

Chapitre 2 :  
Estimateurs à noyau  
en dimension supérieure

A - Le cas multivarié (FV, chap IV)

On suppose à partir de maintenant que  $X \in \mathbb{R}^d$  et  $Y \in \mathbb{R}$  pour la régression (on pourrait avoir  $Y \in \mathbb{R}^p$ , ça ne changerait pas les vitesses de convergence).

$$\hat{f}_X(x) = \frac{1}{mh^d} \sum_{i=1}^m K\left(\frac{X_i - x}{h}\right)$$

$$\hat{\alpha}(x) = \begin{cases} \frac{\sum_{i=1}^m Y_i K\left(\frac{X_i - x}{h}\right)}{\sum_{i=1}^m K\left(\frac{X_i - x}{h}\right)} & \text{si } \sum_{i=1}^m K\left(\frac{X_i - x}{h}\right) \neq 0 \\ \bar{Y} & \text{sinon} \end{cases}$$

$$\hat{\alpha}(x) = \frac{\hat{g}(x)}{\hat{f}(x)} \quad \text{avec} \quad \hat{g}(x) = \frac{1}{mh^d} \sum_{i=1}^m Y_i K\left(\frac{X_i - x}{h}\right)$$

On a  $K: \mathbb{R}^d \rightarrow \mathbb{R}$   
 $h$ : paramètre de lissage

On peut définir (FV, chap IV, p. 78) l'équivalent d'un noyau d'ordre  $k$  en multivarié.

• On obtient sous hypothèses que  $f$  et  $\alpha$  sont  $C^k$  et autres, en convergence presque complète:  $O(h^k) + O\left(\sqrt{\frac{\ln(m)}{mh^d}}\right)$  pour convergence ponctuelle ou uniforme.

• Sous l'hypothèse  $f$  et  $\alpha$  Hölder( $\beta$ ) et d'autres:

convergence p.c.o. | uniforme  
ponctuelle  $O(h^\beta) + O\left(\sqrt{\frac{\ln(m)}{mh^d}}\right)$

• Sous l'hypothèse  $f$  et  $x \in \mathcal{C}^2$  et autres,  
convergence p.c.o.:

$$MSE(\hat{\pi}(x)) = B^2(x)h^4 + V(x) \frac{1}{mh^d} + o\left(h^4 + \frac{1}{mh^d}\right)$$

$$MISE(\hat{\pi}) = B^2h^4 + V \frac{1}{mh^d} + o\left(h^4 + \frac{1}{mh^d}\right)$$

Si on cherche le minimisant AMISE :

$$\frac{\partial AMISE}{\partial h}(h) = 4B^2h^3 - d \frac{V}{m} \frac{1}{h^{d+1}}$$

$$\text{d'où } \frac{\partial AMISE}{\partial h}(h) = 0 \Leftrightarrow \frac{4B^2m}{Vd} = \frac{1}{h^{d+1}}$$

$$\Leftrightarrow h^{d+1} = \frac{Vd}{4B^2m} \Leftrightarrow h = \left(\frac{Vd}{4B^2}\right)^{\frac{1}{d+1}} m^{-\frac{1}{d+1}}$$

Plus  $d$  est grand, plus on prend  $h$  grand.

Pour ce  $h$  précis, on a :

$$MISE(\hat{\pi}) = B^2 \frac{Vd}{4B^2mh^d} + V \frac{1}{mh^d} + o\left(h^4 + \frac{1}{mh^d}\right)$$

$$= \frac{1}{mh^d} \left(\frac{Vd}{4} + V\right) + o\left(\frac{1}{mh^d}\right)$$

$$= \frac{V}{mh^d} \left(\frac{d}{4} + 1\right) + o\left(\frac{1}{mh^d}\right)$$

ou

$$MISE(\hat{\pi}) = \left(B^2 + V \frac{4B^2}{Vd}\right) h^4 + o(h^4)$$

$$= B^2 \left(1 + \frac{4}{d}\right) \left(\frac{Vd}{4B^2}\right)^{\frac{d}{d+1}} m^{-\frac{4}{d+1}} + o\left(m^{-\frac{4}{d+1}}\right)$$

On remarque que la dimension  $d$  dégrade la vitesse de convergence optimale pour la MISE, MSE, AMISE, ...

C'est dû au fléau de la dimension. Il y a une rarefaction des données lorsque  $d$  augmente.

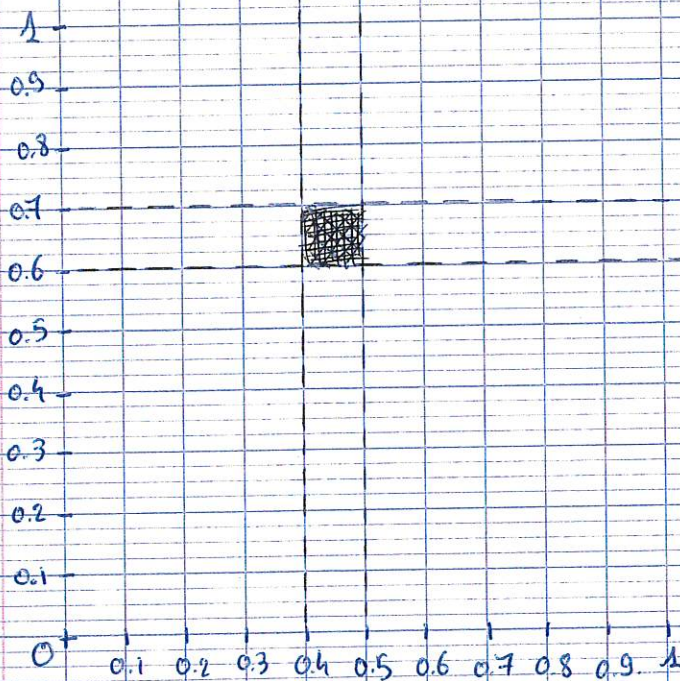
Supposons que l'on ait un échantillon de taille  $m = 10^6$  iid de loi  $U_{[0,1]^d}$ .

Si  $d=1$ :



$N$ , le nombre d'individus sur un segment de longueur  $0.1$  est de loi  $\text{Bin}(m, 0.1)$  dont l'espérance vaut  $100\,000 =$

Si  $d=2$ :



$N$  le nombre d'individus sur un carré de longueur  $0.1$  suit une loi  $\text{Bin}(m, 0.01)$  dont l'espérance est  $m \times 0.01 = 10^4 = 10\,000$

Si  $d > 2$ , le nombre d'individus  $N$  sur une portion (carré) de longueur de côté  $0.1$  a pour loi une  $\text{Bin}(m, 0.1^d)$  d'espérance  $m \times (0.1)^d = m 10^{-d} = 10^{6-d}$ .

$$d=3: m(0.1)^d = 1000$$

$$d=5: m(0.1)^d = 10$$

$$d=4: m(0.1)^d = 100$$

$$d=6: m(0.1)^d = 1$$

Soit on augmente  $m$  (en général impossible), soit on prend  $h$  plus grand, ce qui dégrade la vitesse de convergence.

Stone (1981, 1982) montre que la vitesse optimale d'un estimateur, quel qu'il soit, dans un modèle de régression avec  $X \in \mathbb{R}^d$  et en supposant  $f, \varepsilon \in \mathcal{C}^k$  est donné par  $m^{-\frac{k}{2k+d}}$  en norme  $\mathbb{L}_q, q < \infty$  et  $\left(\frac{m}{\ln(m)}\right)^{-\frac{k}{2k+d}}$  en norme  $\mathbb{L}_\infty$ .

On remarque que notre estimateur à noyau atteint bien ces vitesses optimales :

$$\begin{array}{llll} \text{MSE} \sim C m^{-\frac{k}{4+d}} & k=2 & \text{vitesse} \sim C m^{-\frac{2}{4+d}} \\ \text{AMISE} \sim C m^{-\frac{k}{4+d}} & k=2 & \text{vitesse} \sim C m^{-\frac{2}{4+d}} \end{array}$$

Le problème ne vient donc pas de l'estimateur mais du modèle qui est trop complexe à estimer. Peu considérée pour  $d \geq 4$  ou 5.

On considère alors souvent des modèles semi-paramétriques

Modèle additif :

$$\varepsilon(X) = \mu + \sum_{j=1}^d \varepsilon_j(X_j)$$

Modèle additif étendu :

$$\varepsilon(X) = L\left(\mu + \sum_{j=1}^d \varepsilon_j(X_j)\right) \quad \text{où } L \text{ une fonction connue}$$

Modèle de transformations optimales :

$$\mathbb{E}[T(Y) | X] = \mu + \sum_{j=1}^d \varepsilon_j(X_j) \quad \text{où } T \text{ fonction inconnue}$$

Modèle single index

$$Y = \varepsilon(\langle X, \theta \rangle) + \varepsilon \quad \text{où } \varepsilon, \theta \text{ inconnus}$$

Modèle multi-index :

$$Y = \sum_{j=1}^d \alpha_j \langle X, \theta_j \rangle + \varepsilon \quad \text{où} \quad \begin{array}{|l} \alpha_j \\ \theta_j \end{array} \text{ inconnus.}$$

B - le cas fonctionnel

$X \in (\mathcal{F}, d)$  où

- $\mathcal{F}$  espace de fonction
- $d$  pseudo-distance (même propriété de distance sauf  $d(x, y) = 0 \nRightarrow x = y$ )

Régression : Ferraty et Vieu (2000)  $x \in \mathcal{F}$

$$\hat{x}(x) = \begin{cases} \frac{\sum_{i=1}^m Y_i k\left(\frac{d(X_i, x)}{h}\right)}{\sum_{i=1}^m k\left(\frac{d(X_i, x)}{h}\right)} & \text{si } \sum_{i=1}^m k\left(\frac{d(X_i, x)}{h}\right) \neq 0 \\ \bar{Y} & \text{sinon} \end{cases}$$

Remarque : • En pratique, on dispose seulement d'observations discrétisées de cette courbe. On suppose que le passage discrétisé à courbe est effectué en amont.

• Le fléau de la dimension peut dans un premier temps faire penser que cela ne peut pas fonctionner. C'est vrai si on prend un processus gaussien et d'une norme alors les vitesses de convergence sont en  $(\ln(m))^{-\delta}$ . Mais les courbes que l'on a à traiter contiennent souvent une certaine régularité, une structure, une dynamique particulière qui lui confère un statut particulier, qui font qu'une telle approche peut être tout de même adaptée à condition de choisir  $d$  qui capte l'information utile de  $X$ . (une courbe  $\neq$  ensemble de v.a. iid).

On peut par exemple choisir pour  $d$

$$\bullet d_R(x_1, x_2) = \sqrt{\int (x_1^{(R)}(t) - x_2^{(R)}(t))^2 dt}$$

$$\bullet d_{PCA}^k(x_1, x_2) = \left\| \sum_{j=1}^k \langle x_1 - x_2, \Psi_j \rangle \Psi_j \right\| \text{ où } \Psi_j \text{ } j^{\text{ème}} \text{ fct. propres de l'ACP}$$

$$\bullet d_{\Psi}^k(x_1, x_2) = \left\| \sum_{j=1}^k \langle x_1 - x_2, \Psi_j \rangle \Psi_j \right\| \text{ où } \Psi_j \text{ base donnée (Fonct. Ordonnée)}$$

Le choix de  $d$  en pratique se fait par validation croisée :

$$CV(d) = \sum_{i=1}^m w(x_i) (y_i - \hat{x}_d^{-i}(x_i))^2$$