frac{3 L}{2 n h} \ln\left(\frac{2 M}{\varepsilon}\right)}{\frac{L}{n h} + \frac{1}{3 n h} \sqrt{\frac{3 L}{2 n h} \ln\left(\frac{2 M}{\varepsilon}\right)}} \\ &\geq \ln\left(\frac{2 M}{\varepsilon}\right) \end{aligned} \tag{1}$$

...car $L = \max_{x \in [0,1]} f(x) \geq \int_0^1 f(x) dx = 1 \Rightarrow \sqrt{L} \leq L$. En poursuivant, on obtient :

$$\frac{\rho^2}{2 n \text{Var}(Z_1) + \frac{\rho}{3 n h}} \geq \frac{3}{2 + \frac{2}{3} \sqrt{\frac{3}{2 h n} \ln\left(\frac{2 M}{\varepsilon}\right)}} \ln\left(\frac{2 M}{\varepsilon}\right)$$

On peut choisir h tel que :

$$\frac{1}{n} \ln\left(\frac{2 M}{\varepsilon}\right) \leq h$$

D'où,

$$\frac{\rho^2}{2n \operatorname{Var}(Z_1) + \frac{\rho}{3nh}} \geq \frac{3}{2 + \sqrt{\frac{2}{3}}} \ln\left(\frac{2M}{\varepsilon}\right) \geq \ln\left(\frac{2M}{\varepsilon}\right)$$

On déduit de (1) que :

$$\begin{aligned} \mathbb{P}(|Z_1 + \dots + Z_n| > \rho) &\leq 2 \exp\left\{-\ln\left(\frac{2M}{\varepsilon}\right)\right\} \\ &= \frac{\varepsilon}{M} \end{aligned}$$

$\mathbb{P}(\overline{\Omega}_0) \leq M \times \frac{\varepsilon}{M} = \varepsilon \Rightarrow \mathbb{P}(\Omega_0) \geq \varepsilon$. D'où le résultat. \square

Régression non-paramétrique

On observe $(X_1, Y_1), \dots, (X_n, Y_n)$, et $Y_i = f(X_i) + \xi_i$ pour $i = 1, \dots, n$.

On note f la fonction de régression et ξ_i les erreurs (également appelées les résidus). Ces dernières vérifient les conditions : $\mathbb{E}[\xi_i | X_i] = 0$ et $\operatorname{Var}[\xi_i] = \sigma^2 < +\infty$ pour tout i .

Cas particulier : *Régression à design fixe* : X_1, \dots, X_n sont déterministes : cela correspond le plus souvent à des problèmes de traitement d'image ou de signal.

Les ξ_i sont iid, $\mathbb{E}[\xi_i] = 0$ et $\operatorname{Var}[\xi_i] = \sigma^2$.

- On cherche à estimer la fonction f
- On va mesurer la qualité d'un estimateur \hat{f} de f
 - Soit en utilisant le *risque empirique* (typiquement pour le débruitage) :

$$\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n (\hat{f}(X_i) - f(X_i))^2\right] = \mathbb{E}[\|\hat{f} - f\|_n^2]$$

- Soit en utilisant le *risque intégré* MISE (typiquement pour la prédiction) :

$$\operatorname{MISE}_f[\hat{f}] = \mathbb{E}[\|\hat{f} - f\|_2^2]$$

- On suppose que le design $\{X_i\}_i$ est unidimensionnel régulier sur $[0, 1]$:

$$X_i = \frac{i}{n}, i = 1, \dots, n$$

Réduction au modèle de suites gaussiennes

Afin de présenter les fondements théoriques des méthodes basées sur les projections, il est commode de se ramener à un modèle inérialisé, celui de suites gaussiennes. Pour effectuer cette réduction, rappelons que, dans un modèle de régression, on observe $\{Y_i\}_i$ tels que :

$$Y_i = f\left(\frac{i}{n}\right) + \xi_i, i = 1, \dots, n, \text{ où } \xi_1, \dots, \xi_n \text{ sont iid et } \mathbb{E}[\xi_i] = 0, \operatorname{Var}[\xi_i] = \sigma^2.$$

On suppose que $f \in \mathcal{L}^2[0, 1]$. Soit $\{\phi_m, m \in \mathcal{M}\}$ une base orthonormée de $\mathcal{L}^2[0, 1]$:

$$\langle \phi_m, \phi_{m'} \rangle = \int_0^1 \phi_m(x) \phi_{m'}(x) dx = \delta_{m,m'} = \begin{cases} 0 & \text{si } m \neq m' \\ 1 & \text{sinon} \end{cases}$$

On pose $\theta_m = \langle f, \phi_m \rangle, m \in \mathcal{M}$ — \mathcal{M} étant isomorphe à \mathbb{N} , mais pas forcément égal, cf. cas d'une base d'ondelettes par exemple. D'après Parseval, on écrit :

$$\|\hat{f} - f\|_2^2 = \sum_{m \in \mathcal{M}} (\hat{\theta}_m - \theta_m)^2 \quad \text{où } \hat{\theta}_m = \langle \hat{f}, \phi_m \rangle$$

Conclusion. Estimer f équivaut à estimer la suite $\theta = \{\theta_m, m \in \mathcal{M}\}$. Si $\hat{\theta}$ est un estimateur de θ , alors $\hat{f} = \sum_{m \in \mathcal{M}} \hat{\theta}_m \phi_m$ est un estimateur de f .

Remarque 3. En utilisant les sommes de Riemann, on peut écrire les approximations suivantes :

$$\begin{aligned} \theta_m &= \int_0^1 f(x) \phi_m(x) dx \approx \frac{1}{n} \sum_{i=1}^n f\left(\frac{i}{n}\right) \phi_m\left(\frac{i}{n}\right) \\ \theta_m &\approx \underbrace{\frac{1}{n} \sum_{i=1}^n Y_i \phi_m\left(\frac{i}{n}\right)}_{Z_m} - \underbrace{\frac{\sigma}{\sqrt{n}} \frac{1}{\sqrt{n} \sigma} \sum_{i=1}^n \xi_i \phi_m\left(\frac{i}{n}\right)}_{\varepsilon_m} \end{aligned}$$

Alors, on a, $\forall m \in \mathcal{M}$:

$$Z_m = \theta_m + \frac{\sigma}{\sqrt{n}} \varepsilon_m \quad / \quad m \in \mathcal{M} \quad (2)$$

Le système (2) sous la condition que ε_m soit iid $\mathcal{N}(0, 1)$ est appelé **modèle de suites gaussiennes**.

Important. Les Z_m sont observables. θ_m sont les valeurs à estimer.

Note 4. En utilisant la version de Lyapunov du Théorème central limite, on peut montrer que lorsque n tend vers l'infini, la famille des v.a. $\{\varepsilon_m, m \in \mathcal{M}\}$ converge en loi vers une famille des variables iid Gaussiennes centrées réduites.

Proposition 5. Si ξ_1, \dots, ξ_n sont iid, $\mathbb{E}[\xi_i] = 0$ et $\mathbb{E}[\xi_i^2] = \sigma^2$, alors :

i. $\mathbb{E}[\varepsilon_m] = 0, \forall m \in \mathcal{M}$

ii. On a :

$$\lim_{n \rightarrow \infty} \text{Var}[\varepsilon_m] = 1$$

iii. On a également :

$$\lim_{n \rightarrow \infty} \text{Cov}[\varepsilon_m, \varepsilon_{m'}] = 0 \quad \text{pour } m \neq m'$$

Démonstration.

i. $\mathbb{E}_m = \frac{1}{\sqrt{n} \sigma} \sum_{i=1}^n \mathbb{E}[\xi_i] \phi_m\left(\frac{i}{n}\right) = 0$

ii. On peut écrire :

$$\begin{aligned}
 \text{Var}[\varepsilon_m] &= \frac{1}{n \sigma^2} \text{Var} \left[\underbrace{\sum_{i=1}^n \xi_i \phi_m \left(\frac{i}{n} \right)}_{\text{indépendants}} \right] \\
 &= \frac{1}{n \sigma^2} \sum_{i=1}^n \text{Var} \left[\underbrace{\xi_i \phi_m \left(\frac{i}{n} \right)}_{=\sigma^2 \phi_m \left(\frac{i}{n} \right)^2} \right] \\
 &= \frac{1}{n} \sum_{i=1}^n \phi_m \left(\frac{i}{n} \right)^2 \xrightarrow{n \rightarrow \infty} \int_0^1 \phi_m^2 = 1
 \end{aligned}$$

iii. De la même manière,

$$\begin{aligned}
 \text{Cov}(\varepsilon_m, \varepsilon_{m'}) &= \mathbb{E}[\varepsilon_m \varepsilon_{m'}] \\
 &= \frac{1}{n \sigma^2} \mathbb{E} \left[\sum_{i=1}^n \xi_i \phi_m \left(\frac{i}{n} \right) \sum_{j=1}^n \xi_j \phi_{m'} \left(\frac{j}{n} \right) \right] \\
 &= \frac{1}{n \sigma^2} \sum_{i=1}^n \mathbb{E} \left[\xi_i^2 \phi_m \left(\frac{i}{n} \right) \phi_{m'} \left(\frac{i}{n} \right) \right] \\
 &= \frac{1}{n} \sum_{i=1}^n \phi_m \left(\frac{i}{n} \right) \phi_{m'} \left(\frac{i}{n} \right) \\
 &\xrightarrow{n \rightarrow \infty} \int_0^1 \phi_m \phi_{m'} = 0
 \end{aligned}$$

□

Estimateurs classiques dans le modèle (2)

Estimateurs par projection (*cut-off*)

$$\hat{\theta}_m = \begin{cases} Z_m & \text{si } \|m\| \leq \lambda \\ 0 & \text{si } \|m\| > \lambda \end{cases} = Z_m \mathbf{1}(\|m\| < \lambda)$$

Idée. Comme $f \in \mathcal{L}^2[0, 1]$, on a :

$$\int f^2 < +\infty \Rightarrow \sum_{\mathcal{M}} \theta_m^2 < +\infty \Rightarrow \theta_m \xrightarrow{\|m\| \rightarrow +\infty} 0$$

L'estimateur par *cut-off* estime θ_m par 0 à partir d'un certain indice, ce qui n'induit pas beaucoup d'erreur car $\theta_m \rightarrow 0$.

Estimateurs par filtrage diagonal ($\mathcal{M} = \mathbb{N}$)

L'estimateur par *cut-off* étant un peu brutal, on peut définir $\hat{\theta}_m$ par :

$$\hat{\theta}_m = \alpha_m Z_m \quad \text{où } \alpha_m \in [0, 1] \text{ est une suite donnée}$$

Il existe par exemple :

- *Tikhonov* : $\alpha_m = \frac{1}{1 + \left(\frac{m}{\lambda}\right)^2}$ (années '70)

- *Pinsker* : $\alpha_m = \left(1 - \left|\frac{m}{\lambda}\right|^\beta\right)_+$ (année '80)

β et λ sont des paramètres de ces algorithmes.

Seuillage dur (*hard thresholding*)

$$\hat{\theta}_m = Z_m \mathbf{1}(|Z_m| > \lambda)$$

Seuillage doux (*soft thresholding*)

$$\hat{\theta}_m = |Z_m - \lambda| \operatorname{sign}(Z_m) \mathbf{1}(|Z_m| > \lambda)$$