

République Algérienne Démocratique et Populaire  
Ministère de l'Enseignement Supérieur et de la Recherche Scientifique

**UNIVERSITÉ MOHAMED KHIDER, BISKRA**

FACULTÉ des SCIENCES EXACTES et des SCIENCES de la NATURE et de la VIE

**DÉPARTEMENT DE MATHÉMATIQUES**



Mémoire présenté en vue de l'obtention du Diplôme :

**MASTER en Mathématiques**

Option : **Statistique**

Par

**MANCER Aicha**

Titre :

# Estimation non paramétrique de la fonction de régression

Membres du Comité d'Examen :

Mme. TOUBA Sounia	UMKB	Président
Mlle. KHEIREDDINE Souraya	UMKB	Encadreur
Mme. OUANOUGHY Yasmina	UMKB	Examinateur

Juin 2018

## REMERCIEMENTS

*Je tiens à dire que j'ai réussie à réaliser ce modeste travail grâce à **Dieu** qui m'a donné le pouvoir, la santé, la volonté et le courage d'être arrivé jusqu'au là.*

*A mon très cher père et ma très chère mère.*

*Je remercie en particulier mon encadreur Dr. **Kheireddine Souraya** qui m'a honorés par son encadrement, sa direction, son orientation, sa modestie, sa patience, ses conseils et toutes, ses remarques constructives pour le bon déroulement de mon travail.*

*J'exprime aussi mes remerciements à tous le personnel de département de Mathématiques, à toute personne qui a posée son empreinte de loin ou de près pour faire aboutir ce modeste travail.*

*Merci du fond du cœur*

# Table des matières

Remerciements	i
Table des matières	ii
Liste des figures	iv
Introduction	1
<b>1 Estimation non paramétrique de la densité</b>	<b>3</b>
1.1 Définition d'estimation paramétrique et non paramétrique . . . . .	3
1.2 Estimation non paramétrique de la densité par la méthode du noyau . . . . .	5
1.2.1 Propriétés de l'estimateur . . . . .	7
1.2.2 La consistance de l'estimateur . . . . .	8
<b>2 Estimation non paramétrique de la fonction de régression</b>	<b>12</b>
2.1 L'estimateur non paramétrique de régression . . . . .	12
2.2 Les propriétés de l'estimateur . . . . .	14
2.2.1 La consistance . . . . .	14
2.2.2 Convergence presque complète . . . . .	17
2.3 La normalité asymptotique de l'estimateur . . . . .	18
2.4 Choix du paramètre de lissage . . . . .	18
2.4.1 Critère d'erreur quadratique moyenne de $r_n(x)$ . . . . .	19

<b>3 Simulation</b>	<b>22</b>
3.1 Régression linéaire . . . . .	24
3.1.1 Paramètre de lissage $h$ fixé, $n$ varié . . . . .	24
3.1.2 Choix graphique du paramètre de lissage . . . . .	29
3.2 Régression non linéaire . . . . .	32
3.2.1 Paramètre de lissage $h$ fixé, $n$ varié . . . . .	33
3.2.2 Choix graphique du paramètre de lissage . . . . .	37
<b>Conclusion</b>	<b>41</b>
<b>Bibliographie</b>	<b>42</b>
<b>Annexe A : Rappels</b>	<b>44</b>
<b>Annexe B : Abréviations et Notations</b>	<b>46</b>

# Table des figures

3.1	Régression linéaire : $h$ fixé, $n$ varié et $K$ noyau normale. . . . .	27
3.2	Régression linéaire : $h$ fixé, $n$ varié et $K$ noyau d'Epanechnikov . . . . .	28
3.3	Régression linéaire avec $h$ varié, $n$ fixé et $K$ noyau gaussien. . . . .	31
3.4	Régression linéaire avec $h$ varié, $n$ fixé et $K$ d'Epanechnikov . . . . .	32
3.5	Régression non linéaire : $h$ fixé, $n$ varié et $K$ noyau normale . . . . .	36
3.6	Régression non linéaire : $h$ fixé, $n$ varié et $K$ noyau d'Epanechnikov . . . . .	37
3.7	Régression non linéaire avec $h$ varié, $n$ fixé et $K$ gaussien. . . . .	39
3.8	Régression non linéaire avec $h$ varié, $n$ fixé et $K$ d'Epanechnikov. . . . .	40

# Introduction

La théorie de l'estimation est une des préoccupations majeures des statisticiens. Ainsi l'estimation non paramétrique réelle a reçu un intérêt croissant tant sur le plan théorique que pratique. Elle consiste généralement à estimer à partir des observations une fonction inconnue, élément d'une certaine classe fonctionnelle, telle que la fonction de densité ou la fonction de régression à titre d'exemples.

Les travaux de Rosenblatt (1956)[14] et Parzen (1962)[13] puis de Nadaraya- Watson (1964)[[11], [17]] portant respectivement sur les estimateurs non paramétriques des fonctions de la densité et de la régression, la méthode du noyau a été largement utilisée dans de nombreux travaux.

L'estimation de la fonction de régression est un problème important dans l'analyse des données avec un large gamme d'applications en filtrage et la prévision dans les communications et le contrôle des systèmes, la reconnaissance de formes et de classification...

L'objet de cette mémoire est l'étude d'estimateur non paramétriques de fonction de régression par la méthode du noyau.

Ce travail est subdivisé en trois chapitres :

Dans le premier chapitre nous introduisons la définition de l'estimation fonctionnelle et l'estimateur à noyau de la densité (estimateur de Parzen-Rosenblatt), et ses propriétés asymptotiques.

Nous présentons, au chapitre 2, l'estimateur non paramétrique de la fonction de régression (estimateur de Nadaraya-Watson). Nous traitons ici, les propriétés asymptotiques de cet

estimateur (la convergence en moyenne quadratique et la convergence presque complète) et sa normalité asymptotique.

Finalement, Nous terminons notre mémoire par un troisième chapitre où nous utilisons le logiciel R pour donnons des exemples sur l'estimateur non paramétrique de régression (linéaire et non linéaire) par la méthode du noyau qui expriment l'importance de paramètre de lissage  $h$ , du noyau  $K$ .

# Chapitre 1

## Estimation non paramétrique de la densité

### 1.1 Définition d'estimation paramétrique et non paramétrique

Premièrement, nous appelons modèle statistique, le triplet  $(\mathbb{E}, \mathcal{A}, \mathbb{P})$  où  $\mathbb{E}$  est l'espace des observations (par exemples des réels),  $\mathcal{A}$  une tribu sur  $\mathbb{E}$  et  $\mathbb{P}$  une famille de probabilité sur  $(\mathbb{E}, \mathcal{A})$ .

Soit  $X : \Omega \longrightarrow \mathbb{E}$  une application mesurable. On peut toujours écrire  $\mathbb{P}$  par  $(\mathbb{P}_\theta, \theta \in \Theta)$ .

Soit  $h$  une application de  $\mathbb{P}$  dans  $\Theta'$ . Estimer  $h(P)$  c'est essayer de l'évaluer au vu de l'observation d'un échantillon de la variable aléatoire  $X$  qui est à valeurs dans  $\mathbb{E}$ . Donc, le paramètre à estimer est l'application

$$h : P \longrightarrow \Theta' \quad \text{ou} \quad \Theta \longrightarrow \Theta'$$
$$\theta \longmapsto h(P_\theta)$$

Un estimateur de  $h$  est une fonction  $h_n : x \longmapsto h_n(X_1, \dots, X_n)$  mesurable par rapport à l'observation  $(X_1, \dots, X_n)$ .

**Définition 1.1.1 Estimation paramétrique :** Si l'on sait à priori que  $h$  appartient à une famille paramétrée  $\{h(x, \theta), \theta \in \Theta\}$  où  $\Theta \subset \mathbb{R}^s$  et  $h(\cdot, \cdot)$  est une fonction connue, on parle alors d'estimation paramétrique, car estimer  $h$  revient à estimer le paramètre fini-dimensionnel  $\theta$ .

**Définition 1.1.2 Estimation non paramétrique :** Par contre, si l'on sait seulement que  $h$  appartient à  $\mathbb{P}$  ensemble des lois de probabilités qui est un espace de dimension infinie, alors on dit que l'on fait de l'estimation non paramétrique ou de l'estimation fonctionnelle.

Dans ce qui suit, on suppose que l'on a observé un échantillon  $X_1, X_2, \dots, X_n$  à valeurs dans  $\mathbb{R}^s$  muni de sa tribu borélienne  $\mathcal{B}$ . De plus, on suppose que les  $\{X_i, i = 1, \dots, n\}$  sont indépendantes et identiquement distribuées (*i.i.d.*)  $\mu \in P_0$  une famille de loi sur  $(\mathbb{R}^s; \mathcal{B})$ ;

i) **La densité de probabilité :** Si  $P_0$  est une famille de loi dominée par une loi  $\lambda$ , donc elle admet (théorème de Radon-Nykadim) une densité  $f = \frac{d\mu}{d\lambda}$  c'est un paramètre dans  $L^1$ . Si  $\frac{d\mu}{d\lambda}$  admet une version bornée (respectivement continue et bornée) alors on peut la considérer comme un paramètre dans  $L^2$  (respectivement dans  $C_b(\mathbb{R}^s)$ ).

Enfin, si  $f_\mu$  est différentiable, on définit de nouveaux paramètres fonctionnels : les dérivées partielles de  $f_\mu$ .

ii) **La fonction de répartition :** C'est la fonction définie par

$$F_\mu(x_1, \dots, x_s) = \mu\left(\prod_{i=1}^s ]-\infty, x_i]\right), \quad (x_1, \dots, x_s) \in \mathbb{R}^s.$$

iii) **La fonction des quantiles :** Pour  $s = 1$ , la fonction quantile d'ordre  $p$  définie par

$$F_\mu^{-1}(p) = Q(p) = \inf\{t \in \mathbb{R}; F_\mu(t) \geq p\}, \quad 0 < p < 1.$$

$F_\mu^{-1}$  est un paramètre à valeur dans l'espace de fonctions réelles définies sur  $]0; 1[$  monotones non décroissantes et continues à gauche.

v) **La fonction caractéristique** : Elle est définie par

$$\hat{\mu}(t) = E_{\mu} [\exp \{i \langle t, x \rangle\}] \quad \text{où } t, x \in \mathbb{R}^s,$$

$\hat{\mu}$  est un paramètre dans  $C_b(\mathbb{R}^s)$ .

iv) **Le paramètre de régression** : Supposons que l'on observe un échantillon  $\{(X_i, Y_i), i = 1, \dots, n\}$  d'un couple  $(X, Y)$  à valeurs dans  $\mathbb{R}^{s_1} \times \mathbb{R}^{s_2}$  est soit  $\mu_Y^x, x \in \mathbb{R}^{s_1}$  une famille de versions des lois conditionnelles de  $Y$  sachant  $X = x$  : Toute fonction de la forme  $r : x \mapsto r(\mu_Y^x)$  est un paramètre de régression. Les plus usuels sont :

- 1) L'espérance conditionnelle (qui est la fonction de régression),
- 2) La densité conditionnelle,
- 3) Le mode conditionnel,
- 4) La fonction de répartition conditionnelle,
- 5) Le quantile conditionnel.

## 1.2 Estimation non paramétrique de la densité par la méthode du noyau

Soient  $X_1, \dots, X_n$  des variables aléatoires réelles (v.a) *i.i.d* de densité de probabilité  $f_X$  par rapport à la mesure de Lebesgue sur  $\mathbb{R}$  et de fonction de répartition  $F(x) = \int_{-\infty}^x f(t) dt$ .

Considérons la fonction de répartition empirique

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n 1_{(X_i \leq x)},$$

où  $1(\cdot)$  désigne la fonction indicatrice. D'après la loi forte des grands nombres, presque sûrement,

$$F_n(x) \longrightarrow F(x), \quad \forall x \in \mathbb{R},$$

quand  $n \rightarrow \infty$ . Donc  $F_n$  est un estimateur convergent (consistant) de  $F$  pour tout  $x \in \mathbb{R}$ . Comment peut-on estimer  $f$ ? Une des premières solutions intuitives a été proposée par Rosenblatt (1956)[14]. Pour  $h > 0$  assez petit on a

$$f_X(x) \simeq \frac{F(x+h) - F(x-h)}{2h}.$$

En remplaçant ici  $F$  par son estimateur  $F_n$ , on obtient.

$$f_{n,X}^R(x) = \frac{F_n(x+h) - F_n(x-h)}{2h}.$$

La fonction  $f_{n,X}^R$  est un estimateur de  $f$  appelé estimateur de Rosenblatt. On peut aussi l'écrire sous la forme

$$f_{n,X}^R(x) = \frac{1}{2nh} \sum_{i=1}^n 1_{\{x-h < X_i \leq x+h\}} = \frac{1}{nh} \sum_{i=1}^n K_0\left(\frac{X_i - x}{h}\right),$$

où  $K_0(u) = \frac{1}{2}1_{\{-1 < u \leq 1\}}$ . Parzen(1962)[13] a suggéré une généralisation de cet estimateur :

$$f_{n,X}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{X_i - x}{h}\right), \tag{1.1}$$

où  $K : \mathbb{R} \rightarrow \mathbb{R}$  est une fonction intégrable, positive et telle que  $\int K(u) du = 1$ . C'est l'estimateur à noyau de la densité ou estimateur de Parzen-Rosenblatt. La fonction  $K$  est dite noyau. Et  $h > 0$  un paramètre de lissage "fenêtre" (en anglais "bandwidth") de l'estimateur.

Dans le cadre asymptotique où  $n \rightarrow \infty$  on supposera que la fenêtre  $h$  dépend de  $n$  et on la notera  $h_n$ , la suite  $(h_n)_{n \geq 1}$  tendant vers 0 lorsque  $n \rightarrow \infty$ .

Voici quelques exemples de noyaux classiques :

- Noyau rectangulaire :

$$K_1(x) = \begin{cases} \frac{1}{2}, & \text{si } |x| \leq 1; \\ 0, & \text{si } |x| > 1. \end{cases}$$

- Noyau triangulaire :

$$K_2(x) = \begin{cases} 1 - |x|, & \text{si } |x| \leq 1; \\ 0, & \text{si } |x| > 1. \end{cases}$$

- Noyau d'Epanechnikov ou parabolique :

$$K_3(x) = \begin{cases} \frac{3}{4}(1 - x^2), & \text{si } x \in [-1, 1]; \\ 0, & \text{sinon.} \end{cases}$$

- Noyau quadratique :

$$K_4(x) = \begin{cases} \frac{15}{16}(1 - x^2)^2, & \text{si } x \in [-1, 1]; \\ 0, & \text{sinon.} \end{cases}$$

- Noyau cubique :

$$K_5(x) = \begin{cases} \frac{35}{32}(1 - x^2)^3, & \text{si } x \in [-1, 1]; \\ 0, & \text{sinon.} \end{cases}$$

- Noyau gaussien :

$$K_6(x) = \frac{1}{\sqrt{2\Pi}} \exp\left(\frac{-1}{2}x^2\right), \quad x \in \mathbb{R}.$$

On remarquera que si le noyau  $K$  est positif et si  $X_1, \dots, X_n$  sont fixés, la fonction  $x \rightarrow f_{n,X}(x)$  est une densité de probabilité.

### 1.2.1 Propriétés de l'estimateur

Le pilier des premiers résultats de la convergence de cet estimateur est le théorème de Bochner (1955)[3], rappelé ci dessous :

**Théorème 1.2.1** (*Bochner*)

Soit  $K : (\mathbb{R}^m, \beta^m) \rightarrow (\mathbb{R}, \beta)$  une fonction mesurable, où  $\beta^p$  est la tribu borélienne de  $\mathbb{R}^p$ , vérifiant :

$$\begin{aligned} \exists M \text{ (constante) telle que, } \quad \forall z \in \mathbb{R}^m, |K(z)| \leq M, \\ \int_{\mathbb{R}^m} |K(z)| dz < \infty, \end{aligned}$$

et

$$\|z\|^m |K(z)| \rightarrow 0 \text{ quand } \|z\| \rightarrow \infty.$$

par ailleurs, soit

$g : (\mathbb{R}^m, \beta^m) \rightarrow (\mathbb{R}, \beta)$  une fonction tq

$$\int_{\mathbb{R}^m} |g(z)| dz < \infty,$$

Si  $g$  est continue, et si  $0 < h_n \rightarrow 0$ , quand  $n \rightarrow \infty$  alors :

$$\lim_{n \rightarrow \infty} \frac{1}{h_n^m} \int_{\mathbb{R}^m} K\left(\frac{z}{h_n}\right) g(x - z) dz = g(x) \int_{\mathbb{R}^m} K(z) dz.$$

Si  $g$  est uniformément continue alors la convergence ci dessus est uniforme.

## 1.2.2 La consistance de l'estimateur

L'estimateur à noyau de la densité dépend de deux paramètres la fenêtre  $h$  et le noyau  $K$ . Le noyau  $K$  établit l'aspect du voisinage de  $x$  et  $h$  contrôle la taille de ce voisinage, donc  $h$  est le paramètre prédominant pour avoir de bonnes propriétés asymptotiques, néanmoins le noyau  $K$  ne doit pas être négligé, comme le montre le travail de Parzen (1962)[13] cité ci dessous sur la consistance de cet estimateur. Cette dernière est obtenue, en se basant

sur l'étude asymptotique du biais, de la variance et de la décomposition suivante :

$$E [f_{n,X}(x) - f_X(x)]^2 = Var [f_{n,X}(x)] + [Biais \{f_{n,X}(x)\}]^2.$$

Dans la suite, nous supposons que  $K$  est un noyau vérifiant les conditions suivantes.

**(H.1)**  $K$  est bornée, c'est à dire  $\sup_{x \in \mathbb{R}} |K(x)| < \infty$ ,

**(H.2)**  $\lim_{|x| \rightarrow \infty} |x| K(x) = 0$ , quand  $|x| \rightarrow \infty$ ,

**(H.3)**  $K \in L_1(\mathbb{R})$ , c'est à dire  $\int_{\mathbb{R}} |K(x)| dx < \infty$ ,

**(H.4)**  $\int_{\mathbb{R}} |K(x)| dx = 1$ .

**1) Etude du biais :**

**Proposition 1.2.1** *Sous les hypothèses [(H.1), (H.2), (H.3) et (H.4)] et si  $f_X$  est continue alors*

$$\forall x \in \mathbb{R} \lim_{n \rightarrow \infty} E [f_{n,X}(x)] = f_X(x). \quad (1.2)$$

**Preuve.** En effet :

$$\begin{aligned} E [f_{n,X}(x)] &= E \left[ \frac{1}{nh_n} \sum_{i=1}^n K \left( \frac{x - X_i}{h_n} \right) \right] \\ &= \frac{1}{h_n} E \left[ K \left( \frac{x - X_i}{h_n} \right) \right] \\ &= \frac{1}{h_n} \int_{\mathbb{R}} K \left( \frac{x - t}{h_n} \right) f_X(t) dt. \end{aligned}$$

En posant  $x - t = z$ , on arrive à :

$$E [f_{n,X}(x)] = \frac{1}{h_n} \int_{\mathbb{R}} K \left( \frac{z}{h_n} \right) f_X(x - z) dz.$$

Comme  $K$  et  $f$  vérifient les conditions du théorème de Bochner, et  $\lim_{n \rightarrow \infty} h_n = 0$ ,  $n \rightarrow \infty$ ,

on a alors

$$\lim_{n \rightarrow \infty} \frac{1}{h_n} \int_{\mathbb{R}} K \left( \frac{x - t}{h_n} \right) f_X(t) dt = f_X(x) \int_{\mathbb{R}} K(z) dz.$$

d'où

$$\lim_{n \rightarrow \infty} E[f_{n,X}(x)] = f_X(x).$$

■

Nous constatons que le biais de l'estimateur converge vers zéro quand la fenêtre tend vers zéro, de plus vu son expression, on constate qu'il ne dépend pas du nombre des variables, il dépend surtout du noyau  $K$ .

## 2) Etude de la variance de $f_{n,X}(x)$

**Proposition 1.2.2** *Sous les conditions [(H.1), (H.2), (H.3) et (H.4)] et si  $f_X$  est continue en tout point  $x$  de  $\mathbb{R}$ , alors*

$$\lim_{n \rightarrow \infty} \text{Var}[f_{n,X}(x)] = 0$$

**Preuve.** En effet :

$$\begin{aligned} \text{Var}[f_{n,X}(x)] &= E[f_{n,X}(x)]^2 - [Ef_{n,X}(x)]^2 \\ &\leq E[f_{n,X}(x)]^2 \\ &\leq \frac{1}{n} E \left[ \frac{1}{h_n} K \left( \frac{x - X_i}{h_n} \right) \right]^2 \\ &\leq \frac{1}{nh_n^2} \int K^2 \left( \frac{x-t}{h_n} \right) f_X(t) dt \\ &\leq \frac{1}{nh_n^2} \int K^2 \left( \frac{z}{h_n} \right) f_X(x-z) dz. \end{aligned}$$

Remarquons que (H.1) et (H.3) impliquent que le noyau est de carré intégrable et les hypothèses sur  $h_n$ ,  $K$  et  $f_X$  assurent que :

$$\frac{1}{nh_n^2} \int K^2 \left( \frac{z}{h_n} \right) f_X(x-z) dz \sim \frac{1}{nh_n} f_X(x) \int K^2(z) dz,$$

d'où

$$\lim_{n \rightarrow \infty} \text{Var} [f_{n,X}(x)] = 0, \quad \text{quand } nh_n \rightarrow \infty.$$

■

Ces deux propositions impliquent la convergence en moyenne quadratique et donc, à foriori, la consistance de l'estimateur.

# Chapitre 2

## Estimation non paramétrique de la fonction de régression

### 2.1 L'estimateur non paramétrique de régression

Nous disposons d'un échantillon composé de  $n$  couples indépendants de variables aléatoires  $(X_1, Y_1), \dots, (X_n, Y_n)$  et nous considérons le modèle de régression non paramétrique donné, pour  $i = 1, \dots, n$ , par

$$Y_i = r(X_i) + \varepsilon_i. \quad (2.1)$$

où  $\varepsilon_i$  est l'aléatoire centré et indépendante de  $X_i$  et  $r$  est une application mesurable réelle. La fonction de régression  $r(\cdot) = E[Y|X = \cdot]$ , apportant de l'information sur la relation de dépendance inconnue de  $Y$  et  $X$ , un problème important est l'estimation de  $r$  à partir de l'observation de  $n$  copies  $(X_i, Y_i)$ ,  $i = 1, \dots, n$  qui suivent la même loi que  $(X, Y)$ .

Supposons que  $(X, Y)$  a une densité  $f : (x, y) \rightarrow f(x, y)$  sur  $\mathbb{R}^2$  et que  $f_X : x \rightarrow f_X(x) = \int f(x, y) dy > 0$  (densité de  $X$ ). Alors,

$$\forall x \in \mathbb{R}, r(x) = E[Y|X = x] = \frac{\int y f(x, y) dy}{f_X(x)}.$$

Les densités  $f$  et  $f_X$  sont inconnues mais on peut les estimer via  $\forall (x, y) \in \mathbb{R}^2$ ,

$$f_n(x, y) = \frac{1}{nh_n^2} \sum_{i=1}^n K\left(\frac{X_i - x}{h_n}\right) K\left(\frac{Y_i - y}{h_n}\right),$$

$$f_{n,X}(x) = \frac{1}{nh_n} \sum_{i=1}^n K\left(\frac{X_i - x}{h_n}\right),$$

puis on considère l'estimateur de la régression

$$\forall x \in \mathbb{R}, r_n(x) = \frac{\int y f_n(x, y) dy}{f_{n,X}(x)} \mathbf{1}_{f_{n,X}(x) \neq 0}. \quad (2.2)$$

**Proposition 2.1.1**

Si  $K$  est un noyau d'ordre 1, l'estimateur défini par [2.2] vérifie

$$\begin{aligned} \forall x \in \mathbb{R}, r_n(x) &= \frac{\sum_{i=1}^n Y_i K\left(\frac{X_i - x}{h_n}\right)}{\sum_{i=1}^n K\left(\frac{X_i - x}{h_n}\right)} \mathbf{1}_{\left\{\sum_{i=1}^n K\left(\frac{X_i - x}{h_n}\right) \neq 0\right\}} \\ &= \frac{\sum_{i=1}^n Y_i K_{h_n}(X_i - x)}{\sum_{i=1}^n K_{h_n}(X_i - x)} \mathbf{1}_{\left\{\sum_{i=1}^n K_{h_n}(X_i - x) \neq 0\right\}}, \end{aligned}$$

où  $K_{h_n}(\cdot) = K_{h_n}(\cdot/h)$ .

donc l'estimateur à noyau de la régression est donné par :

$$r_n(x) = \begin{cases} \frac{\sum_{i=1}^n Y_i K_{h_n}(X_i - x)}{\sum_{i=1}^n K_{h_n}(X_i - x)} = \frac{\Phi_{n,X}(x)}{f_{n,X}(x)}, & \text{si } \sum_{i=1}^n K_{h_n}(X_i - x) \neq 0; \\ \frac{1}{n} \sum_{i=1}^n Y_i, & \text{si } \sum_{i=1}^n K_{h_n}(X_i - x) = 0. \end{cases}$$

C'est l'estimateur à noyau introduit par Nadaraya-Watson (Nadaraya, 1964[11] et Watson, 1964[17]).

La construction de cet estimateur dépend de deux paramètres, le paramètre de lissage  $h$  dont le choix est crucial pour obtenir de bonnes propriétés asymptotiques et la noyau  $K$  dont on ne peut pas négliger le rôle pour la réduction du biais.

## 2.2 Les propriétés de l'estimateur

D'une manière analogue aux propriétés asymptotiques de l'estimateur de Parzen Rosenblatt, nous étudions dans cette partie deux modes de convergence, la convergence en moyenne quadratique et la convergence presque complète.

En plus des condition **(H.1-H.4)** sur le noyau  $K$ , nous avons besoin des hypothèses suivantes.

**(H.5)**  $\int_{\mathbb{R}} uK(u) du = 0,$

**(H.6)**  $\int_{\mathbb{R}} u^2K(u) du < \infty.$

**(H.7)**  $K$  est borné, intégrable et à support compact.

### 2.2.1 La consistance

En vu de la décomposition suivante :

$$E[r_n(x) - r(x)]^2 = Var[r_n(x)] + [Er_n(x) - r(x)]^2.$$

L'étude asymptotique du biais et de la variance de l'estimateur de Nadaraya-Watson détermine les conditions suffisantes à la consistance de cet estimateur.

#### 1) Etude asymptotique de la variance

**Proposition 2.2.1** *Sous les hypothèse de la proposition (1.3.2) et si  $EY^2 < \infty$ , alors en chaque point de continuité des fonctions  $r(x)$ ,  $f_X(x)$  et  $\sigma^2(x) = Var(Y|X=x)$*

on a

$$Var[r_n(x)] = \frac{1}{nh_n} \left\{ \frac{\sigma^2(x)}{f_X(x)} \int_{\mathbb{R}} K^2(u) du \right\} (o(1) + 1).$$

où  $f_X(x) > 0$ .

**Preuve.**

Soit la fonction  $\psi(x) = \int_{\mathbb{R}} y^2 f(x,y) dy$ , en se basant sur le lemme de Bochner on a

$$\begin{aligned}
 \text{Var} [\Phi_{n,X}(x)] &= \frac{1}{nh_n^2} \left\{ E \left[ Y^2 K^2 \left( \frac{x-X}{h_n} \right) \right] - \left[ E Y K \left( \frac{x-X}{h_n} \right) \right]^2 \right\} \\
 &= \frac{1}{nh_n} \left\{ \int_{\mathbb{R}} K^2(u) \psi(x - h_n u) du - h_n \left( \int_{\mathbb{R}} K(u) f(x - uh_n) r(x - h_n u) \right)^2 \right\} \\
 &= \frac{1}{nh_n} \psi(x) \int_{\mathbb{R}} K^2(u) du (1 + o(1)),
 \end{aligned}$$

$$E \{ [f_{n,X}(x) - E(f_{n,X}(x))] \{ \Phi_{n,X}(x) - E(\Phi_{n,X}(x)) \} \} = \frac{1}{nh_n} \Phi(x) \int_{\mathbb{R}} K^2(u) du (1 + o(1))$$

et

$$\text{Var} [f_{n,X}(x)] = \frac{1}{nh_n} f_X(x) \int_{\mathbb{R}} K^2(u) du (1 + o(1))$$

posons

$$B_n(x) = \begin{pmatrix} f_{n,X}(x) \\ \Phi_{n,X}(x) \end{pmatrix}$$

et

$$A(x) = \begin{pmatrix} -r(x) \\ \frac{1}{[f_X(x)]^2} \end{pmatrix}.$$

La matrice de variance covariance de  $B_n(x)$  est alors donnée par l'expression suivante

$$\Sigma := \frac{1}{nh_n} \begin{pmatrix} f_X(x) & \Phi(x) \\ \Phi(x) & \psi(x) \end{pmatrix} \int_{\mathbb{R}} K^2(u) du (1 + o(1)).$$

En remarquant, que

$$\begin{aligned}
 \text{Var} [r_n(x)] &= A \Sigma A^t \\
 &= \frac{1}{nh_n} \left( \frac{\psi(x)}{[f_X(x)]^2} - \frac{(\Phi(x))^2}{[f_X(x)]^3} \right) \int_{\mathbb{R}} K^2(u) du (1 + o(1)),
 \end{aligned}$$

où  $A^t$  désigne la transposée de  $A$ , on obtient alors

$$\text{Var} [r_n(x)] = \frac{1}{nh_n} \left\{ \frac{\sigma^2(x)}{f_X(x)} \int_{\mathbb{R}} K^2(u) du \right\} (o(1) + 1).$$

■

## 2) Etude asymptotique du biais

L'étude asymptotique du biais repose sur la proposition suivante.

**Proposition 2.2.2** *Sous les hypothèse de la proposition (2.2.1) et*

a) Si  $|Y| \leq C_1 < \infty$  P.S et si  $nh_n \rightarrow \infty$ , quand  $n \rightarrow \infty$ , alors :

$$Er_n(x) = \frac{E[\Phi_{n,X}(x)]}{E[f_{n,X}(x)]} + O\left(\frac{1}{nh_n}\right).$$

b) Si  $EY^2 < \infty$ ,  $nh_n^2 \rightarrow \infty$ , quand  $n \rightarrow \infty$ , alors :

$$Er_n(x) = \frac{E[\Phi_{n,X}(x)]}{E[f_{n,X}(x)]} + O\left(\frac{1}{\sqrt{nh_n}}\right).$$

Maintenant nous sommes en mesure d'énoncer le resultat suivant.

**Proposition 2.2.3** *Si les condition (H.4), (H.5) et (H.6) sont vérifiées et si  $f_X(\cdot)$  et  $r(\cdot)$  sont le classe  $C^2(\mathbb{R})$  et si  $|Y|$  est borné.*

Alors :

$$E[r_n(x)] - r(x) = \frac{h_n^2}{2} \left\{ \left\{ r''(x) + 2r'(x) \frac{f'_X(x)}{f_X(x)} \right\} \int_{\mathbb{R}} u^2 K(u) du \right\} (1 + o(1)). \quad (2.3)$$

### Remarque 2.2.1

1)- Les conditions (H.4), (H.5) et (H.6) peuvent être remplacées par le noyau  $K$  est d'ordre 2 au sens de Gasser et Müller.

2)-  $o(1)$  dans la relation [2.3] est égale à  $O(h) + O((nh)^{-1})$ .

**Preuve.**

$$\begin{aligned} E[r_n(x)] - r(x) &= \left[ EK\left(\frac{x-X}{h_n}\right) \right]^{-1} \left\{ \int_{\mathbb{R}} \frac{1}{h_n} K\left(\frac{x-t}{h_n}\right) \Phi(t) dt - r(x) \int_{\mathbb{R}} \frac{1}{h_n} K\left(\frac{x-t}{h_n}\right) f(t) dt \right\} \\ &= \left\{ (f(x))^{-1} \left\{ \frac{h_n^2}{2} \Phi''(x) - \frac{h_n^2}{2} r(x) f''(x) \right\} \int_{\mathbb{R}} u^2 K(u) du + \Phi(x) - r(x) f(x) \right\} (1 + o(1)) \end{aligned}$$

comme  $\Phi(x) = r(x) f(x)$ . L'équation précédente peut s'écrire :

$$E[r_n(x)] - r(x) = \left\{ \frac{h_n^2}{2} \left\{ r''(x) + 2r'(x) \frac{f'_X(x)}{f_X(x)} \right\} \int_{\mathbb{R}} u^2 K(u) du \right\} (1 + o(1)). \quad (2.4)$$

D'où

$$\lim_{n \rightarrow \infty} E[r_n(x)] = r(x).$$

■

### 2.2.2 Convergence presque complète

En se basant sur la preuve donnée dans Ferraty et Vieu (2003)[7], nous traitons dans ce paragraphe la convergence presque complète de l'estimateur à noyau de la fonction de régression. Nous gardons quelques conditions précédentes, aux quelles nous rajoutons les hypothèses suivantes.

- $f_X, r$  sont des fonctions continues au voisinage de  $x$ , un point fixé de  $\mathbb{R}$ . (2.5)

La densité  $f_X$  et la variable  $Y$  sont telles que

$$f_X > 0 \quad (2.6)$$

- Le paramètre de lissage  $h_n$  est tel que

$$\lim_{n \rightarrow \infty} h_n = 0 \text{ et } \lim_{n \rightarrow \infty} \frac{\log n}{nh_n} = 0, \quad (2.7)$$

et

$$|Y| < M < \infty, \quad (2.8)$$

où  $M$  est une constante réelle positive.

**Théorème 2.2.1** *Sous les hypothèses (2.5), (2.6), (2.7), (2.8), (H.4) et (H.7), on a :*

$$\lim_{n \rightarrow \infty} r_n(x) = r(x). \quad p.co$$

## 2.3 La normalité asymptotique de l'estimateur

La première démonstration de la normalité asymptotique de l'estimateur est due à Schuster (1972) [15]. On se réfère également aux théorème 1.3 et 1.4 p. 117-120 de Nadaraya (1989) [12] et au théorème 4.2.1 p. 99 de Härdle (1990) [8], qui proposent d'autres méthodes de démonstration. Le noyau  $K$  est supposé borné, à support compact et d'ordre 2. La fenêtre  $h_n$  est choisie égale à  $cn^{-1/5}$ .

**Théorème 2.3.1 Härdle (1990)**

*Supposons  $Y$  bornée ou admettant un moment d'ordre  $l > 2$ . Les fonctions  $f_X(\cdot)$  et  $r(x)$  sont supposées deux fois continûment dérivables sur  $\mathbb{R}$ . A chaque point de continuité de  $\sigma^2(x)$ , tel que  $f_X(x) > 0$ ,*

$$(nh)^{1/2} \{r_n(x) - r(x)\} \xrightarrow{\mathcal{L}} \mathcal{N}(B(x), v^2(x)), \quad (2.9)$$

avec

$$v^2(x) := \frac{\sigma^2(x)}{f_X(x)} \int_{\mathbb{R}} K^2(u) du, \quad (\text{la variance asymptotique}),$$

et

$$B(x) := \left\{ r''(x) + 2r'(x) \frac{f'_X(x)}{f_X(x)} \right\} \times \int_{\mathbb{R}} u^2 K(u) du, \quad (\text{le biais asymptotique}).$$

## 2.4 Choix du paramètre de lissage

Les vitesses de convergence dépendent de deux paramètres : la fonction de noyau  $K$  dont l'efficacité est peu influente et le paramètre de lissage  $h_n$ , dont le choix est crucial

aussi bien pour l'approche ponctuelle que pour la globale que nous exposons ci après.

### 2.4.1 Critère d'erreur quadratique moyenne de $r_n(x)$

L'erreur quadratique moyenne  $MSE$  (mean square error) est une mesure permettant d'évaluer la similarité de  $r_n$  par rapport à la fonction de régression inconnue  $r$ , au point  $x$  donné de  $\mathbb{R}$ .

Notre but est de minimiser

$$MSE(r_n(x)) = E[r_n(x) - r(x)]^2.$$

Le développement de cette expression faite précédemment, nous donne

$$MSE(r_n(x)) = Var[r_n(x)] + [biais(r_n(x))]^2.$$

Nous constatons d'une part que les expressions du *biais* de  $r_n(x)$  et de la variance de  $r_n(x)$  (voir les propositions (2.2.3), (2.2.1)) permettent de conclure qu'une grande valeur de  $h_n$  donne une augmentation du *biais* et une diminution de la variance (estimation fortement biaisée) et qu'un faible paramètre  $h_n$ , donne une diminution du *biais* et une augmentation de la variance (phénomène de sous lissage).

D'autre part, sous les hypothèses de ces mêmes propositions, nous obtenons

$$\begin{aligned} MSE(r_n(x)) &= \frac{h_n^4}{4} \left[ \left( r''(x) + 2r'(x) \frac{f'_X(x)}{f_X(x)} \right) (u^2 K(u)) + o(1) \right]^2 \\ &+ \frac{1}{nh_n} \left( \frac{\sigma^2(x)}{f_X(x)} [K^2(u)] \right) (1 + o(1)), \end{aligned} \quad (2.10)$$

où  $u^p K^q(u) = \int t^p K^q(t) dt$ . Pour trouver donc un compromis entre le *biais* et la variance nous minimisons par rapport à  $h_n$  l'expression de l'erreur quadratique moyenne asympto-

tique  $AMSE$  (asymptotique mean square error) donnée par

$$AMSE [r_n(x)] = \frac{h_n^4}{4} \left\{ r''(x) + 2r'(x) \frac{f'_X(x)}{f_X(x)} \right\}^2 [u^2 K(u)]^2 + \frac{1}{nh_n} \times \frac{\sigma^2(x)}{f_X(x)} [K^2(u)].$$

Comme  $AMSE$  est une fonction convexe. La fenêtre  $h_{opt(r_n(x))}^{MSE} = \arg \min_h [AMSE(r_n(x))]$  est solution de l'équation suivante

$$\frac{\partial}{\partial h_n} \left[ \frac{h_n^4}{4} \left\{ r''(x) + 2r'(x) \frac{f'_X(x)}{f_X(x)} \right\}^2 [u^2 K(u)]^2 + \frac{1}{nh_n} \times \frac{\sigma^2(x)}{f_X(x)} [K^2(u)] \right] = 0.$$

lorsque  $\left[ r''(x) + 2r'(x) \frac{f'_X(x)}{f_X(x)} \right]^2 [u^2 K(u)] \neq 0$ ,

d'où

$$h_{opt(r_n(x))}^{MSE} = n^{-1/5} \left\{ \frac{\frac{\sigma^2(x)}{f_X(x)} [K^2(x)]}{\left\{ \left[ r''(x) + 2r'(x) \frac{f'_X(x)}{f_X(x)} \right] [tK] \right\}^2} \right\}^{1/5}.$$

On s'intéresse maintenant à l'approche globale pour la sélection du paramètre  $h_n$ , pour cela on introduit le critère d'erreur quadratique intégrée moyenne ou  $MISE$  (mean integrated squared error) de  $r_n(x)$ .

$$MISE [r_n(x)] = E \left[ \int_{\mathbb{R}} (r_n(x) - r(x))^2 dx \right],$$

En appliquant le théorème de Fubini, on a

$$MISE [r_n(x)] = \left[ \int_{\mathbb{R}} E (r_n(x) - r(x))^2 dx \right],$$

Sous les mêmes hypothèses que les propositions (2.2.3) et (2.2.1), on a

$$AMISE [r_n(x)] = \frac{h_n^4}{4} \int \left\{ r''(x) + 2r'(x) \frac{f'_X(x)}{f_X(x)} \right\}^2 dx [u^2 K] + \frac{1}{nh_n} \int \frac{\sigma^2(x)}{f_X(x)} dx [K^2(u)].$$

La fenêtre  $h_{opt(r_n(x))}^{MISE}$  minimisant l'*AMISE* du critère global est :

$$h_{opt(r_n(x))}^{MISE} = n^{-1/5} \left\{ \frac{\int \frac{\sigma_I^2(x)}{f_X(x)} [K^2] dx}{\int \left\{ r''(x) + 2r'(x) \frac{f'_X(x)}{f_X(x)} \right\}^2 dx [tK]} \right\}^{1/5}.$$

Un travail similaire se fait pour le choix optimum du paramètre de lissage dans le cas de l'estimateur de Parzen-Rosemblatt, nous obtenons :

$$h_{opt(f_{n,X}(x))}^{MSE} = n^{-1/5} \left\{ \frac{f_X(x) [K^2]}{(f_X''(x))^2 [t^2 K]^2} \right\}^{1/5}, \quad (2.11)$$

$$h_{opt(f_{n,X}(x))}^{MISE} = n^{-1/5} \left\{ \frac{[K^2]}{[t^2 K]^2 \int_{\mathbb{R}} (f_X''(x))^2 dx} \right\}^{1/5}. \quad (2.12)$$

quand  $f_X''(x) \neq 0$ .

En insérant  $h_{opt(f_{n,X}(x))}^{MSE}$  dans  $MSE[f_{n,X}(x)]$  on peut montrer que le taux de convergence est d'ordre  $n^{-4/5}$ , il est plus faible que celui de l'histogramme dont l'ordre est égale à  $n^{-2/3}$ .

Nous notons que l'expression de  $h_n$  optimal, minimisant asymptotiquement les quatre critères d'erreurs à la forme  $Cn^{-1/5}$  alors

$$h_{opt} = Cn^{-1/5}$$

où la constante  $C$  est en fonction de la distribution et de termes aléatoires inconnues.

# Chapitre 3

## Simulation

Dans ce dernier chapitre, nous utilisons le logiciel **R**, pour calculer et représenter graphiquement la fonction de regression et son estimateur en vue de les comparer dans des situations simulées. Il s'agit de l'estimateur proposé par Nadaraya-Watson(1964) et présenté au chapitre 2. Nous donnons des exemples sur cet estimateur qui expriment l'importance de paramètre de lissage  $h$ , du noyau  $K$ .

Ensuite, nous présentons les résultats obtenus pour les différents jeux de données ainsi que pour les différents noyaux  $K$  (noyau Gaussien : à support non compact et noyau Epanichnekov : à support compact), différents valeurs de  $h$  strictement positif ( $h$  fixé ou  $h$  varié), régression linéaire et non linéaire.

Rappelons qu'on suppose que l'on a observé un échantillon  $\{(X_i, Y_i), i = 1, \dots, n\}$  et on veut expliquer la variable aléatoire  $Y_i$  par  $X_i$ . De plus, on suppose que le modèle est donné par l'expression :

$$Y_i = r(X_i) + \varepsilon_i.$$

où  $\varepsilon_i$  est l'aléatoire centré et indépendante de  $X_i$ . Aussi la fonction de regression

$$r(x) = E[Y|X = x] = \frac{\int y f(x, y) dy}{f_X(x)} \quad (3.1)$$

où  $f_X(x)$  est la densité de la variable  $X$ .

Nous avons vu que  $r(x)$  est estimé par la quantité :

$$r_n(x) = \frac{\frac{1}{nh_n} \sum_{i=1}^n Y_i K\left(\frac{X_i - x}{h_n}\right)}{\frac{1}{nh_n} \sum_{i=1}^n K\left(\frac{X_i - x}{h_n}\right)} = \frac{\Phi_{n,X}(x)}{f_{n,X}(x)} \quad (3.2)$$

Il dépend de la taille de l'échantillon  $n$ , et aussi du noyau  $K$  et de la fenêtre  $h_n$  qu'il faut choisir pour calculer  $r_n(x)$ . avec  $\Phi_{n,X}(x)$  est l'estimateur naturel de  $\Phi_X(x)$  :

$$\Phi_{n,X}(x) = \frac{1}{nh_n} \sum_{i=1}^n Y_i K\left(\frac{X_i - x}{h_n}\right)$$

et  $f_{n,X}$  l'estimateur à noyau étudié au chapitre 1 de la densité :

$$f_{n,X}(x) = \frac{1}{nh_n} \sum_{i=1}^n K\left(\frac{X_i - x}{h_n}\right)$$

Dans la suite de ce chapitre, nous supposons que notre modèle à la forme

$$y = r(x) + \varepsilon, \quad \text{où } \varepsilon \longrightarrow \mathcal{N}(0, \sigma^2), \quad (3.3)$$

et on va estimer les 2 fonctions de régression suivantes à l'aide d'un estimateur de Nadaray-Watson :

- Régression linéaire :  $r(x) = 3 + 0.9x + \varepsilon$ .
- Régression non linéaire :  $r(x) = \cos(x) + \varepsilon$ .

on supposons que :  $X$  suit la loi exponentielle de paramètre 0.5 :  $X \sim \mathcal{E}(0.5)$

Nous allons donc étudier les cas suivants dans chaque modèle :

- Paramètre de lissage ou fenêtre  $h$  fixe, noyau normale (noyau à support non compact) et  $n$  varié.
- Paramètre de lissage ou fenêtre  $h$  fixe, noyau d'Epanechnikov (noyau à support non compact) et  $n$  varié.

- $n$  fixe et fenêtre,  $h$  varié (noyau normale).
- $n$  fixe et fenêtre,  $h$  varié (noyau d'Epanechnikov).

## 3.1 Régression linéaire

On veut estimer le modèle linéaire

$$y = 3 + 0.9x + \varepsilon.$$

où  $\varepsilon$  un terme d'erreur de loi  $\mathcal{N}(0, 1)$ .

Dans les résultats graphique de cette section, on a :

- \* La droite noire exprime la fonction de régression  $r(x)$  [Eq.3.1].
- \* La droite en rouge exprime la fonction de régression empirique  $r_n(x)$  [Eq.3.2].

### 3.1.1 Paramètre de lissage $h$ fixé, $n$ varié

En choisissant le paramètre de lissage  $h_n = n^{-1/5}$  (fixé) et  $n$  varié ( $n = 40, 200, 400$ )

#### 3.1.1.1 $K$ à support non compact

Dans ce premier cas, on pose un noyau gaussien  $K(t) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}t^2\right)$ , et on va utiliser le code ci-dessous pour estimer ce modèle, et le resultat graphique obtenu représenté dans la figure [FIG-3.1] :

**Code R utilisé :**

```
rm(list=ls(all=TRUE))
```

```
n=40
```

```
X=rexp(n,0.5)
```

```
E=rnorm(n)
```

```
Y=3+.9*X+E
```

```
# Noyau Normale K(t)
K=function(t){(1/sqrt(2*pi))*exp(-0.5*t^2)}
h=n^-.2
# Initiation
s=100 # taille de l'intervalle [a,b]
a=min(X) #borne inf
b=max(X) # borne sup
x=seq(a,b,length=s) # Intervalle [a,b]
V=numeric(n)
fn=numeric(s)
for(j in 1 :s){
  for(i in 1 :n){ V[i]=K((x[j]-X[i])/h) }
  fn[j]=sum(V)/(n*h)}
# Fonction Hn(.)
W=numeric(n)
Hn=numeric(s)
for(j in 1 :s){
  for(i in 1 :n){ W[i]=K((x[j]-X[i])/h)*Y[i] }
  Hn[j]=sum(W)/(n*h)}
Rn =Hn/fn
# Graphes
op=par(mfrow=c(1,3))
plot(x,Rn,xlab="x", ylab="Rn(x)", main="n=40",type='l',col=2, lwd= 2)
abline(3,.9,lwd= 2)
#####Pour n =200
n=200
X=rexp(n,0.5)
```

```
E=rnorm(n)
Y=3+.9*X+E
h=n^-.2
V=numeric(n)
for(j in 1 :s){
for(i in 1 :n){ V[i]=K((x[j]-X[i])/h) }
fn[j]=sum(V)/(n*h)}
W=numeric(n)
for(j in 1 :s){
for(i in 1 :n){ W[i]=K((x[j]-X[i])/h)*Y[i] }
Hn[j]=sum(W)/(n*h)}
Rn =Hn/fn
plot(x,Rn,xlab="x", ylab="Rn(x)", main="n=200",type='l',col=2, lwd= 2)
abline(3,.9,lwd= 2)
#####Pour n =400
n=400
X=rexp(n,0.5)
E=rnorm(n)
Y=3+.9*X+E
h=n^-.2
V=numeric(n)
for(j in 1 :s){
for(i in 1 :n){ V[i]=K((x[j]-X[i])/h) }
fn[j]=sum(V)/(n*h)}
W=numeric(n)
for(j in 1 :s){
for(i in 1 :n){ W[i]=K((x[j]-X[i])/h)*Y[i] }
```

```
Hn[j]=sum(W)/(n*h)}
Rn =Hn/fn
plot(x,Rn,xlab="x", ylab="Rn(x)", main="n=400",type='l',col=2, lwd= 2)
abline(3,.9,lwd= 2)
par(op)
```

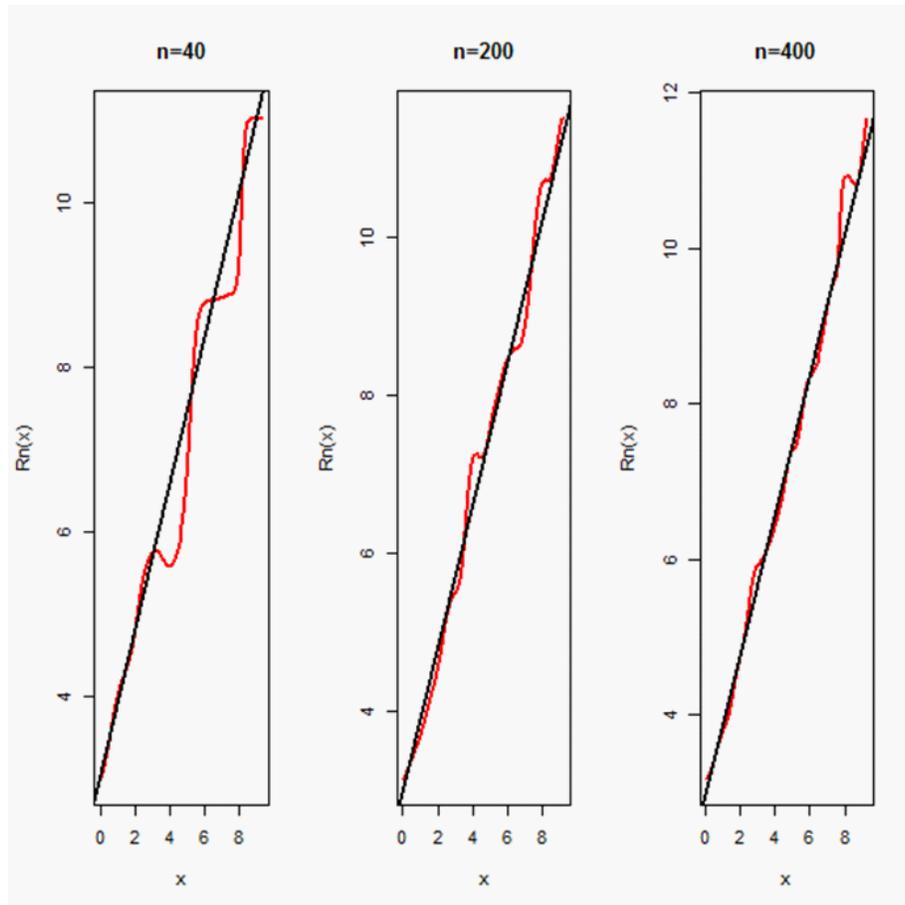


FIG. 3.1 – Régression linéaire :  $h$  fixé,  $n$  varié et  $K$  noyau normale.

L'axe des abscises représente les valeurs des  $x$  et l'axe des coordonnées les valeurs des  $r_n$  (et  $r$ ). Par la comparaison graphique, on remarque que le graphe rouge de  $r_n$  est proche beaucoup à la droite noire de  $r$  dans le troisième graphe, donc ce graphe exprime la convergence de l'estimateur  $r_n$  vers  $r$ .

### 3.1.1.2 $K$ à support compact

Dans ce second cas, on choisit le noyau d'Epanechnikov :  $K(t) = \frac{3}{4}(1 - t^2)$ . Ensuite, on modifie seulement cette partie dans le programme **R** précédent :

```
# Noyau Epanechnikov K(t)
K=function(t){ifelse(abs(t)<1,(3/4)*(1-t^2),0)}
```

On obtient la figure [FIG-3.2], et on arrive à la même conclusion de la convergence de l'estimateur (voire la [FIG-3.1], *i.e.*, convergence de l'estimateur pour  $n$  assez grand).

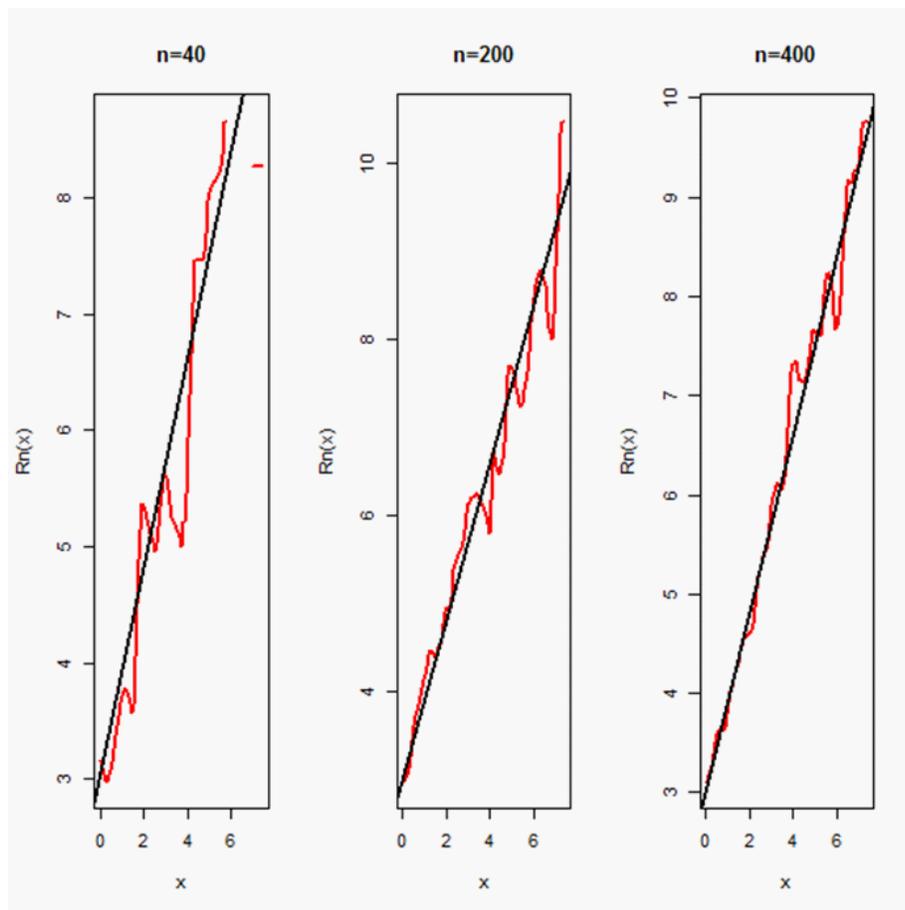


FIG. 3.2 – Régression linéaire :  $h$  fixé,  $n$  varié et  $K$  noyau d'Epanechnikov

### 3.1.2 Choix graphique du paramètre de lissage

Dans cette section, nous prenons le paramètre de lissage dans l'intervalle  $(0; 1)$  et avec des tests graphique en va diterminer le paramètre  $h$  optimal (au sens graphique).

On fixe la taille de l'échantillon  $n = 350$  et le noyau  $K$  est normale, l'estimation obtenue avec les valeurs de  $h$  varié de 0.1 à 0.9 sont données dans la figure [FIG-3.3]. Il est clair que la valeur du  $h$  optimale est de  $h = 0.7$  (ligne 3, colonne 1).

#### Code R

```
n=350
X=rexp(n,0.5)
E=rnorm(n)
Y=3+.9*X+E
# Noyau Normale K(t)
K=function(t){(1/sqrt(2*pi))*exp(-0.5*t^2)}
h=seq(.1,.9,length=9)
# Initiation
s=100 # taille de l'intervalle [a,b]
a=min(X) #borne inf
b=max(X) # borne sup
x=seq(a,b,length=s) # Intervalle [a,b]
V=array(dim=c(n,s,9))
fn=array(dim=c(s,9))
W=array(dim=c(n,s,9))
Hn=array(dim=c(s,9))
for(k in 1 :9){
  for(j in 1 :s){
    for(i in 1 :n){ V[i,j,k]=K((x[j]-X[i])/h[k]) }
    fn[j,k]=sum(V[,j,k])/(n*h[k])}}
```

```
# fonction Hn(x)
for(k in 1 :9){
  for(j in 1 :s){
    for(i in 1 :n){ W[i,j,k]=K((x[j]-X[i])/h[k])*Y[i] }
    Hn[j,k]=sum(W[,j,k])/(n*h[k])}
  Rn=array(dim=c(s,9))
  for(k in 1 :9){ Rn[,k]=Hn[,k]/fn[,k]}
# Graphes
x11() # nouvelle fenetre graphique
op=par(mfrow=c(3,3))
for(k in 1 :9){
  plot(x,Rn[,k],xlab="x", ylab="Rn(x)", main=" ",type='l',col=2, lwd= 2)
  abline(3,.9,lwd= 2)
}
par(op)
```

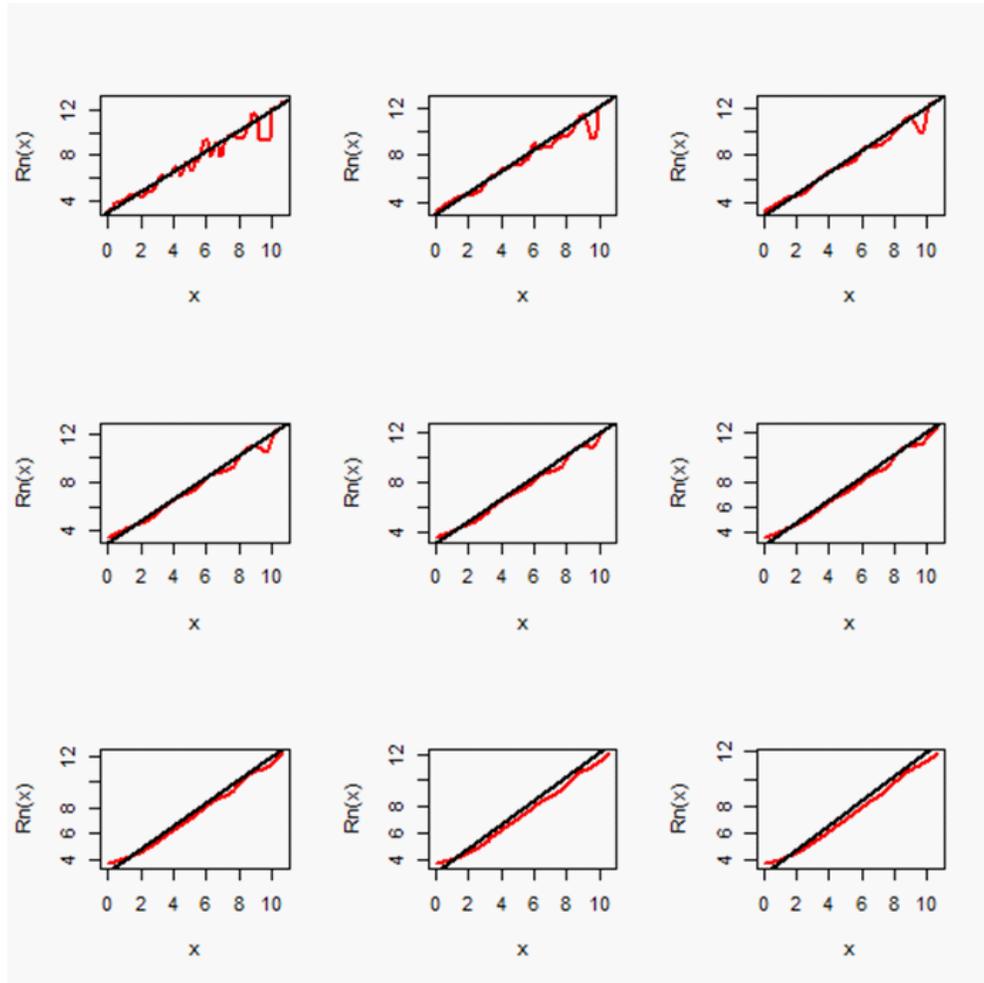
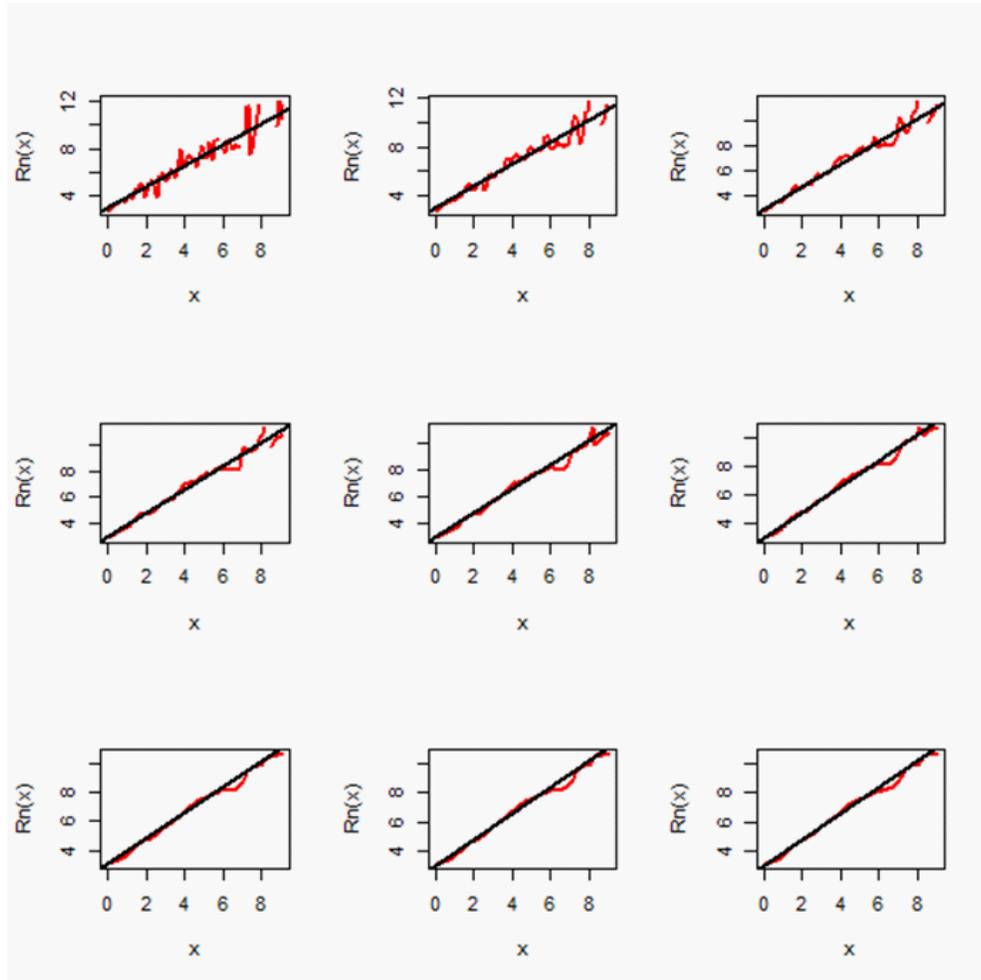


FIG. 3.3 – Régression linéaire avec  $h$  varié,  $n$  fixé et  $K$  noyau gaussien.

Identique aux choix précédents, mais on change le noyau :  $K(t) = \frac{3}{4}(1 - t^2)$  (noyau d'Epanechnikov). On obtenu la figure [FIG-3.4] qui explique l'estimation obtenue avec les valeurs de  $h$  varié de 0.1 à 0.9. Il est claire que la valeur du  $h$  optimale est de  $h = 0.9$  (ligne 3, colonne 3).


 FIG. 3.4 – Régression linéaire avec  $h$  varié,  $n$  fixé et  $K$  d'Epanechnikov

## 3.2 Régression non linéaire

Dans cette section, nous allons répéter les mêmes étapes que dans la régression linéaire mais avec un modèle non linéaire :

$$y = \cos x + \varepsilon. \quad (3.4)$$

où  $\varepsilon$  un terme d'erreur de loi  $\mathcal{N}(0, 1)$ .

Toujours, la ligne noire exprime la fonction de régression théorique  $r(x)$  [Eq.3.1] et la ligne rouge exprime la fonction de régression empirique  $r_n(x)$  donnée par l'équation [Eq.3.2].

### 3.2.1 Paramètre de lissage $h$ fixé, $n$ varié

Dans ce cas, on choisit le paramètre de lissage  $h_n = n^{-1/5}$  (fixé),  $n$  varié ( $n = 40, 200, 400$ ) et  $K$  est un noyau gaussien  $K(t) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}t^2\right)$ .

**Code R :**

```
rm(list=ls(all=TRUE))

n=40

X=rexp(n,0.5)
E=rnorm(n)
Y=cos(X)+E

# Noyau Normale K(t)
K=function(t){(1/sqrt(2*pi))*exp(-0.5*t^2)}

h=n^-.2

# Initiation
s=100 # taille de l'intervalle [a,b]
a=min(X) #borne inf
b=max(X) # borne sup
x=seq(a,b,length=s) # Intervalle [a,b]
V=numeric(n)
fn=numeric(s)
for(j in 1 :s){
  for(i in 1 :n){ V[i]=K((x[j]-X[i])/h) }
  fn[j]=sum(V)/(n*h)}

# Fonction Hn(.)
W=numeric(n)
Hn=numeric(s)
for(j in 1 :s){
  for(i in 1 :n){ W[i]=K((x[j]-X[i])/h)*Y[i] }
```

```
Hn[j]=sum(W)/(n*h)}
Rn =Hn/fn
# Graphes
x11() # nouvelle fenetre graphique
op=par(mfrow=c(1,3))
plot(x,Rn,xlab="x", ylab="Rn(x)", main="n=40",type='l',col=2, lwd= 2)
lines(x,cos(x),lwd= 2)
#####Pour n =200
n=200
X=rexp(n,0.5)
E=rnorm(n)
Y=cos(X)+E
h=n^-.2
V=numeric(n)
for(j in 1 :s){
for(i in 1 :n){ V[i]=K((x[j]-X[i])/h) }
fn[j]=sum(V)/(n*h)}
W=numeric(n)
for(j in 1 :s){
for(i in 1 :n){ W[i]=K((x[j]-X[i])/h)*Y[i] }
Hn[j]=sum(W)/(n*h)}
Rn =Hn/fn
plot(x,Rn,xlab="x", ylab="Rn(x)", main="n=200",type='l',col=2, lwd= 2)
lines(x,cos(x),lwd= 2)
#####Pour n =400
n=400
X=rexp(n,0.5)
```

```
E=rnorm(n)
Y=cos(X)+E
h=n^-.2
V=numeric(n)
for(j in 1 :s){
for(i in 1 :n){ V[i]=K((x[j]-X[i])/h) }
fn[j]=sum(V)/(n*h)}
W=numeric(n)
for(j in 1 :s){
for(i in 1 :n){ W[i]=K((x[j]-X[i])/h)*Y[i] }
Hn[j]=sum(W)/(n*h)}
Rn =Hn/fn
plot(x,Rn,xlab="x", ylab="Rn(x)", main="n=400",type='l',col=2, lwd= 2)
lines(x,cos(x),lwd= 2)
par(op)
```

On obtenu la figure [FIG-3.5], On remarque la même conclusion pour le cas non linéaire que le cas linéaire (*i.e.*, convergence de l'estimateur pour  $n$  assez grand).

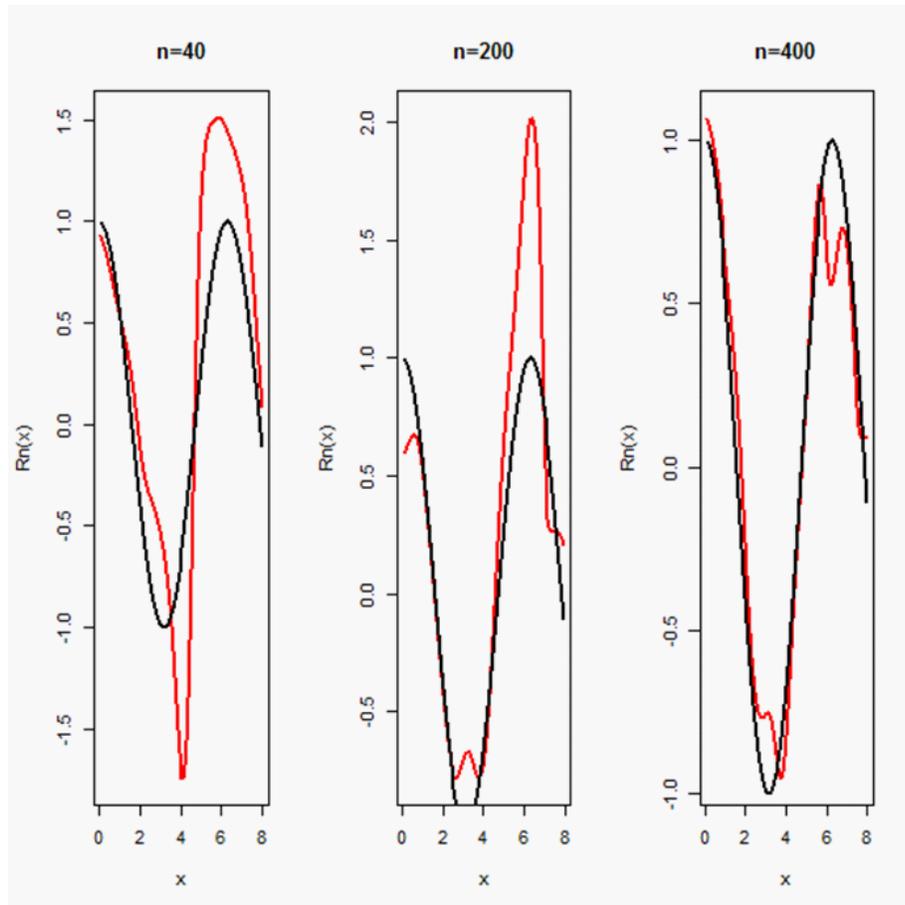


FIG. 3.5 – Régression non linéaire :  $h$  fixé,  $n$  varié et  $K$  noyau normale

Dans ce second cas, on choisit le noyau d'Epanechnikov :  $K(t) = \frac{3}{4}(1 - t^2)$ . Ensuite, on modifie seulement cette partie dans le programme **R** précédent :

```
# Noyau Epanechnikov K(t)
K=function(t){ifelse(abs(t)<1,(3/4)*(1-t^2),0)}
```

On obtient la figure [FIG-3.6], et on arrive à la même conclusion de la convergence de l'estimateur (voir la [FIG-3.5], *i.e.*, convergence de l'estimateur pour  $n$  assez grand).

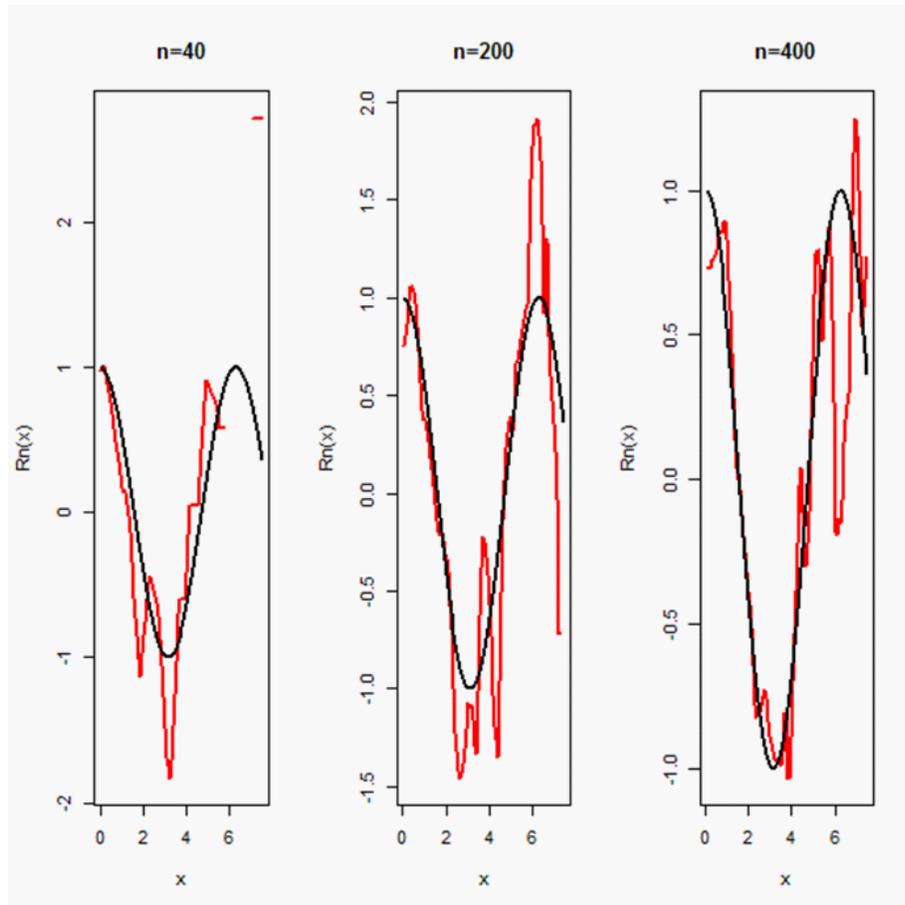


FIG. 3.6 – Régression non linéaire :  $h$  fixé,  $n$  varié et  $K$  noyau d'Epanechnikov

### 3.2.2 Choix graphique du paramètre de lissage

Dans cette partie, on va prendre le paramètre de lissage dans l'intervalle  $(0,1)$  de même façon que pour la régression linéaire, et avec des tests graphique en va diterminer le paramètre  $h$  optimal (au sens graphique).

On fixe la taille de l'échantillon  $n = 350$  et le noyau  $K$  est normale, l'estimation obtenue avec les valeurs de  $h$  varié de 0.1 à 0.9 sont données dans la figure [FIG-3.7]. Il est clair que la valeur du  $h$  optimale est de  $h = 0.5$  (ligne 2, colonne 2).

**Code R :**

```
n=350
X=rexp(n,0.5)
```

```
Y=cos(X)+E
# Noyau Normale K(t)
K=function(t){(1/sqrt(2*pi))*exp(-0.5*t^2)}
h=seq(.1,.9,length=9)
# Initiation
s=100 # taille de l'intervalle [a,b]
a=min(X) #borne inf
b=max(X) # borne sup
x=seq(a,b,length=s) # Intervalle [a,b]
V=array(dim=c(n,s,9))
fn=array(dim=c(s,9))
W=array(dim=c(n,s,9))
Hn=array(dim=c(s,9))
for(k in 1 :9){
  for(j in 1 :s){
    for(i in 1 :n){ V[i,j,k]=K((x[j]-X[i])/h[k]) }
    fn[j,k]=sum(V[,j,k])/(n*h[k])}}
# fonction Hn(x)
for(k in 1 :9){
  for(j in 1 :s){
    for(i in 1 :n){ W[i,j,k]=K((x[j]-X[i])/h[k])*Y[i] }
    Hn[j,k]=sum(W[,j,k])/(n*h[k])}}
Rn=array(dim=c(s,9))
for(k in 1 :9){ Rn[,k]=Hn[,k]/fn[,k]}
# Graphes
x11() # nouvelle fenetre graphique
op=par(mfrow=c(3,3))
```

```

for(k in 1 :9){
plot(x,Rn[,k],xlab="x", ylab="Rn(x)", main=" ",type='l',col=2, lwd= 2)
lines(x,cos(x),lwd= 2)
}
par(op)
    
```

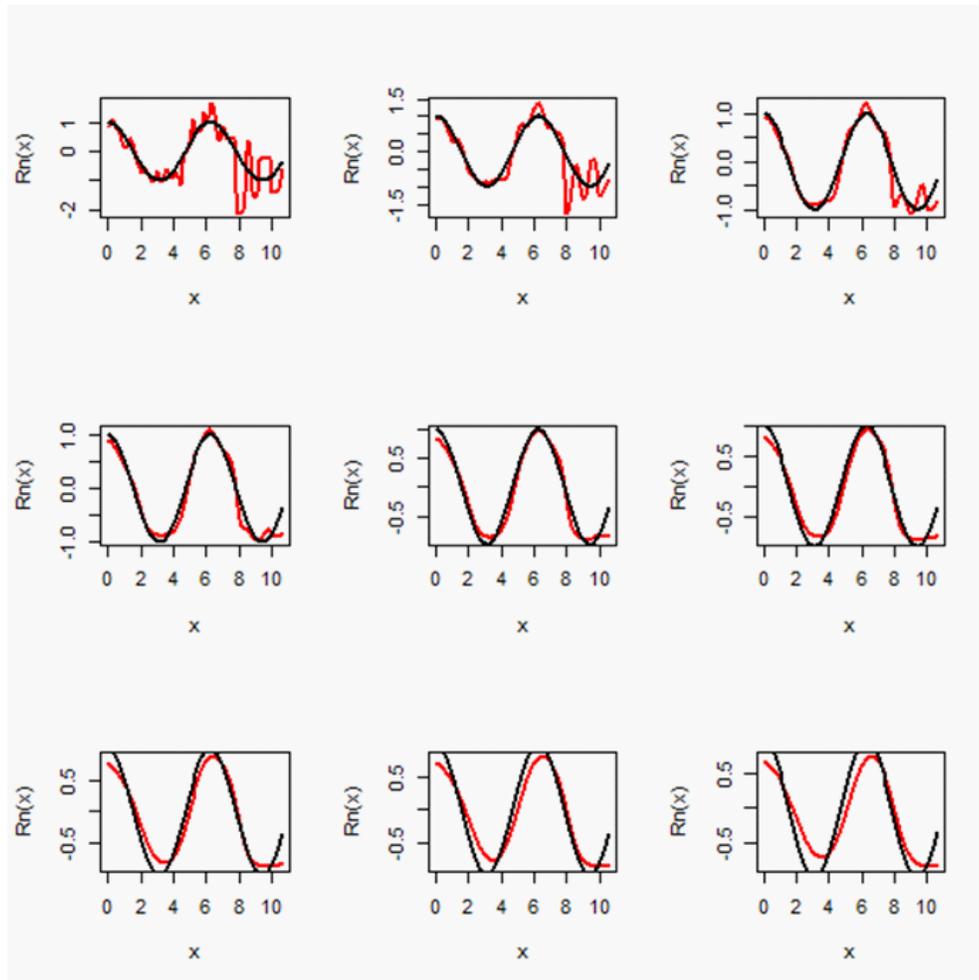


FIG. 3.7 – Régression non linéaire avec  $h$  varié,  $n$  fixé et  $K$  gaussien.

Si nous gardons le même modèle non linéaire  $y = \cos x + \varepsilon$ , mais avec le noyau d'Epanechnikov. On note, que la valeur du  $h$  optimale est de  $h = 0.9$  (ligne 3, colonne 3, voir la FIG-3.8).

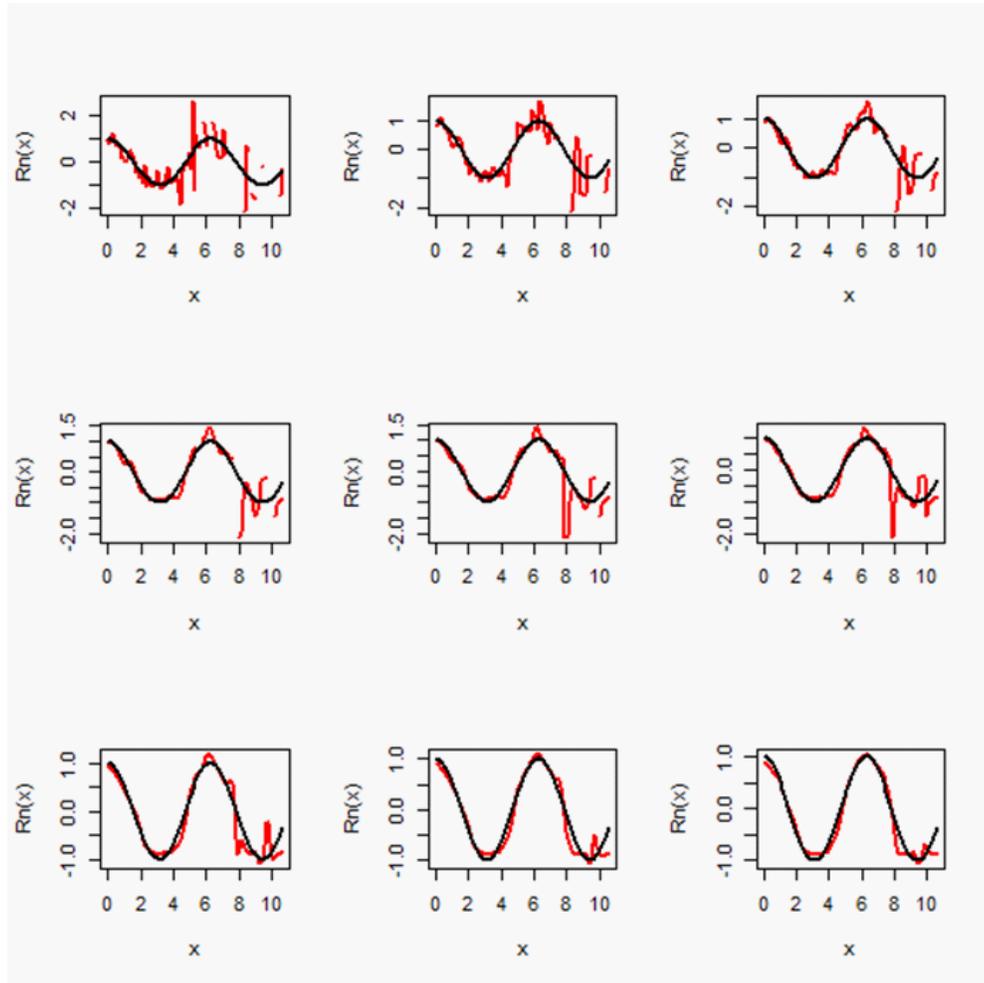


FIG. 3.8 – Régression non linéaire avec  $h$  varié,  $n$  fixé et  $K$  d'Epanechnikov.

Finallement, ce chapitre montre l'importance de paramètre de lissage  $h$  et du noyau  $K$  dans l'estimation non paramétrique de la régression linéaire et non linéaire. Mais à noter que le choix de  $h$  est plus crucial que le choix de noyau.

# Conclusion

Dans ce mémoire, on a présenté la méthode d'estimation à noyau, qui permettant d'effectuer de la régression non paramétrique. Ce travail a montré que la méthode d'estimation de régression non paramétrique est simple et peut être très utile dans plusieurs situations. Par exemple, dans l'analyse des données, lorsque l'on désire comprendre et observer les relations qui existent entre les variables.

Dans la régression non paramétrique, la méthode du noyau joue un grand rôle. Pour que son soit plus utilisée par les praticiens, il est nécessaire que les programmes informatiques permettant d'appliquer ces méthodes soient facilement accessibles et assez simples d'utilisation. Cela favorise aussi les échanges entre statisticiens et utilisateur.

L'estimateur à noyau de la régression non paramétrique dépend de deux paramètres le noyau  $K$  et le paramètre de lissage  $h$ . Dans la pratique, on utilisé le logiciel **R** pour présenté des exemples sur cet estimateur, et à travers les résultats obtenus, nous concluons que : le noyau  $K$  est peu influence sur l'estimateur, par contre le paramètre  $h$  est un grand influence, et dont le choix est crucial .

# Bibliographie

- [1] Benchoulak, H. (2012). Bandes de confiance pour les fonctions de densité et de régression.
- [2] Blondin, D. (2004). Lois limites uniformes et estimation non-paramétrique de la régression (Doctoral dissertation, Université Pierre et Marie Curie-Paris VI).
- [3] Bochner, S. (1955). Harmonic Analysis and the Theory of probability. University of Chicago Press, Chicogo, Illinois.
- [4] Bosq, D. and Lecoutre, J. P., (1987). Théorie de l'Estimation Fonctionnelle. Economica, Paris.
- [5] Cruz, C. M. R. T. D. (1995). Estimation fonctionnelle : applications aux tests d'adéquation et de paramètre constant (Doctoral dissertation).
- [6] Durrieu, G., Nguyen, T. M. N., & Sow, M. (2009). Comparaison d'estimateurs de régression non paramétriques : application en valvometrie. In 41èmes Journées de Statistique, SFdS, Bordeaux.
- [7] Ferraty, F. et Vieu, P. (2002/2003). Statistique fonctionnelle : Modèles Non paramétrique de Régréssion. Cours de DEA.
- [8] Härdle, W. (1990). Applied Nonparametric Regression. Cambridge University Press, Cambridge.
- [9] Kebabi, K. (2014). Estimation non-paramétrique de la fonction de régression.
- [10] Matias, C., & Atelier, S. F. D. S. (2013). Introductiona la statistique non paramétrique. Laboratoire Statistique et Génome, Évry, 2014.

- [11] Nadaraya, E. A. (1964). On estimating Regression. *Theory. Probab. Applic.*, 9, 141-142.
- [12] Nadaraya, E. A. (1989). *Nonparametric Estimation of Probability Densities and Regression Curves*. Kluwer, Dordrecht.
- [13] Parzen, E. (1962). On estimation of a probability density function and mode. *The annals of mathematical statistics*, 33(3), 1065-1076.
- [14] Rosenblatt, M. (1956). Remarks on some nonparametric estimates of a density function. *The Annals of Mathematical Statistics*, 832-837.
- [15] Schuster, E. F. (1972). Joint asymptotic distribution of the estimated regression function at a finite number of points. *Annals of Mathematical Statistics*, 43.1, 84-88.
- [16] Tsybakov, A. B. (2003). *Introduction à l'estimation non paramétrique* (Vol. 41). Springer Science & Business Media.
- [17] Watson, G. S. (1964). Smooth Regression analysis. *Sankhyà Ser. A*, 26, 359-372.

# Annexe A : Rappels

**A.1- Convergence presque complète :** On dit que la suite de variables aléatoires réelles  $(X_n)_{n \in \mathbb{N}}$  converge presque complètement vers une variable aléatoire  $X$  lorsque  $n \rightarrow \infty$  (et on note  $\lim_{n \rightarrow \infty} X_n = X$  *p.co*), si et seulement si :

$$\forall \varepsilon > 0, \quad \sum_{n \in \mathbb{N}} P[|X_n - X| > \varepsilon] < \infty.$$

**A.2- Convergence presque sûre :** On dit que la suite de v.a  $(X_n)_{n \in \mathbb{N}}$  définie sur  $(\Omega, \mathcal{A}, P)$ , converge presque sûrement (*p.s.*) vers la variable aléatoire  $X$  définie sur  $(\Omega, \mathcal{A}, P)$ , si

$$P \left\{ \omega \in \Omega : \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega) \right\} = 1$$

Dans ce cas, on note  $\lim_{n \rightarrow \infty} X_n = X$  *p.s* ou  $X_n \xrightarrow{p.s} X$  lorsque  $n \rightarrow \infty$ .

**A.3- Convergence en probabilité :** On dit que la suite de v.a  $(X_n)$  converge en probabilité vers une v.a.  $X$  si, pour tout  $\varepsilon > 0$  :

$$P[|X_n - X| < \varepsilon] \rightarrow 1 \quad \text{quand } n \rightarrow \infty, \quad \text{on écrit alors, } X_n \xrightarrow{P} X.$$

**A.4- Convergence en loi :** On dit que la suite de v.a  $(X_n)$ , de fonction de répartition  $F_n$ , converge en loi vers une v.a  $X$  de fonction de répartition  $F$ , si la suite  $(F_n(x))$  converge vers  $F(x)$  en tout point  $x$  où  $F$  est continue :  $X_n \xrightarrow{\mathcal{L}} X$ , quand  $n \rightarrow \infty$ .

**A.5- Convergence en moyenne quadratique :** On dit que la suite de v.a  $(X_n)$  converge en moyenne quadratique vers une v.a  $X$  si :

$$E |X_n - X|^2 \longrightarrow 0, \quad \text{quand } n \longrightarrow \infty, \quad \text{on écrit alors, } X_n \xrightarrow{m.q.} X.$$

**A.6- Théorème (Lois des grands nombres) :** Si  $(X_1, \dots, X_n)$  est un échantillon provenant d'une v.a.  $X$  telle que  $E |X| < \infty$ , alors :

$$\bar{X}_n \xrightarrow{\mathcal{P}} E(X) \quad \text{quand } n \longrightarrow \infty, \quad (\text{loi faible})$$

$$\bar{X}_n \xrightarrow{p.s.} E(X) \quad \text{quand } n \longrightarrow \infty, \quad (\text{loi forte}).$$

**A.7- Définition (Biais d'un estimateur) :** Un estimateur  $\hat{\theta}_n$  de  $\theta$  est dite sans biais si pour tout  $\theta \in \Theta$  et tout entier positif  $n$  :  $E(\hat{\theta}_n) = \theta$ .

De même,  $\hat{\theta}_n$  est dite asymptotiquement sans biais si pour tout  $\theta \in \Theta$  :

$$E(\hat{\theta}_n) \longrightarrow \theta, \quad \text{quand } n \longrightarrow \infty.$$

La quantité :  $E(\hat{\theta}_n) - \theta$ , est appelée le biais de l'estimateur  $\hat{\theta}_n$ .

**A.8- Définition ( $o(h_n)$ ,  $O(h_n)$ ) :** Soit  $x_n$  et  $y_n$  deux suites de nombres réels. Alors, lorsque  $n \rightarrow \infty$ ,

$$i) \quad x_n = O(y_n) \Leftrightarrow \limsup_{n \rightarrow \infty} |x_n/y_n| < \infty,$$

$$ii) \quad x_n = o(y_n) \Leftrightarrow \lim_{n \rightarrow \infty} |x_n/y_n| = 0.$$

# Annexe B : Notations et Abréviations

Les différentes notations et abréviations utilisées tout au long de ce mémoire sont expliquées ci-dessous :

$h := h_n$	Paramètre de lissage ou Fenêtre
$h_{opt}$	Fenêtre optimale
$K(\cdot)$	Noyau
$\mathcal{A}$	Tribu
$\mathbb{E}$	Espace des observations
$\mathcal{B}$	Tribu borélienne
$\mathbb{P}$	une famille de probabilité
$\Theta$	Ensemble des paramètres
$L^1$	Espace des fonctions intégrables
<i>iid</i>	Indépendantes et identiquement distribuées
$E$	Espérance de probabilité
<i>Biais</i>	Biais d'un estimateur
<i>Var</i>	Variance d'un estimateur
$f_X$	Densité de $X$
$F$	Fonction de répartition

$F^{-1}$	Fonction des quantiles
$f_{n,X}$	Estimateur de $f$
$v.a$	Variable aléatoire
$r$	Fonction de regression
$r_n$	Estimateur de $r$
$1(\cdot)$	Fonction indicatrice
$\ \cdot\ $	Norme euclidienne
$\xrightarrow{\mathcal{P}}$	Convergence en probabilité.
$\xrightarrow{\mathcal{L}}$	Convergence en loi.
$\xrightarrow{p.s.}$	Convergence presque sûre.
$\xrightarrow{m.p}$	Convergence en moyenne d'ordre $p$ .
$p.co$	Convergence presque complète.
$MSE$	L'erreur quadratique moyenne (mean square error).
$AMSE$	L'erreur quadratique moyenne asymptotique (asymptotic mean square error).
$MISE$	L'erreur quadratique intégrée moyenne (mean integrated squared error).