

I. Séance Précédente :

$X_1, \dots, X_n$  iid de densité  $f$  sur  $[0, 1]$

$K \geq 1$  un entier,  $h = 1/K$ ,  $C_k = [(k-1)h; kh[$ .

Estimateur par histogramme :

$$\hat{f}_{n,h}(x) = \frac{1}{nh} \sum_{i=1}^n \mathbb{1}_{C_k}(X_i) \quad \text{si } x \in C_k; k=1, \dots, K.$$

On introduit les notations

$$p_k = \int_{C_k} f(x) dx \quad \hat{p}_k = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{C_k}(X_i)$$

On a vu que

$$E \left[ \|\hat{f}_{n,h} - f\|_2^2 \right] = \underbrace{\|f\|_2^2 - \frac{1}{h} \sum_{k=1}^K p_k^2}_{\text{Carré du biais}} + \underbrace{\frac{1}{nh} \sum_{k=1}^K p_k(1-p_k)}_{\text{Variance}}$$

De plus si  $f \in \text{Lip}(L)$ , alors

$$E \left[ \|\hat{f}_{n,h} - f\|_2^2 \right] \leq L^2 h^2 + \frac{1}{nh} \quad (*)$$

II. Remarques (1) Le minimum par rapport à  $h$  de (\*) est

atteint lorsque  $h = h_{\text{opt}} = (2nL^2)^{-1/3}$ . Dans ce cas, le risque est majoré par  $3(L/2)^{2/3} n^{-2/3}$

(2) Si on considère des fonctions  $f$  appartenant à la classe de Hölder :  $H(\beta, L) = \{f : |f(x) - f(y)| \leq L|x-y|^\beta\}$

alors la borne (\*) se transforme en

$$\mathbb{E}[\|\hat{f}_{n,h} - f\|_2^2] \leq L^2 h^{2\beta} + \frac{1}{nh}$$

Le minimum par rapport à  $h$  de cette expression est atteint lorsque  $h_{\text{opt}} = C n^{-\frac{1}{2\beta+1}}$ . Le risque optimal est alors de l'ordre de  $n^{-2\beta/2\beta+1}$ .

③ On voit que le risque optimal ainsi que la fenêtre optimale dépendent des valeurs  $(L, \beta)$  qui sont inconnues en pratique. Pour pallier ce défaut, on introduit des estimateurs adaptatifs. L'idée est de choisir  $h$  en fonction des données.

④ Les résultats précédents se généralisent facilement au cas multidimensionnel :  $X_1, \dots, X_n$  iid de densité  $f: [0,1]^d \rightarrow \mathbb{R}$ . La vitesse optimale est alors obtenue en minimisant par rapport à  $h$  l'expression  $L^2 h^{2\beta} + \frac{1}{nh^d}$ . Cela donne  $h_{\text{opt}} = n^{-1/2\beta+d}$  et un risque de l'ordre de  $n^{-2\beta/2\beta+d}$ .

Cette vitesse est d'autant plus lente que la dimension est grande. De plus, ce ralentissement se fait de façon exponentielle, car  $n^{-2\beta/2\beta+d} = \exp\left\{-\frac{2\beta \ln n}{2\beta+d}\right\}$  et donc

si l'on augmente  $d$  par  $\Delta d$ , il faut augmenter  $n$  t.g.

$$\ln n_{\text{new}} \approx \ln n_{\text{old}} + \Delta d \Rightarrow n_{\text{new}} = n_{\text{old}} e^{\Delta d}.$$

On appelle ce phénomène "curse of dimensionality"

On y reviendra plus tard dans ce cours.

### III. Adaptation par minimisation de l'estimateur sans biais du risque

#### ① Idee principale

Supposons maintenant qu'on a  $M$  fenêtrés potentielles:

$$h_1, h_2, \dots, h_M$$

Comment trouver celle qui donne le meilleur estimateur?

$$\text{Notons } L(h, f) = \|\hat{f}_{n,h} - f\|_2^2$$

On a vu que

$$\begin{aligned} L(h, f) &= \|f\|_2^2 - \frac{1}{h} \sum_{k=1}^K p_k^2 + \frac{1}{nh} \sum_{k=1}^K p_k (1-p_k) \\ &= \|f\|_2^2 + J(h). \end{aligned}$$

La meilleure fenêtre est celle qui minimise le risque, ce qui équivaut à minimiser  $J(h)$ :

$$h^* \in \arg \min_{h \in \{h_1, \dots, h_M\}} J(h).$$

Or  $h^*$  est impossible à calculer à cause de non-disponibilité de la densité  $f$  et donc des  $p_k$ .

Pour contourner cette difficulté, on remplace  $J$  par un estimateur sans biais.

Proposition: Pour toute fenêtre déterministe  $h$ , l'expression

$$\hat{J}_n(h) = \frac{2}{(n-1)h} - \frac{n+1}{(n-1)h} \sum_{k=1}^K \hat{p}_k^2$$

fournit un estimateur sans biais de  $J(h)$ .

## Preuve.

Comme on l'a déjà remarqué  $n\hat{p}_k \sim B(n, p_k)$

[la variable aléatoire  $n\hat{p}_k$  suit la loi binomiale de paramètres  $n$  et  $p_k$ .]

Par conséquent  $E[n\hat{p}_k] = np_k$  et  $\text{Var}[n\hat{p}_k] = np_k(1-p_k)$ .

On en déduit que

$$E[\hat{p}_k^2] = \text{Var}[\hat{p}_k] + (E[\hat{p}_k])^2 = \frac{p_k(1-p_k)}{n} + p_k^2.$$

Par conséquent,

$$\begin{aligned} E[\hat{J}_n(h)] &= \frac{2}{(n-1)h} - \frac{n+1}{(n-1)h} \sum_{k=1}^K \left( \frac{p_k - p_k^2}{n} + p_k^2 \right) \\ &= \frac{2}{(n-1)h} - \frac{n+1}{n(n-1)h} \sum_{k=1}^K p_k - \frac{n+1}{(n-1)h} \sum_{k=1}^K \frac{n-1}{n} p_k^2 \\ &= \frac{1}{nh} - \frac{n+1}{nh} \sum_{k=1}^K p_k^2 = J(h) \quad \blacksquare \end{aligned}$$

## ② Algorithme :

- On a  $X_1, \dots, X_n$ .
- On se donne un paramètre  $a > 1$  (typiquement  $a=1,2$ )
- On pose  $h_1 = \frac{1}{n}$ ,  $h_m = \frac{1}{[na^{1-m}]}$   $m=1, \dots, M$
- On calcule  $\hat{J}_n(h_1), \dots, \hat{J}_n(h_M)$   
et  $\hat{m} \in \arg \min_m \hat{J}_n(h_m)$
- On pose  $\hat{f}(x) = \hat{f}_{n, h_{\hat{m}}}(x)$ .

### ③ Pourquoi ça marche ?

Ecrivons la perte sous la forme :

$$L(h, f) = \|f\|_2^2 + \hat{J}_n(h) + \xi_n(h).$$

On a

$$\begin{aligned} E[\xi_n(h)] &= E[L(h, f)] - \|f\|_2^2 - E[\hat{J}_n(h)] \\ &= E[L(h, f)] - \|f\|_2^2 - J(h) = 0. \end{aligned}$$

Donc  $\{\xi_n(h) : h \in \{h_1, \dots, h_M\}\}$  est un vecteur aléatoire de moyenne nulle. Posons

$$V_n(h) = \text{Var}[\xi_n(h)]$$

Lemme 1. Quelle que soit la variable aléatoire  $\tilde{h}$  à valeurs dans  $\{h_1, \dots, h_M\}$ , on a

$$E[\xi_n(\tilde{h})] \leq \sqrt{E[V_n(\tilde{h})]} \times M^{1/2}.$$

Preuve.

$$E[\xi_n(\tilde{h})] = E\left[V_n(\tilde{h})^{1/2} \times \frac{\xi_n}{\sqrt{V_n}}(\tilde{h})\right]$$

$$\leq E\left[V_n(\tilde{h})^{1/2} \times \max_m \frac{|\xi_n|}{\sqrt{V_n}}(h_m)\right]$$

Cauchy-Schw.

$$\leq E[V_n(\tilde{h})]^{1/2} \times E\left[\max_m \frac{\xi_n^2}{V_n}(h_m)\right]^{1/2}$$

$$\leq E[V_n(\tilde{h})]^{1/2} \times \left(\sum_{m=1}^M \underbrace{E\left[\frac{\xi_n^2}{V_n}(h_m)\right]}_{=1}\right)^{1/2} \quad \blacksquare$$

On a donc

$$\mathbb{E}[L(\hat{h}, f)] = \|f\|_2^2 + \mathbb{E}[\hat{J}_n(\hat{h})] + \mathbb{E}[\xi_n(\hat{h})]$$

$$\leq \|f\|_2^2 + \mathbb{E}[\hat{J}_n(h^*)] + \sqrt{M \cdot \mathbb{E}[V_n(\hat{h})]}.$$

Cette inégalité est vraie pour tout  $h^* \in \{h_1, \dots, h_M\}$ ,

car  $\hat{h}$  est par définition le minimiseur de  $\hat{J}_n$ . Par conséquent, l'inégalité ci-dessus reste vraie pour

$$h^* = \arg \min_{h \in \{h_1, \dots, h_M\}} J(h) = \arg \min_h \mathbb{E}[L(h, f)]$$

On a donc

$$\begin{aligned} \mathbb{E}[L(\hat{h}, f)] &\leq \|f\|_2^2 + \mathbb{E}[\hat{J}_n(h^*)] + \sqrt{M \mathbb{E}[V_n(\hat{h})]} \\ &= \mathbb{E}[L(h^*, f)] - \mathbb{E}[\xi_n(h^*)] + \sqrt{M \mathbb{E}[V_n(\hat{h})]} \\ &\leq \mathbb{E}[L(h^*, f)] + \sqrt{M \mathbb{E}[V_n(h^*)]} \end{aligned}$$

Cela nous conduit vers l'inégalité oracle suivante:

$$\mathbb{E}[L(\hat{h}, f)] \leq \underbrace{\min_{h \in \{h_1, \dots, h_M\}} \mathbb{E}[L(h, f)]}_{\text{risque de l'oracle}} + \underbrace{\sqrt{M} \cdot \sqrt{\max_m V_n(h_m)}}_{\text{terme résiduel}}$$

Le choix de  $h_1, \dots, h_M$  proposé dans l'algorithme garantit

que  $M = \text{Const} \cdot \ln n$  est à croissance lente. On peut

également démontrer que  $V_n(h_m) \leq C/n$  pour une constante  $C$  bien choisie.

④ A quoi ressemble  $\xi_n(h)$  ?

$$\begin{aligned}
 \xi_n(h) &= L(h, f) - \|f\|_2^2 - \hat{J}_n(h) \\
 &= \|\hat{f}_{n,h}\|_2^2 - 2\langle \hat{f}_{n,h}, f \rangle - \hat{J}_n(h) \\
 &= \sum_{k=1}^K \frac{1}{h} \hat{P}_k^2 - 2 \sum_{k=1}^K \int_{C_k} \frac{\hat{P}_k}{h} f(x) dx - \hat{J}_n(h) \\
 &= \frac{1}{h} \sum_{k=1}^K \hat{P}_k^2 - \frac{2}{h} \sum_{k=1}^K \hat{P}_k P_k - \frac{2}{(n-1)h} + \frac{n+1}{(n-1)h} \sum_{k=1}^K \hat{P}_k^2 \\
 &= \frac{2n}{(n-1)h} \sum_{k=1}^K \hat{P}_k^2 - \frac{2}{h} \sum_{k=1}^K [(\hat{P}_k - P_k) P_k + P_k^2] - \frac{2}{(n-1)h} \sum_{k=1}^K P_k \\
 &= \frac{2n}{(n-1)h} \sum_{k=1}^K \left[ \hat{P}_k^2 - \underbrace{\frac{n-1}{n} P_k^2 - \frac{1}{n} P_k}_{\mathbb{E}[\hat{P}_k^2]} \right] - \frac{2}{h} \sum_{k=1}^K P_k (\hat{P}_k - P_k)
 \end{aligned}$$

Posons

$$\xi_n^{(1)}(h) = \frac{2n}{(n-1)h} \sum_{k=1}^K (\hat{P}_k^2 - \mathbb{E}[\hat{P}_k^2])$$

$$\xi_n^{(2)}(h) = \frac{2}{h} \sum_{k=1}^K P_k (\hat{P}_k - P_k).$$

Lemme Quelle que soit la valeur  $h > 0$ , on a

$$\begin{aligned}
 \text{Var}[\xi_n^{(2)}(h)] &= \frac{4}{nh^2} \sum_{k=1}^K P_k^3 (1 - P_k) \\
 &\leq (4/n) \times \sup_{x \in [0,1]} f(x)^2.
 \end{aligned}$$

Preuve Il suffit de remarquer que  $\xi_n^{(2)}(h) = \frac{2}{nh} \sum_{i=1}^n Z_i$

où les variables aléatoires  $Z_i$  sont iid de variance

$\sum_{k=1}^K P_k^3 (1 - P_k)$ . Pour conclure, on utilise  $1 - P_k \leq 1$  et  $P_k \leq h \cdot \sup_x f(x)$ .