

# Module : Statistique Non Paramétrique

Master 2 - Statistique (2021-2022)

Département de Mathématiques

Enseignant : Pr. YAHIA Djabrane

## Contenu du module

Ce cours de SNP contient les sections suivantes :

- Estimation non paramétrique et théorie Asymptotique
- Tests non paramétriques
- Mesures d'association

## Références :

1. A. B. Tsybakov. Introduction à l'estimation non-paramétrique, Springer-Verlag, Berlin, 2004.
2. D. Bosq. Nonparametric statistics for stochastic processes, Springer-Verlag, 1996.
3. Wand, M.P., Jones, M.C. Kernel Smoothing, London : Chapman and Hall, 1995.
4. Silverman, B.W. Density Estimation for Statistics and Data Analysis, London : Chapman and Hall, 1986.

# Chapitre 1

## Généralités sur l'estimation NP

### 1.0.1 Introduction

La statistique paramétrique est le cadre "classique" de la statistique. On dispose d'un échantillon  $X_1, \dots, X_n$  d'observations issu d'une population  $X$ . On veut estimer une fonction ou quantité relative à cette population (moyenne, variance, densité, distribution,...) à partir de l'échantillon  $X_1, \dots, X_n$ . On suppose que la fonction à estimer est connue à un vecteur de paramètres près.

**Exemple 1.0.1** Soit  $(X_1, \dots, X_n)$  échantillon *i.i.d* de distribution  $N(m, \sigma^2)$  : estimer  $m$  et  $\sigma^2$ , cela revient à estimer la densité

$$f(x) = \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-1}{2} \left(\frac{x - m}{\sigma}\right)^2\right).$$

*Il s'agit d'une estimation paramétrique.*

Mais souvent ;

- On ne suppose pas de forme paramétrique pour la fonction à estimer. Exemple : on veut étudier la surface moyenne du logement  $Y$  en fonction du salaire  $X$  :  $m(x) = E[Y|X = x]$
- On s'autorise toutes les formes a priori ou on ne fait aucune hypothèse sur la forme/nature/type de la distribution des variables aléatoires. Exemple :  $\mathcal{F} = \{f : [0, 1] \rightarrow R, \quad f \text{ croissante}\}$ .

C'est le cas de SNP. Par exemple :  $X$  va de densité  $f$  : estimer  $f$ .

## SNP : Quand l'utiliser ?

- Quand on n'arrive pas à ajuster correctement les observations avec une distribution paramétrique,
- Quand on n'a aucune idée de modèle, ou qu'on ne veut pas avoir un a priori sur le modèle,
- Quand on ne sait pas combien de composantes on veut mettre dans un modèle.
- Quand le nombre de variables est trop grand (problème de grande dimension) et qu'un modèle paramétrique est nonutilisable car il aurait de toutes façons trop de paramètres.

## Avantages et inconvénients

1. Moins d'a priori sur les observations,
2. Modèles plus généraux, donc plus robuste.
3. Vitesses de convergence plus lentes, il faut plus de données pour obtenir une précision équivalente.

## 1.1 Estimation de la Fonctions de Répartition

On observe  $X_1, \dots, X_n$  échantillon issu d'une v.a. réelle, de fonction de répartition (fdr)  $F$  :

$$F(x) = P(X \leq x).$$

L'estimateur naturel de la fdr  $F$  est la fdr empirique notée  $F_n$  définie par

$$\begin{aligned} F_n(x) &= \frac{1}{n} \{\text{nbr } X_i \leq x\} = \frac{1}{n} \sum_{i=1}^n 1_{(X_i \leq x)} \\ &= \begin{cases} 0, & X_{(1)} > x \\ i/n, & X_{(i)} \leq x < X_{(i+1)} \\ 0, & X_{(n)} < x \end{cases} \end{aligned}$$

avec  $\min(X_i) = X_{(1)} \leq \dots \leq X_{(n)} = \max(X_i)$  les statistiques d'ordres associées à l'échantillon.

Il est clair que  $F_n$  est un estimateur non paramétrique de la fdr  $F$ .

**Propriété 1.1.1 1) Biais :** on a

$$E(F_n) = E\left(\frac{1}{n}\sum_{i=1}^n 1_{(X_i \leq x)}\right) = P(X \leq x) = F(x)$$

Donc, *Biais*  $(F_n) = E(F_n) - F = 0$ .

2) **Variance :**

$$\begin{aligned} E(F_n^2) &= E\left(\frac{1}{n}\sum_{i=1}^n 1_{(X_i \leq x)}\right)^2 = \frac{1}{n^2}E\left(\sum 1_{(X_i \leq x)}^2 + \sum_{i \neq j} 1_{(X_i \leq x)} 1_{(X_j \leq x)}\right) \\ &= \frac{1}{n}F(x) + \frac{n-1}{n}F^2(x). \end{aligned}$$

D'où,

$$\text{Var}(F_n) = E(F_n^2) - E^2(F_n) = \frac{1}{n}F(x)(1 - F(x)).$$

3) Comme  $\text{Var}(F_n) \rightarrow 0$ , alors  $F_n \rightarrow F$  en Probabilités.

4) D'après le Théorème de la limite centrale :

$$\sqrt{n} \left( \frac{F_n(x) - F(x)}{\sqrt{F(x)(1 - F(x))}} \right) \rightarrow N(0, 1) \text{ en loi.}$$

5) De plus, d'après le théorème de Glivenko-Contelli :

$$\sup_x |F_n(x) - F(x)| \rightarrow 0 \text{ P.s.}$$

**Exercice 1.1.1** Soit  $(X_n)$  une suite de variables aléatoires indépendante et identiquement distribués, de fonction de survie notée  $G$ , telle que  $G(x) = 1 - F(x) = P(X > x)$ , où  $F$  est la fonction de distribution de  $X$ . On considère la fonction empirique  $G_n$  définie par :

$$G_n(t) = \frac{1}{n} \sum_{j=1}^n 1_{(X_j > t)}, \quad t \in \mathbb{R}$$

1) Quelle est la loi de  $nG_n(t)$  ? la loi limite de  $\sqrt{n}G_n(t)$  ?

2) Montrer la convergence en probabilités de  $G_n(t)$  vers  $G(t)$  lorsque  $n \rightarrow \infty$ .

3) Montrer que cette convergence est aussi presque sûre et on norme infinie.

4) Montrer que la quantité  $\|G_n(t) - G(t)\|_\infty$  est indépendante de  $G(t)$ .

**Exercice 1.1.2** Soit  $F$  une distribution sur  $R$  et soit  $\theta \in R_+$  un paramètre inconnu. On dispose d'un échantillon  $X_1, \dots, X_n$  de fonction de répartition

$$P(X \leq x) = F_\theta(x) = F(x - \theta).$$

Considérons la variable aléatoire  $Y_n = \sum_{i=1}^n 1_{(X_i > 0)}$ .

1) Pour  $n \in \mathbb{N}^*$  fixé, montrer que  $Y_n$  suit une loi Binomiale de paramètres  $n, p$  à préciser.

2) Montrer que la loi limite de  $\frac{1}{\sqrt{n}}(Y_n - np)$  est gaussienne. Préciser la moyenne et la variance de cette loi limite.

3) Déterminer la loi limite des deux statistiques suivantes :

$$V_n = \left(\frac{Y_n}{n}\right)^2 \quad \text{et} \quad W_n = \frac{n}{Y_n}.$$

# Chapitre 2

## Estimateur NP de la Densité

Soit  $X$  une variable aléatoire de densité de probabilité inconnue  $f$ . Supposons que nous avons  $n$  observations  $X_1, X_2, \dots, X_n$  provenant de  $X$ . Le problème consiste à trouver un estimateur pour la fonction  $f$  à partir de cet échantillon issu de  $X$ . Pour cela, l'approche non paramétrique est plus adéquate lorsqu'on ne possède aucune information précise sur la forme et la classe de la vraie densité. Dans cette approche, ce sont les observations qui vont nous permettre de déterminer un estimateur pour la densité  $f$ .

On observe  $X_1, \dots, X_n$  échantillon issu d'une v.a. réelle  $X$ , de fonction de répartition (fdr)  $F$  :

$$F(x) = P(X \leq x)$$

et de densité de probabilité  $f$  :

$$P(X \in [a, b]) = \int_a^b f(x) dx, \quad \forall a, b \in \mathbb{R}.$$

Nous supposons de plus, que  $f$  est deux fois continûment différentiable.

## 2.1 Estimation de la densité par histogramme

En statistique, l'histogramme est une représentation graphique de la répartition d'une variable aléatoire  $X$  (Pearson, 1895). Supposons que  $f$  est à support compact inclus dans  $[0, 1]$ . Soit  $C_1, \dots, C_m$  une partition uniforme de  $[0, 1[$ :

$$C_k = \left[ \frac{k-1}{m}, \frac{k}{m} \right[, \quad k = 1, \dots, m.$$

**Remarque 2.1.1** Dans le cas général où  $f$  est à support  $[a, b]$ , on peut poser

$$\forall m \in \mathbb{N}^*, \quad h = \frac{b-a}{m}.$$

On définit alors les classes  $C_k$  comme suit :

$$C_k = [a + (k-1)h, a + kh[ \quad \text{pour } j = 1, \dots, m-1 \quad \text{et } C_m = [b-h, b].$$

Il est clair que  $f$  est bien approchée par des fonctions en escalier, constantes par morceaux sur les intervalles  $C_j$ . Posons  $h = \frac{1}{m}$  et on approche  $f$  par la fonction

$$f_h(x) = \sum_{k=1}^m \frac{p_k}{h} I_{C_k(x)}.$$

avec  $p_k = \int_{C_k} f(x) dx = E(I_{C_k(X)})$ . Alors, il est naturel d'estimer  $p$  par

$$\hat{p}_k = \frac{1}{n} \sum_{j=1}^n I_{C_k(X_j)}, \quad k = 1, \dots, m.$$

Observons que chaque  $\hat{p}_k$  représente la proportion des observations  $X_j$  se trouvant dans l'intervalle  $C_k$ . Par substitution, nous définissons l'estimateur de  $f$  par histogramme à  $m$  classes comme suit :

$$\hat{f}_h(x) = \sum_{k=1}^m \frac{\hat{p}_k}{h} I_{C_k(x)}. \quad (2.1)$$

**Remarque 2.1.2** On dit que chaque  $C_k$  est une classe de longueur (ou fenêtre)  $h$ . La hauteur des

rectangles représente les fréquences absolues (nombre d'observations dans chaque classe) ou bien il s'agit des fréquences relatives comme dans la figure suivante.

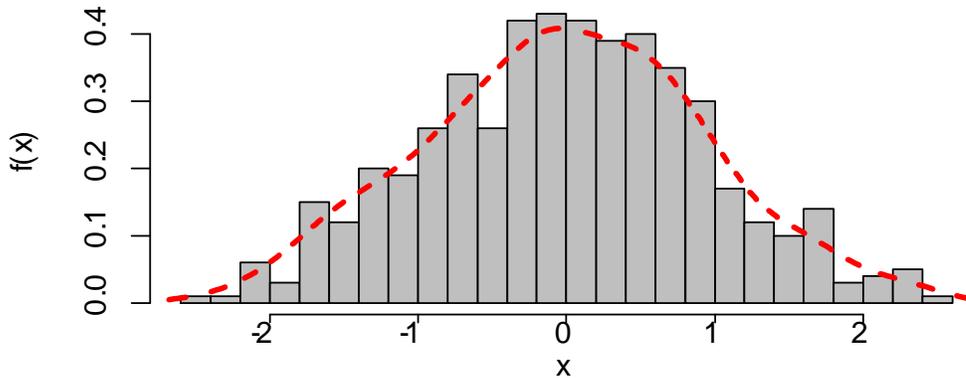


FIG. 2.1 – Estimation par histogramme basée sur un échantillon de taille  $n = 500$ ,  $X \sim N(0, 1)$ ,  $m = 30$ ,  $h \simeq 0.2$

### 2.1.1 Propriétés asymptotique de l’histogramme

Il est clair que la qualité d’ajustement par histogramme dépend fortement de la fenêtre  $h$ . Il est naturel donc, d’étudier le risque quadratique de  $\hat{f}_h$  au point  $x \in [0, 1]$  comme étant l’erreur quadratique moyenne ( $MSE$ ) définit par :

$$\begin{aligned} MSE(\hat{f}_h) &= E\left(\left(\hat{f}_h(x) - f(x)\right)^2\right) = \left(E\hat{f}_h(x) - f(x)\right)^2 + Var\left(\hat{f}_h(x)\right) \\ &= Bias^2\left(\hat{f}_h(x)\right) + Var\left(\hat{f}_h(x)\right) \end{aligned}$$

Pour tout  $x \in C_k$ , on a

$$\hat{f}_h(x) = \frac{\hat{p}_k}{h} = \frac{1}{nh} \sum_{j=1}^n I_{C_k}(X_j) = \frac{W_k}{nh},$$

avec  $W_k$  est la somme de  $n$  variables indépendantes de loi de Bernoulli de paramètre  $p_k$  :

$$P(I_{C_k}(X_j) = 1) = P(X_j \in C_k) = \int_{C_k} f(x) dx = p_k.$$

Donc,  $\forall x \in C_k$  :

$$E\left(\hat{f}_h(x)\right) = \frac{p_k}{h} \text{ et } Var\left(\hat{f}_h(x)\right) = \frac{np_k(1-p_k)}{n^2h^2} = \frac{p_k(1-p_k)}{nh^2}.$$

En déduit de ce dernier résultat que

$$MSE\left(\hat{f}_h\right) = \left(\frac{p_k}{h} - f(x)\right)^2 + \frac{p_k(1-p_k)}{nh^2}. \quad (2.2)$$

Afin d'avoir une évaluation globale valable pour tout point  $x \in [0, 1]$ , on considère le risque quadratique intégré moyen (*MISE*) :

$$\begin{aligned} MISE\left(\hat{f}_h\right) &= \int MSE\left(\hat{f}_h(x)\right) dx = E\left(\int \left(\hat{f}_h(x) - f(x)\right)^2\right) \\ &= \int \text{Biais}^2\left(\hat{f}_h(x)\right) dx + \int Var\left(\hat{f}_h(x)\right) dx. \end{aligned}$$

Premièrement, pour le biais, on a

$$\begin{aligned} \int \text{Biais}^2\left(\hat{f}_h(x)\right) dx &= \sum_{k=1}^m \int_{C_k} \left(\frac{p_k}{h} - f(x)\right)^2 dx \\ &= \sum_{k=1}^m \frac{p_k^2}{h^2} \int_{C_k} dx - 2 \sum_{k=1}^m \frac{p_k}{h} \int_{C_k} f(x) dx + \sum_{k=1}^m \int_{C_k} f^2(x) dx \\ &= \sum_{k=1}^m \frac{p_k^2}{h} - 2 \sum_{k=1}^m \frac{p_k^2}{h} + \int_0^1 f^2(x) dx, \end{aligned}$$

puisque, on a

$$\int_{C_k = [\frac{k-1}{m}, \frac{k}{m}[} dx = \frac{1}{m} = h, \quad \int_{C_k} f(x) dx = p_k \text{ et } \sum_{k=1}^m \int_{C_k} f^2(x) dx = \int_0^1 f^2(x) dx.$$

Nous obtenons donc,

$$\int \text{Biais}^2\left(\hat{f}_h(x)\right) dx = \int \left(E\hat{f}_h(x) - f(x)\right)^2 dx = \int f^2(x) dx - \frac{1}{h} \sum_{k=1}^m p_k^2. \quad (2.3)$$

Pour le terme de la variance, on a

$$\begin{aligned}
\int \text{Var}(\hat{f}_h(x)) dx &= \sum_{k=1}^m \int_{C_k} \text{Var}(\hat{f}_h(x)) dx \\
&= \sum_{k=1}^m \frac{p_k(1-p_k)}{nh^2} \int_{C_k} dx = \frac{1}{nh} \sum_{k=1}^m p_k(1-p_k) \\
&= \frac{1}{nh} \sum_{k=1}^m p_k - \frac{1}{nh} \sum_{k=1}^m p_k^2 \\
&= \frac{1}{nh} \sum_{k=1}^m \int_{C_k} f(x) dx - \frac{1}{nh} \sum_{k=1}^m p_k^2 \\
&= \frac{1}{nh} \int f(x) dx - \frac{1}{nh} \sum_{k=1}^m p_k^2 = \frac{1}{nh} - \frac{1}{nh} \sum_{k=1}^m p_k^2. \tag{2.4}
\end{aligned}$$

**Théorème 2.1.1** Soit  $X_1, \dots, X_n$  un échantillon issu d'une v.a. réelle  $X$ , de densité  $f$  et  $\hat{f}_h$  est l'estimateur de  $f$  par histogramme, où  $m = 1/h$  classes, alors :

$$\text{MISE}(\hat{f}_h) = \int f^2(x) dx + \frac{1}{nh} - \frac{1}{h} \left(1 + \frac{1}{n}\right) \sum_{k=1}^m p_k^2.$$

**Preuve.** La preuve du résultat est détaillée dans (2.3) et (2.4). ■

**Théorème 2.1.2** Supposons que la densité  $f$  de  $X$  est deux fois continûment différentiable et s'annule en dehors de l'intervalle  $[0, 1]$ . Sous la condition que

$$h := h_n, \text{ et } h \rightarrow 0 \text{ quand } n \rightarrow \infty$$

Alors, lorsque  $n \rightarrow \infty$ ,

$$\text{MISE}(\hat{f}_h) = \frac{h^2}{12} \int f'(x)^2 dx + \frac{1}{nh} + O(h^3) + O\left(\frac{1}{n}\right). \tag{2.5}$$

**Preuve.** Du Théorème 1, on a

$$\begin{aligned}
\text{MISE}(\hat{f}_h) &= \int f^2(x) dx + \frac{1}{nh} - \frac{1}{h} \left(1 + \frac{1}{n}\right) \sum_{k=1}^m p_k^2 \\
&= \sum_{k=1}^m \left\{ \int_{C_k} f^2(x) dx - \frac{p_k^2}{h} \right\} + \frac{1}{nh} \left(1 - \sum_{k=1}^m p_k^2\right)
\end{aligned}$$

Il est clair que,

$$\begin{aligned}
\int_{C_k} f^2(x) dx - \frac{p_k^2}{h} &= \int_{C_k} f^2(x) dx - \frac{1}{h} \left( \int_{C_k} f(x) dx \right)^2 \\
&= \int_{C_k} f^2(x) dx - \frac{2}{h} \left( \int_{C_k} f(x) dx \right)^2 - \frac{1}{h} \left( \int_{C_k} f(x) dx \right)^2 \\
&= \int_{C_k} f^2(x) dx - \frac{2}{h} \int_{C_k} f(x) dx \int_{C_k} f(t) dt + \frac{1}{h^2} \int_{C_k} \left( \int_{C_k} f(x) dx \right)^2 dx
\end{aligned}$$

donc

$$\begin{aligned}
\int_{C_k} f^2(x) dx - \frac{p_k^2}{h} &= \int_{C_k} \left\{ f^2(x) - \frac{2}{h} f(x) \int_{C_k} f(t) dt + \frac{1}{h^2} \left( \int_{C_k} f(x) dx \right)^2 \right\} dx \\
&= \int_{C_k} \left( f(x) - \frac{1}{h} \int_{C_k} f(t) dt \right)^2 dx = \int_{C_k} \left( \frac{1}{h} \int_{C_k} f(x) dt - \frac{1}{h} \int_{C_k} f(t) dt \right)^2 dx \\
&= \frac{1}{h^2} \int_{C_k} \left( \int_{C_k} \{f(x) - f(t)\} dt \right)^2 dx,
\end{aligned}$$

Puisque  $f$  est deux fois continûment différentiable, alors pour tout  $x, t \in C_k$  :

$$f(x) - f(t) = (x - t) f'(\theta_k) + O(h^2),$$

avec  $\theta_k$  désigne l'extrémité gauche de l'intervalle  $C_k$ . Donc,

$$\begin{aligned}
\int_{C_k} f^2(x) dx - \frac{p_k^2}{h} &= \frac{1}{h^2} \int_{C_k} \left( \int_{C_k} \{(x - t) f'(\theta_k) + O(h^2)\} dt \right)^2 dx \\
&= \frac{1}{h^2} f'(\theta_k)^2 \int_{C_k} \left( \int_{C_k} (x - t) dt \right)^2 dx + O(h^4).
\end{aligned}$$

En utilisant le changement de variable  $(x, t) = (\theta_k + yh, \theta_k + zh)$ , on obtient

$$x, t \in [\theta_k, \theta_k + h[, \quad dx = hdy, \quad dt = hdz \quad \text{et} \quad y, z \in [0, 1[,$$

donc

$$\begin{aligned} \int_{C_k} \left( \int_{C_k} (x-t) dt \right)^2 dx &= \int_{C_k} \left( \int_{C_k} x dt - \int_{C_k} t dt \right)^2 dx \\ &= h^5 \int \left( \int (y-z) dz \right)^2 dy = \frac{h^5}{12}. \end{aligned}$$

Par conséquence,

$$\int_{C_k} f^2(x) dx - \frac{p_k^2}{h} = \frac{h^3}{12} f'(\theta_k)^2 + O(h^4) = \frac{h^2}{12} \int_{C_k} f'(x)^2 dx + O(h^4).$$

Donc,

$$\begin{aligned} MISE(\hat{f}_h) &= \sum_{kj=1}^m \left( \int_{C_k} f^2(x) dx - \frac{p_k^2}{h} \right) + \frac{1}{nh} - \frac{1}{nh} \sum_{kj=1}^m p_k^2 \\ &= \frac{h^2}{12} \int f'(x)^2 dx + \frac{1}{nh} + O(h^3) + O\left(\frac{1}{n}\right). \end{aligned}$$

■

**Corollaire 2.1.1** *Ce résultat nous permet de calculer la fenêtre  $h$  optimale notée  $h_{opt}$ , en minimisant la MISE asymptotique :*

$$\begin{aligned} h_{opt} &= \arg \min_h AMISE(\hat{f}_h) = \arg \min_h \frac{h^2}{12} \int f'(x)^2 dx + \frac{1}{nh} \\ &= \left( \frac{n}{6} \int f'(x)^2 dx \right)^{-1/3} \simeq Cn^{-1/3}. \end{aligned}$$

**Remarque 2.1.3** *Cette fenêtre optimale est en général incalculable (donc inutilisable, du point de vue pratique), car la densité  $f$  (ainsi que sa dérivée  $f'$ ) est inconnue. Cependant, cette fenêtre optimale est de l'ordre de  $n^{-1/3}$  pour  $n$  assez grand.*

## 2.2 Estimation de la densité par noyau

Par construction, les histogrammes ne sont pas des fonctions continues, ni lisse. Si on cherche donc, à estimer une fonction densité (supposée deux fois continûment différentiable), il sera plus naturel que l'estimateur possède les mêmes propriétés de continuité et de différentiabilité et lisse en plus, donc plus proche de la vraie densité que l'estimateur par histogramme.

Soit  $x \in R$  et  $h > 0$ . Si l'on suppose que  $x$  est le centre d'une classe de l'histogramme et que  $h$  est la longueur des classes, l'estimateur de la densité  $f(x)$  par histogramme peut s'écrire :

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{j=1}^n I_{(|x-X_j| \leq \frac{h}{2})} = \frac{1}{nh} \sum_{j=1}^n I_{(|\frac{x-X_j}{h}| \leq \frac{1}{2})}.$$

Une façon de généraliser les histogrammes consiste donc à utiliser la formule ci-dessus pour tout  $x \in R$  et pas seulement pour les centres des classes. Par conséquent, en remplaçant  $I_{(|z| \leq \frac{1}{2})}$  par une fonction quelconque  $K$ , on obtient l'estimateur

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{j=1}^n K\left(\frac{x - X_j}{h}\right).$$

A partir de la définition d'une densité de probabilité et en utilisant la distribution empirique comme estimateur NP de la distribution  $F$ , on aura pour  $h$  assez petite ( $h \rightarrow 0$ , quand  $n \rightarrow \infty$ ) :

$$f(x) = \frac{F(x+h) - F(x-h)}{2h},$$

ce qui implique

$$\begin{aligned}
f_n(x) &= \frac{F_n(x+h) - F_n(x-h)}{2h} \\
&= \frac{1}{2h} \left( \frac{1}{n} \sum_{j=1}^n I_{(X_j \leq x+h)} - \frac{1}{n} \sum_{j=1}^n I_{(X_j \leq x-h)} \right) \\
&= \frac{1}{2nh} \sum_{j=1}^n \left( I_{\left(\frac{X_j-x}{h} \leq 1\right)} - I_{\left(\frac{X_j-x}{h} \leq -1\right)} \right) = \frac{1}{2nh} \sum_{j=1}^n I_{\left(-1 \leq \frac{X_j-x}{h} \leq 1\right)} \\
&= \frac{1}{nh} \sum_{j=1}^n \frac{1}{2} I_{\left(-|\frac{X_j-x}{h}| \leq 1\right)}.
\end{aligned}$$

Cette dernière peut être réécrite, en ses points de continuité, sous la forme suivante :

$$f_n(x) = \frac{1}{nh} \sum_{j=1}^n K\left(\frac{x - X_j}{h}\right) \quad (2.6)$$

avec  $h := h_n$  est le paramètre de lissage (fenêtre ou bandwidth en anglais) choisi en fonction de  $n$  telle que

$$\lim_{n \rightarrow \infty} h = 0,$$

et  $K$  est la fonction des poids (noyau, kernel en anglais) où :

$$K(t) := \frac{1}{2} I_{(|t| < 1)}.$$

Ce dernier est l'estimateur à noyau uniforme dit de Rosenblatt, proposé pour la première fois en 1956 par Rosenblatt. Six ans après, cet estimateur a été généralisé par Parzen (1962). A partir de cette date, cet estimateur a pris le nom de l'estimateur de Parzen-Rosenblatt. L'idée de l'estimateur par la méthode du noyau consiste donc, à évaluer la densité  $f(x)$  au point  $x$  en comptant le nombre d'observations tombées dans un certain voisinage de  $x$  sur  $\mathbb{R}$ .

**Définition 2.2.1** Soit  $K : \mathbb{R} \rightarrow \mathbb{R}$  une fonction (un noyau),  $h > 0$  (une fenêtre), l'estimateur à noyau de la densité  $f$  est

$$f_n(x) = \frac{1}{nh} \sum_{j=1}^n K\left(\frac{x - X_j}{h}\right).$$

Les noyaux les plus utilisés dans l'estimation de la densité de probabilité sont donnés dans le tableau suivant :

TAB. 2.1 – Noyaux usuels	
Noyau	Fonction $K(t)$
Rectangulaire	$\frac{1}{2}1_{( t <1)}$
Triangulaire	$(1 -  t )1_{( t <1)}$
Gaussien	$\frac{1}{\sqrt{2\pi}}e^{-\frac{t^2}{2}}, t \in R$
Quartique	$\frac{15}{16}(1 - t^2)^2 1_{( t <1)}$
Epanechnikov	$\frac{3}{4}(1 - t^2)1_{( t <1)}$

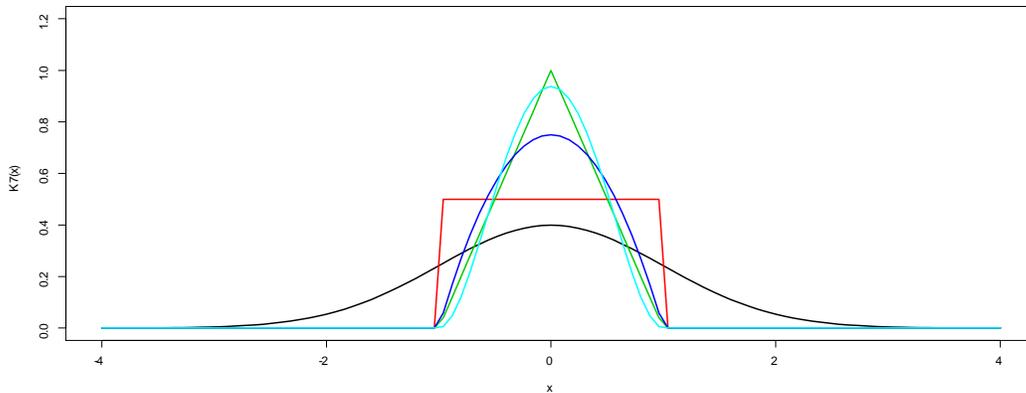


FIG. 2.2 – Courbes des noyaux usuels

**Proposition 2.2.1** *Si  $K$  est un noyau positif et  $\int K(t) dt = 1$ . Alors,  $f_n$  est une densité de probabilité. De plus,  $f_n$  a les mêmes propriétés de continuité et de différentiabilité que  $K$ .*

**Preuve.** On a,  $f_n(x) > 0$  et

$$\begin{aligned} \int f_n(x) dx &= \frac{1}{nh} \int \sum_{j=1}^n K\left(\frac{x - X_j}{h}\right) dx = \frac{1}{nh} \sum_{j=1}^n \int hK(t) dt \\ &= \int K(t) dt = 1. \end{aligned}$$

avec le ch.v ( $x = X_j + th$ ). Donc  $f_n$  est bien une densité de probabilité. ■

## 2.2.1 Propriétés Asymptotique

Supposons que le noyau  $K$  vérifie les conditions suivantes :

(K1)  $\int K(t)dt = 1.$

(K2)  $K$  symétrique autour de zéro, c.à.d  $\int tK(t)dt = 0.$

(K3)  $K$  possède un moment d'ordre 2 fini, c.à.d  $\int t^2K(t)dt < \infty.$

(K4)  $\int t^2 |K(t)| dt < \infty.$

**Proposition 2.2.2** *Si la densité  $f$  est bornée et que  $f''$  existe et bornée aussi. Sous (K1-K3) :*

$$|\text{Biais}(f_n)| \leq C_1 h^2, \quad \text{avec } C_1 = \frac{1}{2} \sup_z |f''(z)| \int t^2 |K(t)| dt.$$

*Si, de plus, la condition (K4) est satisfaite, alors*

$$\text{Var}(f_n) \leq \frac{C_2}{nh}, \quad \text{avec } C_2 = \sup_z f(z) \int \int K(t)^2 dt.$$

**Preuve.** Pour le biais, on a

$$\begin{aligned} E(f_n) &= E\left(\frac{1}{nh} \sum_{j=1}^n K\left(\frac{x - X_j}{h}\right)\right) = \frac{1}{nh} \sum_{j=1}^n E\left(K\left(\frac{x - X_j}{h}\right)\right) \\ &= \frac{1}{nh} \sum \int K\left(\frac{y - x}{h}\right) f(y) dy \\ &= \int K(t) f(x + th) dt \quad (\text{ch.v : } y = x + th, \quad dy = hdt) \end{aligned}$$

En effectuons un développement limité d'ordre 2 à  $f$ , il vient

$$\begin{aligned} E(f_n) &= \int K(t) \left\{ f(x) + thf'(x) + \frac{t^2 h^2}{2} f''(\theta) \right\} dt \quad \theta \in [x, x + th] \\ &= f(x) \int K(t) dt + hf'(x) \int tK(t) dt + \frac{h^2}{2} \int t^2 f''(\theta) K(t) dt \\ &= f(x) + \frac{h^2}{2} \int t^2 f''(\theta) K(t) dt. \end{aligned}$$

Alors,

$$\begin{aligned}
 |\text{Biais}(f_n)| &= |E(f_n) - f| = \frac{h^2}{2} \left| \int t^2 f''(\theta) K(t) dt. \right| \\
 &\leq \frac{h^2}{2} \int |t^2 f''(\theta) K(t)| dt. \\
 &\leq \frac{h^2}{2} \sup_z |f''(z)| \int t^2 |K(t)| dt =: C_1 h^2.
 \end{aligned}$$

Deuxièmement, les variables aléatoires  $Y_i = K\left(\frac{x-X_j}{h}\right)$  sont i.i.d. donc

$$\begin{aligned}
 \text{Var}(f_n) &= \frac{1}{n^2 h^2} \sum_{j=1}^n \text{Var}\left(K\left(\frac{x-X_j}{h}\right)\right) \\
 &\leq \frac{1}{nh^2} E\left(K\left(\frac{x-X_j}{h}\right)^2\right)
 \end{aligned}$$

(ch.v :  $y = x + th, \quad dy = hdt$ ),

$$\begin{aligned}
 \text{Var}(f_n) &\leq \frac{1}{nh^2} \int K^2\left(\frac{y-x}{h}\right) f(y) dy \\
 &= \frac{1}{nh} \int K^2(t) f(x+th) dt \\
 &\leq \frac{1}{nh} \sup_z f(z) \int \int K(t)^2 dt =: \frac{C_2}{nh}.
 \end{aligned}$$

■

**Remarque 2.2.1** On déduit de cette Proposition que le risque MSE de  $f_n$  admet la majoration suivante :

$$\text{MSE}(f_n) = \text{Biais}^2(f_n) + \text{Var}(f_n) \leq C_1^2 h^4 + \frac{C_2}{nh}.$$

La valeur de la fenêtre  $h$  qui minimise ce majorant du MSE est

$$h_{\text{opt}} = (C_2/4C_1^2)^{1/5} n^{-1/5}.$$

En injectant cette valeur dans l'expression du MSE on obtient :  $\text{MSE}(f_n) \leq C.n^{-4/5}$ . Cela montre

que la vitesse de convergence de l'estimateur à noyau est de  $n^{-4/5}$ . Elle est donc meilleure que la vitesse  $n^{-2/3}$  obtenue pour les histogrammes.

**Remarque 2.2.2** Lorsque la fenêtre  $h$  est très petit ( $h \rightarrow 0$ ), le biais de l'estimateur à noyau est très petit face à sa variance et c'est cette dernière qui détermine la vitesse de convergence du risque quadratique. Dans ce type de situation, l'estimateur est très volatile et on parle de sous-lissage (*under-smoothing*, en anglais). En revanche, lorsque  $h$  grandit, la variance devient petite et c'est le biais qui devient dominant. L'estimateur est alors très peu variable et est de moins à moins influencé par les données. On parle alors d'un effet de sur-lissage (*over-smoothing* en anglais). En pratique, il est primordial de trouver la bonne dose de lissage qui permet d'éviter le sous-lissage et le sur-lissage.

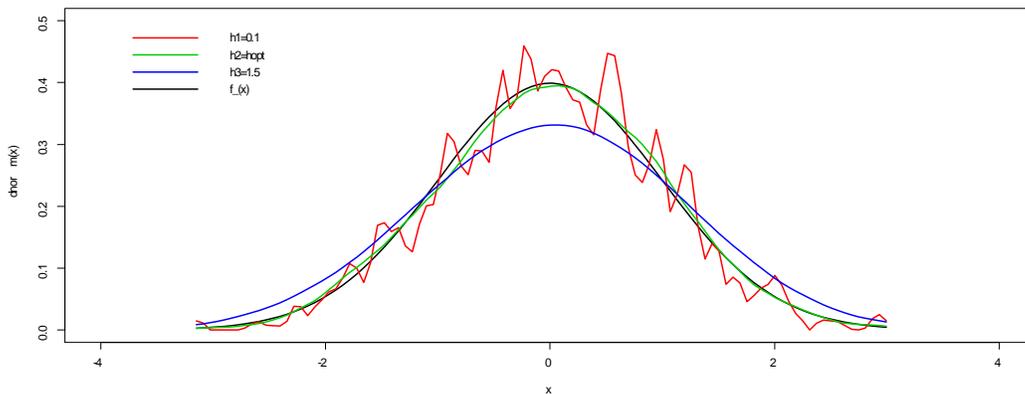


FIG. 2.3 – Biais, Var et choix de  $h$  : sur et sous lissage.

**Proposition 2.2.3** Supposons que  $f \in C^2$  de carré intégrable,  $h = h_n$ , telle que

$$h \rightarrow 0 \quad \text{et} \quad nh \rightarrow \infty, \quad \text{quand } n \rightarrow \infty.$$

De plus, la fonction noyau  $K$  est supposée (densité, bornée, symétrique, de moment d'ordre 4 fini).

Alors,

$$\begin{aligned} \text{i) } \text{Biais}(f_n(x)) &= E(f_n(x)) - f(x) = \frac{h^2}{2} \mu_2(K) f''(x) + o(h^2), \\ \text{ii) } \text{Var}(f_n(x)) &= \frac{1}{nh} f(x) R(K) + o\left(\frac{1}{nh}\right), \end{aligned}$$

avec  $\mu_2(K) := \int t^2 K(t) dt$  et  $R(g) := \int g^2(t) dt$ .

**Preuve.** Pour i), on a

$$\begin{aligned} E(f_n) &= \int K(t) f(x+th) dt \\ &= \int K(t) \left\{ f(x) + th f'(x) + \frac{t^2 h^2}{2} f''(x) + o(h^2) \right\} dt \quad \text{DT au voisinage de } x \\ &= f(x) + \frac{h^2}{2} f''(x) \int t^2 K(t) dt + o(h^2). \end{aligned}$$

De même, pour ii)

$$\begin{aligned} \text{Var}(f_n(x)) &= \frac{1}{n^2 h^2} \sum_{j=1}^n \text{Var} \left( K \left( \frac{x - X_j}{h} \right) \right) \\ &= \frac{1}{nh^2} E \left( K \left( \frac{x - X_j}{h} \right)^2 \right) - \frac{1}{nh^2} E \left( K \left( \frac{x - X_j}{h} \right) \right)^2 \\ &= \frac{1}{nh^2} \int K^2 \left( \frac{y-x}{h} \right) f(y) dy - \frac{1}{n} \left( \frac{1}{h} \int K \left( \frac{y-x}{h} \right) f(y) dy \right)^2 \end{aligned}$$

En utilisant le changement de variable ( $y = x + th$ ,  $dy = h dt$ ) et un développement de Taylor de  $f$  au voisinage de  $x$  d'ordre 1 :

$$\begin{aligned} \text{Var}(f_n(x)) &= \frac{1}{nh} \int K^2(t) f(x+th) dt - \frac{1}{n} \left( \int K(t) f(x+th) dt \right)^2 \\ &= \frac{1}{nh} \int K^2(t) (f(x) + o(1)) dt - \frac{1}{n} \left( \int K(t) (f(x) + o(1)) dt \right)^2 \\ &= \frac{1}{nh} f(x) \int K^2(t) dt + o(1/nh) - o(1/n) \\ &= \frac{1}{nh} f(x) R(K) + o(1/nh), \end{aligned}$$

car  $1/n = o(1/nh)$ . Le résultat est donc démontré. ■

**Notation 2.2.1** 1)  $A_n = o(B_n) \Leftrightarrow \frac{A_n}{B_n} \rightarrow 0$

2)  $A_n = O(B_n) \Leftrightarrow \lim \frac{A_n}{B_n} < \infty$ .

**Remarque 2.2.3** L'estimateur  $f_n(x)$  est asymptotiquement sans biais, i.e.,

$$\lim_{n \rightarrow \infty} E(\hat{f}_n(x)) = f(x).$$

## 2.2.2 Erreur quadratique moyenne (MSE) et intégré (MISE)

Sous les conditions sur  $K$  et en supposant que la densité de probabilité  $f$  avait toutes les dérivées (continues) nécessaires. On peut obtenir facilement les approximations suivantes pour la  $MSE$  et la  $MISE$  :

$$\begin{aligned} MSE(f_n(x)) &= E[(f(x) - f_n(x))^2] = \text{Biais}^2(f_n(x)) + \text{Var}(f_n(x)) \\ &= \frac{h^4}{4} \mu_2^2(K) f''(x)^2 + \frac{1}{nh} f(x) R(K) + o\left(\frac{1}{nh}\right) + o(h^2). \end{aligned}$$

D'ou la MSE asymptotique :

$$AMSE(f_n(x)) = \frac{h^4}{4} \mu_2^2(K) f''(x)^2 + \frac{1}{nh} f(x) R(K).$$

Par minimisation en  $h$  de la  $AMSE$ , on trouve la fenêtre optimale locale, notée  $h_{opt}$  :

$$\begin{aligned} h_{opt} &= \text{Arg} \min_h AMSE(f_n(x)) \\ &= \left( \frac{f(x) R(K)}{\mu_2^2(K) f''(x)^2} \right)^{1/5} n^{-1/5} \simeq C n^{-1/5}. \end{aligned} \tag{2.7}$$

Par intégration, la AMSE intégrée notée AMISE est donnée par :

$$\begin{aligned}
 AMISE(f_n(x)) &= \int AMSE(f_n(x))dx \\
 &= \frac{h^4}{4}\mu_2^2(K) \int f''(x)^2 dx + \frac{1}{nh}R(K) \\
 &= \frac{h^4}{4}\mu_2^2(K)R(f'') + \frac{1}{nh}R(K).
 \end{aligned}$$

Il est clair que cette quantité ne dépend pas de  $x$ , d'où le  $h$  optimal calculer par minimisation de la AMISE est globale ;

$$\begin{aligned}
 h_{opt}^* &= Arg \min_h AMISE(f_n(x)) \\
 &= \left( \frac{R(K)}{\mu_2^2(K)R(f'')} \right)^{1/5} n^{-1/5} \simeq Cn^{-1/5}.
 \end{aligned} \tag{2.8}$$

**Exercice 2.2.1** Calculer  $R(K)$ ,  $\mu_2^2(K)$  et  $R(f'')$  dans les cas suivants :

- 1)  $K$  noyau gaussien et  $f$  densité normale,  $N(\mu, \sigma^2)$ .
- 2)  $K$  noyau d'Epanechnikov et  $f$  densité exponentielle,  $Exp(\lambda)$ .

En déduire le  $h_{opt}$  et  $h_{opt}^*$  en fonction de  $n$  dans les deux cas.

### 2.2.3 Choix du noyau

Lorsqu'on définit un estimateur à noyau, on a non-seulement le choix de la fenêtre  $h$ , mais aussi celui du noyau  $K$ . Pour désigner un noyau optimal dans l'estimation NP de la densité, il suffit d'insérer la valeur du  $h_{opt}^*$  dans la AMISE :

$$\begin{aligned}
 AMISE(f_n(x)) &= \frac{(h_{opt}^*)^4}{4}\mu_2^2(K)R(f'') + \frac{1}{nh_{opt}^*}R(K) \\
 &= \frac{1}{4} \left( \frac{R(K)}{\mu_2^2(K)R(f'')} \right)^{4/5} n^{-4/5}\mu_2^2(K)R(f'') + \left( \frac{R(K)}{\mu_2^2(K)R(f'')} \right)^{-1/5} n^{-4/5}R(K) \\
 &= \frac{5}{4} \{ \mu_2^2(K)R(f'') R^4(K) \}^{1/5} n^{-4/5}.
 \end{aligned} \tag{2.9}$$

Cette dernière expression ne dépend plus de  $h$  ou de  $x$ , elle dépend seulement du noyau  $K$ . La minimisation de (2.9) par rapport à  $K$  donne comme solution le noyau dite d'Epanechnikov (1969) :

$$K_{opt}(t) = \frac{3}{4} (1 - t^2) I_{(|t| < 1)}.$$

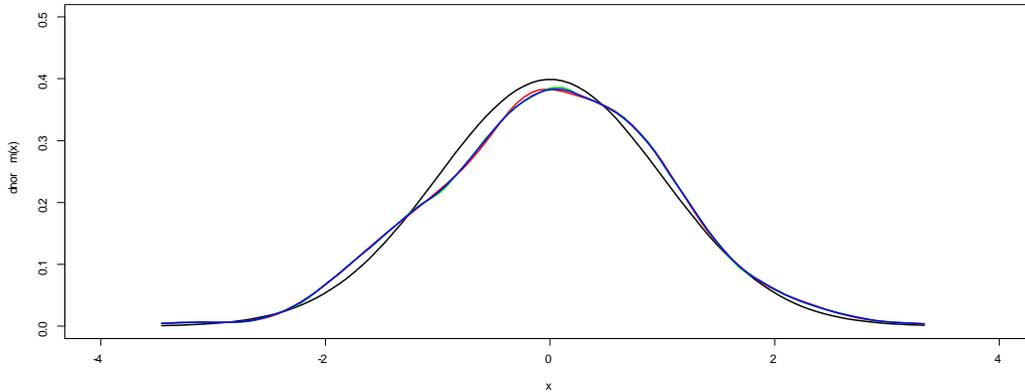


FIG. 2.4 – Effet du noyau : estimation d'une densité normale par différents noyaux.

**Définition 2.2.2** L'efficacité relative d'un noyau  $K$  par rapport à  $K_{opt}$  est donnée par :

$$eff(K) = \frac{AMISE(K_{opt})}{AMISE(K)} = \left( \frac{\mu_2^2(K_{opt})R^4(K_{opt})}{\mu_2^2(K)R^4(K)} \right)^{1/5} \leq 1. \quad (2.10)$$

**Exercice 2.2.2** Vérifier les résultats du tableau suivant :

Noyau	$eff(K)$
Epanechnikov	1.000
Gaussien	0.951
Uniform	0.930
Triangulaire	0.986
Quartique	0.994

**Remarque 2.2.4** Il est clair que le choix du noyau n'influe pas trop dans le cas des noyaux symétriques.

## 2.2.4 Choix du paramètre de lissage

Contrairement au noyau  $K$ , la fenêtre  $h$  à un rôle important dans l'estimation à noyau. Elle détermine la qualité d'estimation et contrôle le lissage (sur lissage, sous lissage) d'une façon très sensible. En effet, le choix de  $h$  tourne au tour de la condition :

$$h \rightarrow 0 \text{ et } nh \rightarrow \infty.$$

Des faibles valeurs de  $h$  impliquent un sous-lissage ( $Biais \rightarrow 0$ ) mais la variance augmente, et lorsque  $h$  grandit, la variance devient petite et c'est le biais qui augmente, donc sur-lissage de l'estimateur.

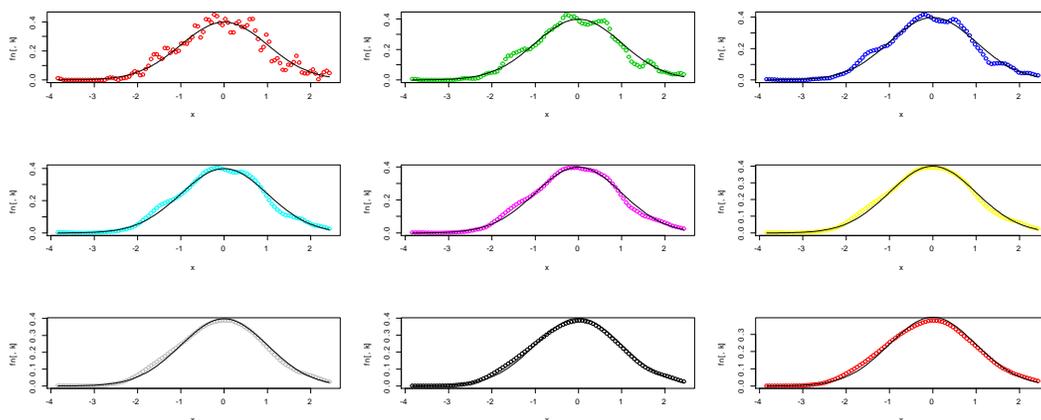


FIG. 2.5 – Effet de la fenêtre  $h$  : estimation d'une densité normale par différentes valeurs de  $h$ .

De plus, puisque  $h_{opt}$  et  $h_{opt}^*$  dépendent des quantités inconnues  $(f, f'')$ , donc pratiquement ne sont plus calculables. Plusieurs méthodes ont été développées pour résoudre ce problème : Validation croisée, Plug-in et la règle de référence. C'est cette dernière qu'on va donner en détail dans la suite.

### Règle de référence à une loi normale :

Silverman (1986) a proposé de se référer à une loi normale pour le calcul de  $h_{opt}^*$ . Soit  $X_1, X_2, \dots, X_n$  une suite de variables aléatoires de densité de probabilité  $f$ , supposons que  $f$  appartient à une

famille de distributions normales,  $N(\mu; \sigma^2)$ . Alors  $f(x) = \frac{1}{\sigma} \varphi\left(\frac{x-\mu}{\sigma}\right)$ , avec

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}, \quad f''(x) = \frac{1}{\sigma^3} \varphi''\left(\frac{x-\mu}{\sigma}\right) \quad \text{et} \quad \varphi''(x) = \frac{1}{\sqrt{2\pi}} (x^2 - 1) e^{-\frac{x^2}{2}}$$

La quantité inconnue  $R(f'')$  s'écrit alors

$$\begin{aligned} R(f'') &= \int f''(x)^2 dx \\ &= \frac{1}{\sigma^6} \int \left\{ \varphi''\left(\frac{x-\mu}{\sigma}\right) \right\}^2 dx = \frac{1}{\sigma^5} \int \left\{ \varphi''(v) \right\}^2 dv \end{aligned}$$

Nous avons :

$$\begin{aligned} \varphi(v) &= \frac{1}{\sqrt{2\pi}} e^{-\frac{v^2}{2}} \Rightarrow \varphi'(v) = -\frac{v}{\sqrt{2\pi}} e^{-\frac{v^2}{2}} \\ \Rightarrow \varphi''(v) &= \frac{1}{\sqrt{2\pi}} (v^2 - 1) e^{-\frac{v^2}{2}}. \end{aligned}$$

Alors,

$$\begin{aligned} R(f'') &= \frac{1}{\sigma^5} \int \left\{ \frac{1}{\sqrt{2\pi}} (v^2 - 1) e^{-\frac{v^2}{2}} \right\}^2 dv \\ &= \frac{1}{\sigma^5} \frac{1}{\sqrt{2\pi}} \left\{ \int v^4 e^{-v^2} dv - 2 \int v^2 e^{-v^2} dv + \int e^{-v^2} dv \right\} \\ &= \frac{1}{\sigma^5} \frac{1}{\sqrt{2\pi}} \left\{ -\frac{1}{2} \int v^2 e^{-v^2} dv + \int e^{-v^2} dv \right\} \\ &= \frac{1}{\sigma^5} \frac{1}{\sqrt{2\pi}} \left\{ -\frac{1}{2} \int \frac{\mu^2}{2} e^{-\frac{\mu^2}{2}} \frac{1}{\sqrt{2}} d\mu + \int \frac{1}{\sqrt{2}} e^{-\frac{\mu^2}{2}} d\mu \right\} \quad \text{avec } \mu = \sqrt{2}v \\ &= \frac{1}{\sigma^5} \frac{1}{2\pi} \left\{ -\frac{1}{4} \sqrt{\pi} + \sqrt{\pi} \right\} = \frac{1}{\sigma^5} \frac{3}{8\sqrt{\pi}}. \end{aligned}$$

Donc, l'expression du paramètre de lissage optimal devient

$$h_{opt}^* = \left( \frac{8\sqrt{\pi}R(k)}{3(\mu_2(k))^2} \right)^{\frac{1}{5}} \hat{\sigma} n^{-\frac{1}{5}}$$

où  $\hat{\sigma}^2$  est la variance empirique (estimateur sans biais) de la variance  $\sigma^2$  de  $X$  :

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2, \quad \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

**Remarque 2.2.5 (cas particulier)** *i) Si  $K$  est un noyau d'Epanechnikov, alors*

$$h_{opt}^* = 2.34 \hat{\sigma} n^{-1/5}.$$

*ii) Si  $K$  est un noyau gaussien, alors*

$$h_{opt}^* = 1.06 \hat{\sigma} n^{-1/5}.$$

*iii) Si  $K$  est un noyau Quartique, i.e.,  $K(t) = \frac{15}{16} (1-t^2)^2 I_{(|t|<1)}$ , donc*

$$h_{opt}^* = 2.78 \hat{\sigma} n^{-1/5}.$$

**Exercice 2.2.3** *Montrer les trois résultats précédents.*

## 2.2.5 Comportement Asymptotique

Dans cette section, nous allons étudier le comportement asymptotique de  $f_n$  l'estimateur à noyau de la densité  $f$  (Normalité asymptotique, convergence presque sure est uniforme).

**Théorème 2.2.1** *Supposons que  $h := h_n$ , est tel que*

$$h \rightarrow 0, \quad nh^3 \rightarrow 0 \quad \text{et} \quad nh \rightarrow \infty, \quad \text{quand } n \rightarrow \infty.$$

$f$  est dérivable sur un voisinage de  $x \in \mathbb{R}$  (où  $f(x) > 0$ ) de dérivées bornées. Alors,

$$\sqrt{2nh} \left( \frac{f_n(x) - f(x)}{\sqrt{f_n(x)}} \right) \rightarrow N(0, 1) \text{ en loi, quand } n \rightarrow \infty.$$

**Corollaire 2.2.1** *Sous les mêmes conditions du théorème, l'intervalle de confiance au niveau  $(1 - \alpha)\%$  de  $f$  est*

$$\left[ f_n(x) - z_{1-\alpha/2} \sqrt{\frac{f_n(x)}{2nh}}; f_n(x) - z_{1-\alpha/2} \sqrt{\frac{f_n(x)}{2nh}} \right],$$

où  $z_{1-\alpha/2}$  est le quantile d'ordre  $(1 - \frac{\alpha}{2})$  de la loi  $N(0, 1)$ .

**Théorème 2.2.2** *Si  $f$  est continue sur le voisinage de  $x$  et si*

$$h \rightarrow 0, \quad \frac{\log n}{nh} \rightarrow 0 \text{ et } nh \rightarrow \infty, \text{ quand } n \rightarrow \infty.$$

*Supposons que le noyau  $K$  est d'ordre 2, borné, intégrable et à support compact. Alors,*

$$f_n(x) \rightarrow f(x), \quad Ps, \text{ quand } n \rightarrow \infty.$$

*Si de plus,  $f$  est deux fois continuellement dérivable. Alors,*

$$f_n(x) - f(x) = O(h^2) + O\left(\frac{\log n}{nh}\right).$$

**Théorème 2.2.3** *Sous les mêmes conditions du théorème précédente. Supposons de plus que  $\Omega$  est un compact de  $\mathbb{R}$  sur le quelle  $f$  est deux fois continuellement dérivable et que le noyau  $K$  est Lipschitzien sur  $\Omega$  :*

$$\exists c < \infty, \forall x, y \in \Omega : |K(x) - K(y)| \leq c|x - y|.$$

*Alors,*

$$\sup_{x \in \Omega} |f_n(x) - f(x)| = O(h^2) + O\left(\frac{\log n}{nh}\right), \quad Ps.$$