

Chapitre 1 :

Estimateur à moynen pour Variables réelles

A - Historique

On se limite ici au cas de l'estimation de la densité (une évolution similaire s'est faite pour l'estimateur de la fonction de régression à partir du régressogramme)

• Modèle paramétrique :

On suppose $f_x = f_\theta$ connue à un paramètre $\theta \in \mathbb{R}^p$ puis et on estime f_x en remplaçant θ par son estimateur.
 $\forall x \in \mathbb{R}, \hat{f}_x(x) = \hat{f}_\theta(x)$

• Viennent ensuite les modèles non-paramétriques.

* Un premier estimateur non paramétrique est l'histogramme. Il consiste à estimer la densité au travers d'une fonction constante par morceaux (correspondant aux classes). Sur chaque classe, la densité f_x est estimée par la densité associée à cette classe (concentration des observations au sein de la classe).

Cet estimateur très rudimentaire présente un assez grand nombre de défauts.

↳ il est discontinu et constant par morceaux

↳ il dépend fortement du choix des classes (leurs nombres et leurs bornes).

$$F_x(x) = \mathbb{P}(X \leq x)$$

$$X_1, \dots, X_m \stackrel{iid}{\sim} X$$

Comment estimer F_x de manière non-paramétrique ?

Remarque : $\forall x \in \mathbb{R}$, $Z_x = \mathbb{1}_{X \leq x}$, on a $\forall x \in \mathbb{R}$, $F_x(x) = \mathbb{E}[Z_x]$

Donc, on pose naturellement, avec la méthode des moments,
 $\forall x \in \mathbb{R}$, $\hat{F}_x(x) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{X_i \leq x}$

Propriétés : 1 - $\hat{F}_x(x)$ est un estimateur sans biais de $F_x(x)$

2 - $\forall x \in \mathbb{R}$, $\text{var}(\hat{F}_x(x)) = \frac{F(x)(1-F(x))}{m}$

3 - $\forall x \in \mathbb{R}$, $\hat{F}_x(x) \xrightarrow[m \rightarrow +\infty]{L^2} F(x)$

4 - $\forall x \in \mathbb{R}$, $\hat{F}_x(x) \xrightarrow[m \rightarrow +\infty]{\text{p.s.}} F(x)$ (LFGN).

5 - $\forall x \in \mathbb{R}$, $\sqrt{m}(\hat{F}_x(x) - F(x)) \xrightarrow[m \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, F(x)(1-F(x)))$

Démonstration :

1 - $\mathbb{E}[\hat{F}_x(x)] = \frac{1}{m} \sum_{i=1}^m \mathbb{E}[\mathbb{1}_{X_i \leq x}] = \frac{1}{m} \sum_{i=1}^m \mathbb{P}(X_i \leq x) = \mathbb{P}(X_1 \leq x) = F(x)$

2 - $\text{var}(\hat{F}_x(x)) = \frac{1}{m^2} \text{var}\left(\sum_{i=1}^m \mathbb{1}_{X_i \leq x}\right) \stackrel{\text{II}}{=} \frac{1}{m^2} \sum_{i=1}^m \text{var}(X_i \mathbb{1}_{X_i \leq x})$
 $= \frac{1}{m} \text{var}(X_1) = \frac{F(x)(1-F(x))}{m}$ car $X_1 \sim \mathcal{B}(F(x))$.

3 - $\mathbb{E}[\hat{F}_x(x)^2] = \mathbb{E}\left[\frac{1}{m^2} \sum_{i=1}^m \mathbb{1}_{X_i \leq x}\right] + \frac{2}{m^2} \mathbb{E}\left[\sum_{i=2}^m \sum_{j=1}^{i-1} \mathbb{1}_{X_i \leq x} \mathbb{1}_{X_j \leq x}\right]$
 $= \frac{1}{m} F(x) + \frac{2}{m^2} F(x)^2 \sum_{i=1}^{m-1} i = \frac{F(x)}{m} + \frac{m-1}{m} F(x)^2$

d'où :

$\mathbb{E}[|\hat{F}_x(x) - F(x)|^2] = \mathbb{E}[\hat{F}_x(x)^2] - 2\mathbb{E}[\hat{F}_x(x)]F(x) + F(x)^2$
 $= \frac{F(x)}{m} + \frac{m-1}{m} F(x)^2 - F(x)^2$

$= \frac{1}{m} (F(x) - F(x)^2) \xrightarrow[m \rightarrow +\infty]{} 0$

d'où $F_x(x) \xrightarrow[m \rightarrow +\infty]{L^2} F(x)^m$ pour tout $x \in \mathbb{R}$.

4 - LFGN et 5 - TCL.

5 - Estimateur de la densité par fenêtre mobile.

Tout part du fait que $f(x) = \lim_{h \rightarrow 0} \frac{F_x(x+h) - F_x(x-h)}{2h}$.

$$\begin{aligned} \text{En effet, } \frac{F_x(x+h) - F_x(x-h)}{2h} &= \frac{F_x(x+h) - F_x(x) + F_x(x) - F_x(x-h)}{2h} \\ &= \frac{1}{2} \frac{F_x(x+h) - F_x(x)}{h} + \frac{1}{2} \frac{F_x(x) - F_x(x-h)}{h} \\ &\xrightarrow{h \rightarrow 0} \frac{1}{2} (f(x) + f(x)) = f(x). \end{aligned}$$

Il suffit donc de prendre h assez petit et d'estimer F_x (avec l'estimateur précédent) pour avoir un estimateur de

$\frac{f}{\sqrt{x}}$. On a alors :

$$\begin{aligned} \forall h > 0, \quad \hat{f}_x(x) &= \frac{\hat{F}_x(x+h) - \hat{F}_x(x-h)}{2h} = \frac{1}{2hm} \sum_{i=1}^m \mathbb{1}_{X_i \leq x+h} - \mathbb{1}_{X_i \leq x-h} \\ &= \frac{1}{2hm} \sum_{i=1}^m \mathbb{1}_{x-h < X_i \leq x+h} \\ &= \frac{1}{m} \sum_{i=1}^m \frac{1}{h} w\left(\frac{X_i - x}{h}\right) \quad \text{où } w(t) = \frac{1}{2} \mathbb{1}_{[-1;1]}(t) \quad \forall t \in \mathbb{R}, \end{aligned}$$

Définition : h est appelé paramètre de lissage (cf. exemple)

Exemple : On prend $x = (2, 3, 3.5, 6, 7.5)$. On choisit $h = 1$.

Remarque : Comme X est continue, $IP(X_i = x_0) = 0$ pour tout $x_0 \in \mathbb{R}$ et donc on considère souvent de manière équivalente l'estimateur par fenêtre mobile :

$$\hat{f}_x(x) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2h} \mathbb{1}_{[-1;1]}\left(\frac{X_i - x}{h}\right) \stackrel{\text{p.s.}}{=} \hat{f}_x(x)$$

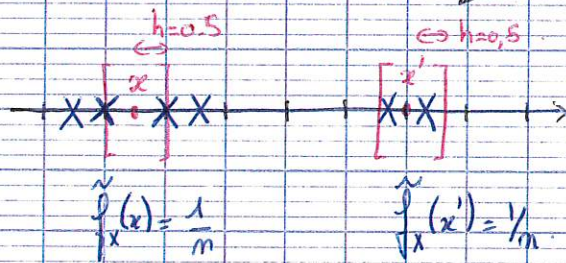
Exemple : On reprend les mêmes données

Exemple : Retour aux données calculées galaxies
 Calculer les estimateurs par fenêtres mobiles
 pour des valeurs de $h = 0,1, 10, 100, 1000, 10000, 100000$,
 Choisir le paramètre de lissage h . Choisir le plus
 adapté. Quels défauts reste-t-il à gommer ?

D'après mes expérimentations, on remarque :

* Avantage : • ne dépend plus du choix des classes (leurs
 nombres, leurs bornes) mais simplement des valeurs
 autour de x (et donc du choix de h)

* Défauts : • discontinu, continu par morceaux
 • même poids donné à tous les points se
 trouvant dans l'intervalle $[x-h; x+h]$



→ on voudrait
 pondérer en fonction
 de la distance entre
 x_i et x .

Rappel : Propriété théorique d'un estimateur

A - Estimateur Ponctuel

1 - Biais d'un estimateur \hat{f} de f au point x

$$\text{Biais}_f(\hat{f}(x)) = \mathbb{E}_f[\hat{f}(x)] - f(x)$$

Un estimateur est dit (asymptotiquement) sans biais
 si son biais est 0 (tend vers 0 quand $n \rightarrow +\infty$).

2 - Variance d'un estimateur \hat{f} de f au point x

$$\text{var}_f(\hat{f}(x)) = \mathbb{E}_f[(\hat{f}(x) - \mathbb{E}_f[\hat{f}(x)])^2]$$

3 - Erreur quadratique moyenne (Mean Squared Error)

$$\begin{aligned} \text{MSE}(\hat{f}(x)) &= \mathbb{E}_f[(\hat{f}(x) - f(x))^2] \\ &= \text{var}_f(\hat{f}(x)) + (\text{Biais}_f(\hat{f}(x)))^2 \end{aligned}$$

Un estimateur est dit consistant en moyenne quadratique si son MSE tend vers 0 lorsque $n \rightarrow +\infty$.

B - Globales

1 - Integrated Mean Squared Error (ou Mean Integrated Squared Error)

$$MISE(\hat{f}) = \int \mathbb{E}[(\hat{f}(x) - f(x))^2] dx$$

on peut inverser \int et \mathbb{E} par Fubini-Tonelli d'où les deux mots équivalents.

2 - Integrated Squared Error

$$ISE(\hat{f}) = \int (\hat{f}(x) - f(x))^2 dx$$

3 - Mean Sum of Squared Error

$$MSSE(\hat{f}) = \mathbb{E}\left[\sum_{j=1}^m (\hat{f}(x_j) - f(x_j))^2\right]$$

(*) on remplace parfois les x_j par des points déterministes t_i formant une grille de discrétisation $t = (t_1, \dots, t_p)$.

4 - Sum of Squared Error

$$SSE(\hat{f}) = \sum_{j=1}^m (\hat{f}(x_j) - f(x_j))^2$$

Proposition: Propriétés de l'estimateur de la densité par fenêtre mobile.

1 - $\forall \epsilon, h = h_m \xrightarrow{m \rightarrow +\infty} 0$ alors, $\forall x \in \mathcal{D}_{f_x}$, $\hat{f}_x(x)$ est un estimateur sans biais de $f_x(x)$.

2 - $\forall \epsilon, h = h_m \xrightarrow{m \rightarrow +\infty} 0$ et $mh_m \xrightarrow{m \rightarrow +\infty} +\infty$. Soit $x \in \mathcal{D}_{f_x}$, si $f_x(x) > 0$ alors on a $\hat{f}_x(x) \xrightarrow{m \rightarrow +\infty} f_x(x)$.

Démonstration:

$$1 - \mathbb{E}\left[\hat{f}_x(x)\right] = \mathbb{E}\left[\frac{\hat{F}_x(x+h) - \hat{F}_x(x-h)}{2h}\right] = \frac{F_x(x+h) - F_x(x-h)}{2h} \xrightarrow{m \rightarrow +\infty} f_x(x)$$

$$2 - \text{var}\left(\hat{f}_x(x)\right) = \text{var}\left(\frac{1}{m} \sum_{i=1}^m \frac{1}{2h} \mathbb{1}_{[x-h, x+h]} \left(\frac{X_i - x}{h}\right)\right)$$

$$\begin{aligned} \text{var}(\hat{f}_X(x)) &= \frac{1}{4m^2h^2} \text{var}\left(\sum_{i=1}^m \mathbb{1}_{[-1,1]}\left(\frac{X_i - x}{h}\right)\right) \\ &\stackrel{\text{iid}}{=} \frac{1}{4m^2h^2} \sum_{i=1}^m \text{var}\left(\mathbb{1}_{[-1,1]}\left(\frac{X - x}{h}\right)\right) \end{aligned}$$

$$\begin{aligned} \text{Or } \text{var}\left(\mathbb{1}_{[-1,1]}\left(\frac{X-x}{h}\right)\right) &= \mathbb{E}\left[\mathbb{1}_{[-1,1]}^2\left(\frac{X-x}{h}\right)\right] - \mathbb{E}\left[\mathbb{1}_{[-1,1]}\left(\frac{X-x}{h}\right)\right]^2 \\ &= \mathbb{P}\left(-1 \leq \frac{X-x}{h} \leq 1\right) - \left[\mathbb{P}\left(-1 \leq \frac{X-x}{h} \leq 1\right)\right]^2 \\ &= \mathbb{P}(x-h \leq X \leq x+h) - \mathbb{P}(x-h \leq X \leq x+h)^2 \\ &= F_X(x+h) - F_X(x-h) - (F_X(x+h) - F_X(x-h))^2 \end{aligned}$$

$$\begin{aligned} \text{var}(\hat{f}_X(x)) &= \frac{1}{4m^2h^2} (F_X(x+h) - F_X(x-h)) (1 - (F_X(x+h) - F_X(x-h))) \\ &= \frac{1}{2nh} \left(\frac{F_X(x+h) - F_X(x-h)}{2h} \right) \left(1 - \frac{2h}{2h} \frac{F_X(x+h) - F_X(x-h)}{2h} \right) \\ &\underset{n \rightarrow \infty}{\sim} \frac{1}{2nh} f_X(x) \xrightarrow{n \rightarrow \infty} 0 \quad \text{car } nh_m \xrightarrow{n \rightarrow \infty} +\infty \end{aligned}$$

Cela implique directement que

$$\text{MSE}(\hat{f}_X(x)) = \text{var}(\hat{f}_X(x)) + \text{Biais}(\hat{f}_X(x))^2 \xrightarrow{n \rightarrow \infty} 0$$

Remarque : A n fixé, on a intérêt à prendre

- h petit pour rendre le terme de biais petit
- h grand pour rendre le terme de variance grand.

Il faudra donc trouver un équilibre entre le biais et la variance. Prendre h trop petit risquent à avoir un surajustement aux données (biais faible et grande variance). Estimateur proche de 0 sauf près des observations où l'on observe des pics. Prendre h trop grand risquent à avoir un sous-ajustement aux données (biais fort et variance faible). Estimation de la densité

trop aplati, les variations de la densité sont sous-représentées

B - Estimateur à noyaux lisses de la fonction de densité

I. Définition

Ideé : On veut remplacer l'indicatrice $\mathbb{1}_A$ par une fonction plus lisse k , appelée noyau possédant les propriétés suivantes :

$$\hookrightarrow \forall x \in \mathbb{R}, k(x) \geq 0 \quad (\text{positivité})$$

$$\hookrightarrow \int k(x) dx = 1 \quad (\text{densité})$$

$$\hookrightarrow \int_{\mathbb{R}} x k(x) dx = 0 \quad (\text{symétrie})$$

$$\hookrightarrow k \text{ maximale en } 0 \text{ et } k(x) \xrightarrow{|x| \rightarrow +\infty} 0 \text{ en décroissant}$$

Définition : Estimateur à noyau de la densité : Parzen (1964)

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

où k est la fonction définissant le poids de chaque observation (appelé noyau, en anglais "kernel").

h est le paramètre de lissage (en anglais bandwidth ou smoothing parameter).

Exemple : 1^{er} exemple de noyau : $k(t) = \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}}$

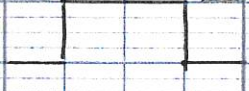
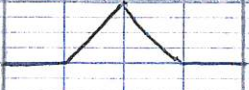
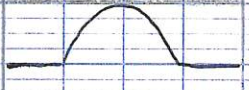
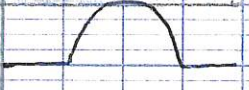

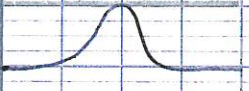
Dans ce cas précis,

$$\hat{f}_X(x) = \frac{1}{n} \sum_{i=1}^m \frac{1}{\sqrt{2\pi}h} e^{-\frac{(X_i - x)^2}{2h^2}} = \frac{1}{\sqrt{2\pi}h^2} \frac{1}{n} \sum_{i=1}^m e^{-\frac{(x - X_i)^2}{2h^2}}$$
$$= \frac{1}{n} \sum_{i=1}^m \frac{1}{\sqrt{2\pi}h} e^{-\frac{(X_i - x)^2}{2h^2}}$$

$\underbrace{\hspace{10em}}_{\hat{f}_X(x; h^2)}$

(voir exemple en code R).

Quelques exemples de noyau

noyau k	$k(t) =$	
Uniforme	$\frac{1}{2} \mathbb{1}_{\{ t \leq 1\}}$	
Triangulaire	$(1 - t) \mathbb{1}_{\{ t \leq 1\}}$	
Epanechnikov	$\frac{3}{4} (1 - t^2) \mathbb{1}_{\{ t \leq 1\}}$	
Quartic	$\frac{15}{16} (1 - t^2)^2 \mathbb{1}_{\{ t \leq 1\}}$	
Triweight	$\frac{35}{32} (1 - t^2)^3 \mathbb{1}_{\{ t \leq 1\}}$	
Gaussien	$\frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}}$	
Cosinus	$\frac{\pi}{4} \cos\left(\frac{t\pi}{2}\right) \mathbb{1}_{\{ t \leq 1\}}$	

Remarques : 1- On choisit un noyau qui soit maximal en 0 et dont la valeur diminue lorsque l'on s'écarte de 0 afin de moduler l'importance donnée aux observations en fonction de leur proximité avec le point x où l'on veut estimer f_x . Plus une observation sera proche, plus son poids sera élevée (car $\frac{x_i - x}{h}$ petit), et réciproquement.

2- Le fait d'utiliser un noyau lisse permet d'obtenir un estimateur lisse de f_x qui n'est plus constant par morceaux (si k ne l'est pas)

3- h joue le même rôle que précédemment pour l'estimateur à fenêtre mobile. Prendre h trop grand fait que tous les $\frac{x_i - x}{h}$ sont proche de 0 et donc que

$$\mathbb{E}_f \left[\hat{f}_x(x) \right] = \int_{\mathbb{R}} k(t) R(x+th) dt + \frac{h^2}{2} f''(x) \mu_2(k)$$

$$\Leftrightarrow \text{Biais}(\hat{f}_x(x)) = \frac{h^2}{2} f''(x) \mu_2(k) + \int_{\mathbb{R}} k(t) R(x+th) dt$$

$$\text{Car } \int_{\mathbb{R}} k(t) R(x+th) dt = \int_{\mathbb{R}} t^2 k(t) \underbrace{\frac{R(x+th)}{t^2 h^2}}_{\xrightarrow{h \rightarrow 0} 0} dt$$

On peut montrer que :

$$\int_{\mathbb{R}} k(t) R(x+th) dt = o(h^2) \quad (\text{voir annexe 1 \& \grave{a} la fin du cours})$$

d'où :

$$\text{Biais}(\hat{f}_x(x)) = \frac{h^2}{2} f''(x) \mu_2(k) + o(h^2)$$

$$2 - \text{Var}(\hat{f}_x) = \text{var} \left(\frac{1}{m} \sum_{i=1}^m \frac{1}{h} k \left(\frac{X_i - x}{h} \right) \right)$$

$$\stackrel{\text{iid}}{=} \frac{1}{m^2} \sum_{i=1}^m \frac{1}{h^2} \text{var} \left(k \left(\frac{X - x}{h} \right) \right)$$

$$= \frac{1}{m} \text{var} \left(\frac{1}{h} k \left(\frac{X - x}{h} \right) \right)$$

On a :

$$\mathbb{E}_f \left[\frac{1}{h^2} k^2 \left(\frac{X - x}{h} \right) \right] = \int_{-\infty}^{+\infty} \frac{1}{h^2} k^2 \left(\frac{y - x}{h} \right) f_x(y) dy \quad t = \frac{y - x}{h}$$

$$= \int_{-\infty}^{+\infty} \frac{1}{h} k^2(t) f_x(x+th) dt$$

$$= \int_{-\infty}^{+\infty} \frac{1}{h} k^2(t) f_x(x) dt + \int_{-\infty}^{+\infty} \frac{1}{h} k^2(t) (f_x(x+th) - f_x(x)) dt$$

$$= \frac{1}{h} f_x(x) R(k) + \int_{-\infty}^{+\infty} \frac{1}{h} k^2(t) (f_x(x+th) - f_x(x)) dt$$

d'après convergence dominée : $\int_{-\infty}^{+\infty} \frac{1}{h} k^2(t) (f_x(x+th) - f_x(x)) dt = o(1)$
(voir annexe 2 \& la fin du cours).

d'où :

$$\mathbb{E}_f \left[\frac{1}{h^2} k^2 \left(\frac{X - x}{h} \right) \right] = \frac{1}{h} f_x(x) R(k) + \frac{1}{h} o(1)$$

$$\text{D'où } \text{var}_f(\hat{f}_x(x)) = \frac{1}{m} \left[\frac{1}{h} f_x(x) R(k) + o\left(\frac{1}{h}\right) - \left(f_x(x) + \frac{h^2}{2} f_x''(x) \mu_2(k) + o(h^2) \right)^2 \right]$$

$$\begin{aligned} \text{var}_f(\hat{f}_x(x)) &= \frac{1}{m} \left[\frac{1}{h} f_x(x) R(k) + o\left(\frac{1}{h}\right) - \left(f_x(x) + O(h^2) \right)^2 \right] \\ &= \frac{1}{mh} f_x(x) R(k) + o\left(\frac{1}{mh}\right) + O\left(\frac{1}{m}\right) \quad \text{mais } O\left(\frac{1}{m}\right) = o\left(\frac{1}{mh}\right) \\ &= \frac{1}{mh} f_x(x) R(k) + o\left(\frac{1}{mh}\right) \end{aligned}$$

Donc, en déduisant aisément:

$$\text{MSE}(\hat{f}_x(x)) = \frac{h^4}{4} (f_x''(x))^2 \mu_2^2(k) + o(h^4) + \frac{1}{mh} f_x(x) R(k) + o\left(\frac{1}{mh}\right)$$

Remarque: Si $h = h_m \xrightarrow{m \rightarrow +\infty} 0$ et $mh_m \xrightarrow{m \rightarrow +\infty} +\infty$ alors on a $\text{MSE}(\hat{f}_x(x)) \xrightarrow{m \rightarrow +\infty} 0$.

Propriété: Convergence presque-complète

Hypothèse: * S CIR, S compact

* $f_x \in \mathcal{C}^2$ au voisinage de x , (resp au voisinage de S)
 * k a support compact, symétrique, de carré intégrable, d'intégrale 1 et bornée (et éventuellement $|k(x) - k(y)| < C|x-y|^2, \forall x, y \in S$).

* $f_x(x) > 0$ (resp. $\inf_{x \in S} f_x(x) > 0$)

* $h = h_m \xrightarrow{m \rightarrow +\infty} 0$ et $\frac{m}{\ln(m)} \xrightarrow{m \rightarrow +\infty} +\infty$

Si les hypothèses sont valides, alors

$$1 - \hat{f}_x(x) - f_x(x) = o_{p.com}(h^2) + o_{p.com}\left(\sqrt{\frac{\ln(m)}{mh}}\right)$$

$$2 - \sup_{x \in S} |\hat{f}_x(x) - f_x(x)| = o_{p.com}(h^2) + o_{p.com}\left(\sqrt{\frac{\ln(m)}{mh}}\right)$$

Démonstration: admise.

tous les points auraient des poids similaires (\rightarrow estimateur aplatis de la densité). On parle de "surelissage". Si au contraire h est très petit, $\frac{x_i - x}{h}$ sera très grand sauf quand x_i est très proche de x . (on h aurait donc encore un estimateur proche de 0 sauf autour des observations où l'on observerait des pics). On parle de "sous lissage".

Exemple: Retour aux données galaxies.
Essayer les différents noyaux. On choisira comme paramètre de lissage celui qui est le plus approprié pour la fenêtre mobile. Ensuite faire varier à noyau fixé le paramètre de lissage.

Avantages:

- \hookrightarrow estimateur lisse et mon constant par morceaux (si k l'est).
- \hookrightarrow poids des observations dépend de leur proximité par rapport à x .

Inconvénient:

- \hookrightarrow dépend du noyau k
- \hookrightarrow Δ dépend du paramètre de lissage h ! Δ

II - Propriétés et vitesse de convergence

- Propriété:
- 1- $\forall i$ $k \geq 0$ alors $\hat{f}_x \geq 0$
 - 2- $\forall i$ k est une densité, alors \hat{f}_x est une densité.
 - 3- $\forall i$ k est continue alors \hat{f}_x est continue.
 - 4- $\forall i$ k est différentiable, alors \hat{f}_x est différentiable.

Démonstration: trivial!

Propriétés: Si f_x est \mathcal{C}^2 au voisinage de x et bornée (on peut s'en passer si k est à support compact) Si

$h = \frac{1}{m} \xrightarrow{m \rightarrow \infty} 0$ alors:

1 - Bias $(\hat{f}_x(x)) = h^2 \frac{f_x''(x)}{2} \mu_2(k) + o(h^2)$ où $\mu_2(k) = \int_{\mathbb{R}} t^2 k(t) dt < \infty$

2 - Si $f_x(x) > 0$ et $R(k) = \int k^2(t) dt$ alors:

$$MSE(\hat{f}_x(x)) = \frac{1}{mh} f_x(x) R(k) + \frac{h^4}{4} \left(\frac{f_x''(x)}{2} \right)^2 \mu_2^2(k) + o(h^4) + o\left(\frac{1}{mh}\right)$$

Démonstration:

1 - Calcul de l'espérance de $\hat{f}_x(x)$

$$E_p \left[\hat{f}_x(x) \right] = \frac{1}{m} \sum_{i=1}^m \frac{1}{h} E_p \left[K \left(\frac{X_i - x}{h} \right) \right]$$

$$\stackrel{iid}{=} \frac{1}{m} \sum_{i=1}^m E_p \left[\frac{1}{h} K \left(\frac{X - x}{h} \right) \right] = \frac{1}{h} E_p \left[K \left(\frac{X - x}{h} \right) \right]$$

$$= \int_{\mathbb{R}} \frac{1}{h} K \left(\frac{y - x}{h} \right) f_x(y) dy$$

$$= \int_{\mathbb{R}} K(t) f_x(x + th) dt.$$

Or, d'après le développement de Taylor-Young:

$$f_x(x + th) = f_x(x) + th f_x'(x) + \frac{(th)^2}{2} f_x''(x) + R(x + th)$$

où $\frac{|R(x + z)|}{z^2} \xrightarrow{z \rightarrow 0} 0$ où.

$$E_p \left[\hat{f}_x(x) \right] = \int_{\mathbb{R}} k(t) f_x(x) dt + \int_{\mathbb{R}} k(t) th f_x'(x) dt + \int_{\mathbb{R}} k(t) \frac{(th)^2}{2} f_x''(x) dt$$

$$+ \int_{\mathbb{R}} k(t) R(x + th) dt$$

$$= f_x(x) \underbrace{\int_{\mathbb{R}} k(t) dt}_{=1} + h f_x'(x) \underbrace{\int_{\mathbb{R}} t k(t) dt}_{=0} + \frac{h^2}{2} f_x''(x) \underbrace{\int_{\mathbb{R}} t^2 k(t) dt}_{\mu_2(k)}$$

$$+ \int_{\mathbb{R}} k(t) R(x + th) dt$$

Discussion du comportement biais/variance

- Si h diminue, le biais diminue mais la variance augmente
- Si h augmente, le biais augmente mais la variance diminue
- Pour que la variance tende vers 0, il faut $mh_m \xrightarrow{m \rightarrow +\infty} +\infty$
- Pour que le biais tend vers 0, il faut $h_m \xrightarrow{m \rightarrow +\infty} 0$
- Pour une valeur de h et de k fixée:
 - ↳ plus la courbure de la densité ($=f''(x)$) est grande, plus le biais en x est grand
 - ↳ plus la densité est grande, plus la variance sera grande

Proposition : Si f est bornée et \mathcal{C}^2 , $h = h_m \xrightarrow{m \rightarrow +\infty} 0$ et
si $f_x \in L^2(\mathbb{R})$, $f'_x \in L^2(\mathbb{R})$, $f''_x \in L^2(\mathbb{R})$.

Alors :

$$MISE(\hat{f}_x(x)) = \frac{1}{mh} R(k) + \frac{h^4}{4} \int_{\mathbb{R}} (f''_x(x))^2 dx \mu_2^2(k) + o\left(h^4 + \frac{1}{mh}\right)$$

Démonstration : admise.

Maintenant, on cherche un choix optimal de h en fonction des expressions asymptotiques de MSE et MISE (h qui minimisent ces critères).

On note dans ce qui suit AMSE et AMISE les termes asymptotiquement dominants obtenus.

$$\hat{h}_{AMSE} = \underset{h>0}{\operatorname{argmin}} (AMSE(h)) = \underset{h>0}{\operatorname{argmin}} \left(\frac{h^4}{4} \int_{\mathbb{R}} (f''_x(x))^2 dx \mu_2^2(k) + \frac{1}{mh} \int_{\mathbb{R}} f_x(x) R(k) \right)$$

$$\text{Or, } \frac{\partial AMSE(h)}{\partial h} = h^3 \int_{\mathbb{R}} (f''_x(x))^2 dx \mu_2^2(k) - \frac{1}{mh^2} \int_{\mathbb{R}} f_x(x) R(k)$$

$$\text{D'où } \frac{\partial AMSE(h)}{\partial h} = 0 \Leftrightarrow h^5 = \frac{\int_{\mathbb{R}} f_x(x) R(k)}{m \int_{\mathbb{R}} (f''_x(x))^2 dx \mu_2^2(k)} \quad \text{si } \int_{\mathbb{R}} (f''_x(x))^2 dx \neq 0$$

$$\text{D'où } \hat{h}_{AMSE} = m^{-1/5} \sqrt[5]{\frac{\int_{\mathbb{R}} f_x(x) R(k)}{m \int_{\mathbb{R}} (f''_x(x))^2 dx \mu_2^2(k)}}$$

Au passage, on a $\frac{\partial^3 \text{AMSE}(h)}{\partial h^3} \geq 0$.

On obtient de la même manière :

$$\hat{h}_{\text{AMISE}} = m^{-1/5} \sqrt[5]{\frac{R(k)}{\int_{\mathbb{R}} (f_x''(x))^2 dx \mu_2^2(k)}}$$

Remarque : Problème ! Les expressions dépendent de f_x'' et de f_x qui sont inconnues. Une manière de s'en sortir serait de les remplacer par un estimateur (méthode plug-in sur laquelle on reviendra) ce qui est parfois lourd.

On déduit des expressions de \hat{h}_{AMISE} et \hat{h}_{AMISE} les expressions asymptotiques correspondantes (de AMSE et AMISE).

$$\begin{aligned} \bullet \text{AMSE}(\hat{f}_x, \hat{h}_{\text{AMISE}}) &= \frac{1}{m \hat{h}_{\text{AMISE}}} f_x(x) R(k) + \frac{\hat{h}_{\text{AMISE}}^4}{4} (f_x''(x))^2 \mu_2^2(k) \\ &= m^{-1} m^{1/5} \left(\frac{f_x(x) R(k)}{(f_x''(x))^2 \mu_2^2(k)} \right)^{1/5} f_x(x) R(k) + \frac{1}{4} m^{-4/5} \left(\frac{f_x(x) R(k)}{(f_x''(x))^2 \mu_2^2(k)} \right)^{4/5} (f_x''(x))^2 \mu_2^2(k) \\ &= m^{-4/5} (f_x(x) R(k))^{4/5} (f_x''(x))^2 \mu_2^2(k)^{1/5} + \frac{1}{4} m^{-4/5} (f_x(x) R(k))^{4/5} (f_x''(x))^2 \mu_2^2(k)^{1/5} \\ &= \frac{5}{4} m^{-4/5} \left[(f_x(x) R(k))^4 (f_x''(x))^2 \mu_2^2(k) \right]^{1/5} \end{aligned}$$

$$\bullet \text{AMISE}(\hat{f}_x, \hat{h}_{\text{AMISE}}) = \frac{5}{4} m^{-4/5} (R(k)^4 R(f_x'') \mu_2^2(k))^{1/5} \text{ où } R(f_x'') = \int_{\mathbb{R}} (f_x''(x))^2 dx$$

Remarque : Comparaison entre les convergences (en vitesse) optimale :

* Dans le cas paramétrique, on a en général : $E[(\hat{\theta} - \theta)^2] \sim \frac{C}{m}$

et si $\theta \mapsto f_\theta$ est \mathcal{C}^2 et si $\exists h \in L^2(\mathbb{R})$ telle que $\frac{\partial^2 f_\theta}{\partial \theta^2}(x) \leq h(x)$
 on a alors la même vitesse de convergence pour $f_\theta \rightarrow f_0$ ($MSE[\hat{f}_\theta] \sim C_1 m^{-1}$ et $MISE[\hat{f}_\theta] \sim C_2 m^{-1}$).
 * Dans le cas d'un histogramme à classes d'amplitudes égales à h , on a $MSE \sim C_1' m^{-2/3}$ et $MISE \sim C_2' m^{-2/3}$

III. Choix optimal du noyau par rapport à AMSE et de l'AMISE

On voit que le noyau k influence l'expression asymptotique de la MISE et MSE optimales au travers du terme $C(k) = [\mu_2^2(k) R(k)]^{1/5}$

Le noyau optimal est donc celui qui minimise cette quantité tout en vérifiant les hypothèses faites sur le noyau

On montre que ce problème de minimisation sous contrainte a pour solution le noyau d'Epanechnikov (cf Epanechnikov, 1969)

$$k_{opt}(h) = \frac{3}{4} (1-t^2) \mathbb{1}_{\{|t| < 1\}}$$

Afin de comprendre ce que l'on cherche en n'utilisant pas le noyau optimal, on considère la notion d'efficacité d'un noyau k par rapport au noyau optimal k_{opt} qui se définit par :

$$eff(h) = \left(\frac{C(k_{opt})}{C(k)} \right)^{5/4}$$

Remarque : Par définition, $eff(k) \leq 1$

Le rapport s'interprète comme la proportion d'observations

nécessaires pour atteindre avec k_{opt} la même AMISE qu'avec le noyau k .

En effet,

$$AMISE(k) := AMISE(\hat{f}_x, \hat{h}_{AMISE}, k, m)$$

$$AMISE(k_{opt}) := AMISE(\hat{f}_x, \hat{h}_{AMISE}, k_{opt}, m_0)$$

$$AMISE(k) = \frac{5}{4} m^{-4/5} R(f_x'')^{1/5} C(k)$$

$$AMISE(k_{opt}) = \frac{5}{4} m_0^{-4/5} R(f_x'')^{1/5} C(k_{opt})$$

$$AMISE(k) = AMISE(k_{opt}) \Leftrightarrow \left(\frac{m_0}{m}\right)^{4/5} = \frac{C(k_{opt})}{C(k)}$$

$$\Leftrightarrow \frac{m_0}{m} = \left(\frac{C(k_{opt})}{C(k)}\right)^{5/4} = \text{eff}(k) \leq 1.$$

Efficacité relative de quelques noyaux

Noyau	$K(k)$	$\text{eff}(k)$
Epanechnikov	$\frac{3}{4} (1-t^2) \mathbb{1}_{\{ t \leq 1\}}$	1
Quartic	$\frac{15}{16} (1-t^2)^2 \mathbb{1}_{\{ t \leq 1\}}$	0.994
Triweight	$\frac{16}{32} (1-t^2)^3 \mathbb{1}_{\{ t \leq 1\}}$	0.987
Gaussien	$\frac{1}{\sqrt{2\pi}} \exp(-t^2/2)$	0.951
Triangulaire	$\sqrt{2\pi} (1- t) \mathbb{1}_{\{ t \leq 1\}}$	0.986
Uniforme	$\frac{1}{2} \mathbb{1}_{\{ t \leq 1\}}$	0.930

On voit finalement que le choix du noyau n'a pas un si gros impact que cela sur les performances de notre estimateur (sauf peut-être pour uniforme et gaussien). Comme nous le verrons, un mauvais choix de

h est beaucoup plus néfaste qu'un mauvais choix de k .

IV - Choix automatique du paramètre de lissage h .

Nous allons voir maintenant trois méthodes permettant de choisir une valeur "appropriée" du paramètre de lissage en pratique :

- 1- la règle simple de référence à une distribution normale
- 2- la méthode plug-in
- 3- la méthode de validation croisée

1) La règle simple de référence à une distribution normale

Rappelons, comme nous l'avons vu précédemment que

$$\hat{h}_{\text{AISE}} = \left[\frac{R(k)}{R(f_x'') \mu_2^2(k)} \right]^{1/5} n^{-1/5}$$

Si on fait l'hypothèse que la densité que l'on cherche à estimer est de loi normale alors on peut calculer $R(f_x'')$. On a :

$$f_x''(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \left[\frac{(x-\mu)^2}{\sigma^4} - \frac{1}{\sigma^2} \right]$$

D'où :

$$\begin{aligned} R(f_x'') &= \int_{\mathbb{R}} f_x''(x)^2 dx \\ &= \int_{\mathbb{R}} \frac{1}{\sigma^3\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \left(\frac{(x-\mu)^2}{\sigma^4} - \frac{1}{\sigma^2} \right)^2 dx \\ &= \int_{\mathbb{R}} \frac{1}{\sigma^3\sqrt{2\pi}} e^{-\frac{t^2}{2\sigma^2}} \left(\frac{t^2}{\sigma^2} - 1 \right)^2 dt \quad t = x - \mu \\ &= \int_{\mathbb{R}} \frac{1}{2\pi\sigma^6} e^{-\frac{z^2}{2\sigma^2}} \left(\frac{z^2}{2\sigma^2} - 1 \right)^2 \frac{1}{\sqrt{2}} dz \quad z = t\sqrt{2} \end{aligned}$$

$$\begin{aligned}
 R(f_x''') &= \frac{1}{\sqrt{2\pi} \sqrt{2} \sqrt{\sigma}} \int_{\mathbb{R}} \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{z^2}{2\sigma^2}} \left(\frac{z^4}{4\sigma^4} - \frac{z^2}{\sigma^2} + 1 \right) dz \\
 &= \frac{1}{2\sqrt{\pi} \sigma^5} \left[\frac{1}{4} \underbrace{E[U^4]}_{=3} - \underbrace{E[U^2]}_{=1} + 1 \right] \text{ où } U \sim \mathcal{N}(0,1) \\
 &= \frac{1}{2\sigma^5 \sqrt{\pi}} \left[\frac{3}{4} - 1 + 1 \right] \\
 &= \frac{3}{8\sigma^5 \sqrt{\pi}}
 \end{aligned}$$

On obtient alors :

$$\begin{aligned}
 \hat{h}_{NR}^{AFISE} &= \left(\frac{8\sqrt{\pi} \sigma^5}{3} \frac{R(k)}{\mu_2^2(k)} \right)^{1/5} m^{-1/5} \\
 &= \sigma \left(\frac{8\sqrt{\pi} R(k)}{3\mu_2^2(k)} \right)^{1/5} m^{-1/5}
 \end{aligned}$$

On proposera donc de prendre comme paramètre de lissage

$$\hat{h}_{NR} = \left(\frac{8\sqrt{\pi} R(k)}{3\mu_2^2(k)} \right)^{1/5} \hat{\sigma} m^{-1/5} \text{ où } \hat{\sigma} \text{ est un estimateur de l'écart-type } \sigma.$$

Par exemple, $\hat{\sigma} = S_{m-1}$, où $\hat{\sigma} = \min(S_{m-1}, \frac{R}{1.349})$ où R est l'écart interquartiles $Q_{0.75} - Q_{0.25}$ empurqué et S_{m-1} est l'écart-type corrigé.

Remarque : d'où sort $\hat{\sigma} = \min(S_{m-1}, \frac{R}{1.349})$?

Si $X \sim \mathcal{N}(0,1)$, $Q_x(0.75) - Q_x(0.25) = 1.349$

Si $X \sim \mathcal{N}(\mu, 1)$, $Q_x(0.75) - Q_x(0.25) = 1.349$

Si $X \sim \mathcal{N}(\mu, \sigma^2)$, $Q_x(0.75) - Q_x(0.25) = \sigma \times 1.349$.

En effet, si $Z = aX$, $a > 0$

$Q_x(\alpha)$ est solution en x de $F_x(x) = \alpha$.

Donc $Q_z(\alpha)$ solution de $F_z(z) = \alpha = F_x(z/a)$

d'où $\frac{z}{a} = Q_x(\alpha) \Leftrightarrow z = a \times Q_x(\alpha)$.

d'où $Q_z(\alpha) = a Q_x(\alpha)$.

est un estimateur à noyau (4 fois différentiable) de la densité utilisant un paramètre de lissage sous-optimal et construit à partir des $X_{i-1}, X_i, X_{i+1}, \dots, X_n$.

Les méthodes peuvent éventuellement être itérées en utilisant comme paramètre de lissage celui trouvé à l'étape précédente. $h_0 \rightarrow h_1 \rightarrow h_2 \rightarrow \dots$

3) La méthode de validation croisée.

L'idée est assez différente de ce que nous avons fait pour les deux méthodes précédentes. On part pas de l'expression de \hat{h}_{AMISE} mais on estime directement $\text{MISE}(\hat{f}_x) = \int \hat{f}_x^2(x) dx$. Puis on prend comme valeur de h celui qui minimise ce critère.

↳ Estimation de $\text{MISE}(\hat{f}_x) = \int \hat{f}_x^2(x) dx$
 On note $\text{MISE}(h) = \text{MISE}(\hat{f}_x, h)$.

$$\text{MISE}(\hat{f}_x) = \text{IE} \left[\int (\hat{f}_x - f_x)^2(x) dx \right]$$

$$= \text{IE} \left[\int \hat{f}_x^2(x) dx \right] + \int f_x^2(x) dx - 2 \text{IE} \left[\int \hat{f}_x(x) f_x(x) dx \right]$$

On minimise $\text{MISE}(h)$ en h revient à minimiser :

$$\text{MISE}(h) = \int f_x^2(x) dx = \text{IE} \left[\int \hat{f}_x^2(x) dx \right] - 2 \text{IE} \left[\int \hat{f}_x(x) f_x(x) dx \right]$$

Il nous suffit alors de trouver un estimateur sans biais de ce dernier critère :

$$\text{LSCV}(h) = \int \hat{f}_x^2(x) dx - \frac{2}{m} \sum_{i=1}^m \hat{f}_x^{-i}(X_i)$$

$$\text{où } \hat{f}_x^{-i}(X_i) = \frac{1}{m-1} \sum_{j \neq i} \frac{1}{h} K\left(\frac{X_j - X_i}{h}\right)$$

D'autre part, on peut calculer le coefficient $\left(\frac{8\sqrt{\pi}R(k)}{3\mu_2^2(k)}\right)^{1/5}$ en fonction du noyau utilisé :

Noyau	$\frac{8\sqrt{\pi}R(k)}{3\mu_2^2(k)}$
Epanechnikov	2.34
Gaussien	1.06
Triweight	2.78

↳ Avant : facile à calculer (et à implémenter).

↳ Défaut : on part de l'hypothèse que l'on cherche à estimer est de loi normale. Cela ne marchera pas bien si la vraie densité (que l'on estime) est de nature différente.

2) La méthode plug-in.

Elle consiste à estimer la constante $R(f'')$ qui apparaît dans l'expression de \hat{h}_{AMISE} à l'aide d'un estimateur à noyau utilisant un paramètre de lissage sous-optimal.

• Méthode de Sheather et Jones (1991).

$R(\hat{f}_x'') = R(f_x'')$ où \hat{f}_x est un estimateur à noyau de la densité utilisant un noyau deux fois différentiable et un paramètre de lissage h_0 sous-optimal.

• Une autre méthode provient du fait que si

$$\lim_{x \rightarrow \pm\infty} f_x^{(3)}(x) f_x'(x) = \lim_{x \rightarrow \pm\infty} f_x^{(2)}(x) f_x^{(1)}(x) = 0 \quad \text{alors} \quad \mathbb{E}[f_x^{(4)}(x)] = R(f_x^{(2)})$$

On peut ensuite estimer $\mathbb{E}[f_x^{(4)}(x)]$ par $\frac{1}{n} \sum_{i=1}^n (\hat{f}_x^{(4)}(x_i))$ où $\hat{f}_x^{(4)}$

vérifie est un $\hat{f}_x^{(4)}$ de $\mathbb{E}[f_x^{(4)}(x)]$

$\frac{1}{m} \sum_{i=1}^m \hat{f}_x^{-i}(x_i)$ est un estimateur sans biais de type "leave-one-out" de $\mathbb{E} \left[\int \hat{f}_x(x) f_x(x) dx \right]$

Proposition : $\mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \hat{f}_x^{-i}(x_i) \right] = \mathbb{E} \left[\int \hat{f}_x(x) f_x(x) dx \right]$

Démonstration :

$$\begin{aligned} \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \hat{f}_x^{-i}(x_i) \right] &= \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left[\hat{f}_x^{-i}(x_i) \right] \stackrel{iid}{=} \mathbb{E} \left[\hat{f}_x^{-1}(x_1) \right] \\ &= \mathbb{E} \left[\frac{1}{m-1} \sum_{j \neq i} \frac{1}{h} k \left(\frac{X_j - X_1}{h} \right) \right] \stackrel{iid}{=} \frac{1}{m-1} (m-1) \mathbb{E} \left[\frac{1}{h} k \left(\frac{X_2 - X_1}{h} \right) \right] \\ &= \mathbb{E} \left[\frac{1}{h} k \left(\frac{X_2 - X_1}{h} \right) \right] = \int \frac{1}{h} k \left(\frac{x_2 - x_1}{h} \right) f_x(x_1) f_x(x_2) dx_1 dx_2 \end{aligned}$$

Or $\mathbb{E} \left[\int \hat{f}_x(x) f_x(x) dx \right] = \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \int \frac{1}{h} k \left(\frac{X_i - x}{h} \right) f_x(x) dx \right]$

$\stackrel{iid}{=} \mathbb{E} \left[\int \frac{1}{h} k \left(\frac{X_1 - x}{h} \right) f_x(x) dx \right]$
 → on fait $x = X_1 - h \cdot t$ et $k > 0$
 → on regarde les \int et simplifie

$= \mathbb{E} \left[\mathbb{E} \left[\frac{1}{h} k \left(\frac{X_1 - x}{h} \right) \mid X_1 \right] \right] = \mathbb{E} \left[\frac{1}{h} k \left(\frac{X_1 - X_1}{h} \right) \right]$

D'où $\mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \hat{f}_x^{-i}(x_i) \right] = \mathbb{E} \left[\int \hat{f}_x(x) f_x(x) dx \right]$

La méthode "least squares cross validation" consiste donc à prendre $h_{LSCV} = \underset{h>0}{\operatorname{argmin}} LSCV(h)$

V. Estimation du mode

On va proposer une méthode d'estimation non paramétrique du mode θ

Rappel : $\theta \in \operatorname{argmax} \{ t \in \mathbb{R} : f_x(t) \}$

On propose donc de manière assez naturelle l'estimateur du mode suivant

$$\hat{\theta} = \operatorname{argmax} \{t \in \mathbb{R} : \hat{f}_x(t)\}$$

Et argmax n'est pas forcément unique. On prend alors l'une des valeurs qui maximise \hat{f}_x .

En pratique, la méthode la plus basique consiste à estimer \hat{f}_x sur une grille t_1, \dots, t_m de points équidistants d'un pas ε petit et de prendre comme estimateur du mode le point de la grille où \hat{f}_x est maximale

$$\hat{\theta} = \operatorname{argmax} \{j \in \llbracket 1, M \rrbracket, \hat{f}_x(t_j)\}.$$

On peut bien entendu utiliser des méthodes plus raffinées :

↳ Newton-Raphson

↳ méthode du nombre d'or

↳ ...

C - Estimateur à moindres carrés de la fonction de régression

On suppose maintenant que l'on observe pour chaque individu une observation (x_i, y_i) de variables réelles.

On suppose $(x_i, y_i) \stackrel{\text{iid}}{\sim}_{1 \leq i \leq n} (X, Y)$.

I - Définition

On souhaite étudier la manière dont la variable d'intérêt (ou réponse) Y dépend de la variable