

République Algérienne Démocratique et Populaire
Ministère de l'Enseignement Supérieur et de la Recherche Scientifique
Université Abderrahmane Mira de Béjaïa
Faculté des Sciences Exactes
Département de Recherche Opérationnelle



Mémoire Présenté
Pour l'obtention du Diplôme de Master
En Recherche Opérationnelle
Option : Mathématiques Financières

Par : BERKANE Zahra
Et : BAROUDI Lidia

Régression non paramétrique : Méthodes d'estimation

Soutenu à l'Université Abderrahmane Mira de Béjaïa,
Le 02/07/2023, devant le jury composé de :

<i>M^{me}</i> BERNINE Nassima	Présidente	à l'UAMB - Bejaia.
<i>M^{me}</i> AMROUN Sonia	Encadreur	à l'UAMB - Bejaia.
<i>M^{me}</i> DJERROUD Lamia	Examinatrice	à l'UAMB - Bejaia.
<i>M^{me}</i> HARFOUCHE Lynda	Examinatrice	à l'école - amizour.

Année Universitaire 2022 – 2023

Remerciment

Tout d'abord on aimera d'adresser nos plus sincères remerciements au Dieu le tout puissant et le miséricordieux de nous avoir donné la chance, la patience et le courage pour réaliser ce modeste travail.

On tient à exprimer nos vifs remerciements à madame Sonia AMROUNE pour avoir accepté de nous encadrer lors du présent travail en dépit de son emploi du temps très chargé et de la confiance qu'il nous a témoignée, et les précieux conseils qu'elle a bien voulu prodiguer pour cibler les aspects traités dans ce travail.

On tient à remercier aussi mr ASLI, KABYLE, ABAS, BIBI, et madame TABTI et DJERROUD qui nous donnent des conseils très importants en signe de reconnaissance. On adresse nos sincères remerciements à tous les professeurs, intervenants et toutes les personnes qui par leurs paroles, leurs écrits, leurs conseils et leurs critiques ont guidé nos réflexions et ont accepté à nous rencontrer et répondre à nos questions durant la recherches.

Enfin, on ne voudra pas oublier de remercier tous nos amis qu'ont été d'un soutien que ce soit moral ou matériel, et qui ont contribué de près ou de loin à m'aider.

Beaucoup de charité et Bonne chance à tous.

Merci

Dédicace

*Du profond de mon coeur,
je dédie ce travail à tous ceux qui me sont chers.*

*A mes parents (Larbi, Zahia)
qui nous ont soutenu tout le long de mon parcours d'études.*

Mes grands-parents, que Dieu leur fasse miséricorde.

Mon frère : Farid.

Mes soeurs : Fatima, Assia.

*Ma petite cousine Thafsouth que Dieu lui
accorde le paradis.*

Mes très chers et précieux amis (Salhi, Ark, Yasmine, Nesrine, Mira, ...)

Ma binôme Zahra.

*En souvenir des moments heureux passés ensemble, avec mes vœux sincères de réussite,
bonheur, et la santé.*

Enfin, à tous les professeurs et les enseignants qui m'ont appris.

Lidia

Dédicace

Je dédie ce modeste travail À Mes chers parents,

Grâce à leurs tendres encouragements et leurs grands sacrifices, ils ont pu créer le climat affectueux et propice à la poursuite de mes études.

À mes deux frère mourad et mouhand.

Ma soeur Rosa.

A mes chers grands-parents (Setti et Jeddi)et mes cousines .

À tous mes chers amis (Romaiassa, Cylia, Yamina, Djajiga, Sabrina, Imane, Rahma...).

Ma binome Lidia Baroudi.

A mes collègues de l'Université de Abd mira de Bejaia.

*Et à tous ce qui m'ont enseigné au long de
ma vie scolaire.*

Zahra

Table des matières

Liste des figures	IV
Liste des tables	V
Notations	1
Introduction générale	2
1 Méthodes d'estimation de la régression non paramétrique	4
Introduction	4
1.1 Régression paramétrique	4
1.1.1 Modèle de régression linéaire simple	4
1.1.2 Estimation des paramètres par la méthode des moindres carrés	5
1.2 Régression non-paramétrique	5
1.3 Notions sur l'estimateur	6
1.3.1 Estimateur sans biais	6
1.3.2 Estimateur asymptotiquement sans biais	6
1.3.3 Estimateur asymptotiquement uniformément sans biais	6
1.3.4 Espérance mathématique	6
1.3.5 Variance	7
1.3.6 Erreur moyenne quadratique (MSE)	7
1.3.7 Erreur moyenne quadratique intégrée (MISE)	7
1.3.8 Intégrale des erreurs quadratiques (ISE)	8
1.3.9 Convergence en loi	8
1.4 Méthodes d'estimation de la régression non paramétrique	8
1.4.1 Méthode des k-plus proche voisin	8
1.4.2 Méthode du noyau	9
1.4.3 Méthode des splines	13
1.4.4 Méthode des séries orthogonale	15
Conclusion	17
2 Estimation non paramétrique par la méthode du noyau et par la méthode des splines de lissages	18
Introduction	18

2.1	Estimation non paramétrique par la méthode noyau	18
2.2	Propriétés asymptotiques de l'estimateur à noyau	19
	2.2.1 Choix du paramètre de lissage	21
	2.2.2 Normalisation asymptotique	24
2.3	Estimation non paramétrique par la méthode des fonctions splines	24
	2.3.1 Fonction Spline	24
	2.3.2 Fonction Spline Naturelle	24
	2.3.3 Spline Cubique Naturelle	25
	2.3.4 Interpolation par les fonctions splines	25
	2.3.5 Existence et unicité des spline d'interpolation	29
	2.3.6 Splines de lissage	30
	2.3.7 Existence et unicité de la spline de lissage minimisante	31
	2.3.8 Propriétés de l'estimateur splines de lissage	31
	2.3.9 Propriétés asymptotiques de l'estimateur	32
	2.3.10 Choix du paramètre de lissage	33
	Conclusion	34
3	Application	35
3.1	Application sur \mathbb{R}	35
3.2	Etude comparatif	36
	3.2.1 Méthode de noyau	36
	3.2.2 Méthode des splines	36
3.3	Modèles de fonction	37
	3.3.1 Modèle $f_1(x)$	37
	3.3.2 Modèle $f_2(x)$	41
	3.3.3 Cas réel	45
	Conclusion	46
	Conclusion générale	47
	Bibliographie	50

Table des figures

1.1	Les courbes de certains noyau	12
1.2	Un dessin en zig-zag (à gauche) et en splines (à droite)	13
3.1	Estimation de f_1 , $n = 50$	38
3.2	Estimation de f_1 , $n = 100$	38
3.3	Estimation de f_1 , $n = 200$	39
3.4	Estimation de f_1 , $n = 500$	39
3.5	Estimation de f_1 , $n = 1000$	40
3.6	Estimation de f_1 , $n = 2000$	40
3.7	Estimation de f_2 , $n = 50$	42
3.8	Estimation de f_2 , $n = 100$	42
3.9	Estimation de f_2 , $n = 200$	43
3.10	Estimation de f_2 , $n = 500$	43
3.11	Estimation de f_2 , $n = 1000$	44
3.12	Estimation de f_2 , $n = 2000$	44
3.13	Estimation de la courbe d'éoliennes par la méthode du noyau et la méthode des fonctions splines.	46

Liste des tableaux

3.1	ASE donnée par les deux méthodes associée au modèle $f_1(x)$ en fonction de la taille de l'échantillon n	37
3.2	ASE donnée par les deux méthodes associée au modèle $f_2(x)$ en fonction de la taille de l'échantillon n	41
3.3	Données sur les éoliennes.	45
3.4	ASE associée au cas réel.	45

Notations

$f(\cdot)$	Fonction de régression.
$\hat{f}(\cdot)$	Estimateur de fonction de régression.
$m(\cdot)$	Densité.
$M(\cdot)$	Fonction de répartition.
$\hat{m}(\cdot)$	Estimateur de la densité.
$K(\cdot)$	Fonction noyau.
h	Paramètre de lissage ou Fenêtre.
h_{opt}	Fenêtre optimale.
λ	Paramètre de lissage positif.
E	Espérance de probabilité.
Var	Variance d'un estimateur.
<i>Biais</i>	Biais d'un estimateur.
<i>ASE</i>	Erreur moyenne quadratique.
<i>MISE</i>	Erreur moyenne quadratique intégrée.
<i>ISE</i>	Intégrale des erreurs quadratique.
<i>AMISE</i>	Erreur moyenne quadratique intégrée asymptotique.
<i>CV</i>	Validation croisé.
<i>GCV</i>	Validation croisé généralisée.
$N(p)$	Voisinage de p .
ϵ_i	Erreurs aléatoires de moyenne nulle et de variance σ^2 .
$\ \cdot\ $	Norme euclidienne.

Introduction générale

La littérature statistique propose deux types de modèles de régression : la régression paramétriques et la régression non paramétriques. La régression non paramétrique est une méthode d'analyse statistique qui permet d'estimer la relation entre une variable dépendante Y et une variable explicative X sans supposer de forme spécifique pour cette relation. Cette méthode est largement utilisée dans divers domaines tels que l'économie, la biostatistique et les sciences de l'environnement. Son développement remonte au 19^{ème} siècle, et son objectif est d'estimer la dépendance entre les variables sans contraindre sa forme.

Soit les observations suivantes $(x_1, y_1), \dots, (x_n, y_n)$. Le modèle de régression étudié dans ce travail est de la forme :

$$y_i = f(x_i) + \epsilon_i, i = 1, \dots, n,$$

ou f est la fonction de régression que l'on cherche à estimer et les $(\epsilon_i)_{i=1}^n$ sont des erreurs aléatoires supposées que est de loi normal $(0, \sigma^2)$.

Plusieurs méthodes ont été proposées pour estimer la fonction de régression f , la première méthode rencontrée dans la littérature est la méthode du noyau ; également connue sous le nom de régression lissage par l'opérateur du noyau. L'estimateur noyau proposé par [Nadaraya (1964)][25] et de [Watson (1964)][48] repose sur l'utilisation d'une fonction noyau K et d'une fenêtre h , sous la forme suivante qui est un estimateur linéaire par rapport Y :

$$\hat{f} = \frac{\sum_{i=1}^n Y_i K\left(\frac{X_i - x}{h_n}\right)}{\sum_{i=1}^n K\left(\frac{X_i - x}{h_n}\right)}.$$

Une autre méthode rencontrée dans la littérature est la méthode des fonctions splines, cette méthode est une technique d'analyse numérique utilisée pour l'interpolation. Elle permet de décrire une courbe à travers un ensemble fini de points donnés. Initialement introduite par [Whittaker(1923)][49] et développée par [Schoenberg(1964)][34], elle est largement utilisée dans le calcul scientifique et l'approximation.

Les splines de lissage sont désormais couramment utilisées en infographie pour représenter des courbes et des surfaces. Les travaux de [Wahba][41] ont démontré leur utilité pour résoudre des problèmes d'estimation statistique, ce qui les a rendues populaires dans des domaines tels que l'analyse de données de croissance, la médecine et l'économie. L'estimateur par la méthode des splines de lissage est défini comme étant le fonction \hat{f}_α

qui réalise le minimum de $s(\alpha)$, qui s'écrit sous la forme

$$\hat{f}_\alpha = A_\alpha Y;$$

avec A_α une matrice lisse qui ne dépend que des x_i , $i = 1, \dots, n$, l'estimateur est linéaire par rapport au vecteur Y .

Ce travail est subdivisé en trois chapitres :

Dans le premier chapitre nous introduirons le modèle de régression non paramétrique ainsi que les méthodes d'estimations de cette fonction (k-plus proche voisin, noyau, spline, et la méthode des séries orthogonal).

Dans le deuxième chapitre, nous présentons les méthodes du noyau et la méthode des splines en détails ainsi que les propriétés des estimateurs.

Dans le troisième chapitre, nous présentons les résultats obtenus par la simulation effectuée sur les exemples de fonction et un cas réel pour comparer les deux méthodes (noyau et spline), le critère utilisé est celui de ASE.

Ce mémoire se termine par une conclusion générale et quelques perspectives de recherche, suivie d'une bibliographie.

1

Méthodes d'estimation de la régression non paramétrique

Introduction

La théorie de l'estimation est une des branches les plus basiques de la statistique. Cette théorie est habituellement divisée en deux composantes principales : régression paramétrique et régression non paramétrique.

1.1 Régression paramétrique

La régression paramétrique consiste à supposer que f appartient à une famille de fonctions continues ou discrète qui peuvent être décrites par un certain nombre de paramètres réels. Le but de la régression paramétrique linéaire simple est d'expliquer une variable Y à l'aide d'une variable X . La variable Y est appelée variable dépendante (à expliquer) et X variable indépendante (explicative). Elle consiste à estimer des paramètres réels finis (nombres finis).

1.1.1 Modèle de régression linéaire simple

Dans un modèle de régression paramétrique ; la fonction de régression est :

- (i) De forme explicite.
- (ii) Peut s'écrire en fonction d'un nombre réduit de paramètres.

Considérons un échantillon formé d'une suite de n couple (x_i, y_i) ; $i = 1, \dots, n$ à valeurs dans R telque :

$$y_i = a + bx_i + \epsilon_i ; \quad (1.1)$$

où les erreurs aléatoires ϵ_i sont non corrélées, de moyenne nulle et de variance σ^2 .

Pour estimer y_i on cherche alors à déterminer les meilleures valeurs de a et b compte tenu d'un critère.

1.1.2 Estimation des paramètres par la méthode des moindres carrés

Les paramètres sont estimés avec la méthode des moindres carrés, son principe est de minimiser la somme des carrés des résidus en fonction de a et b .

$$\min_{a,b} Q(a,b) = \min \sum_{i=1}^n \epsilon_i^2 = \min \sum_{i=1}^n (y_i - a - bx_i)^2 ; \quad (1.2)$$

1- Calcul des estimateurs :

La première dérivé par rapport à a et b est :

$$\begin{cases} \frac{\partial Q(a,b)}{\partial a} = -2 \sum_{i=1}^n (y_i - a - bx_i) = 0 ; & (*) \\ \frac{\partial Q(a,b)}{\partial b} = -2 \sum_{i=1}^n x_i (y_i - a - bx_i) = 0 . & (**) \end{cases}$$

L'équation (*) donne :

$$\hat{a} = \bar{y} - b\bar{x} ;$$

l'équation (**) donne :

$$\hat{b} = \frac{\sum_{i=1}^n (x_i - \bar{x}) * (y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} ;$$

telle que :

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i. \quad \text{et} \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i.$$

1.2 Régression non-paramétrique

C'est un outil statistique permettant de d'écrire une relation entre une variable dépendante et une ou plusieurs variables explicatives, le problème de l'estimation non paramétrique consiste, dans la majeure partie des cas, à estimer à partir des observations, une fonction inconnue élément d'une certaine classe fonctionnelle. Plus particulièrement, on parle d'estimation non paramétrique lorsque celle-ci ne se ramène pas à l'estimation

d'un nombre fini de paramètres réels associés à la loi de l'échantillon. Quand on veut écrire l'effet d'une variable sur un évènement, faire le moins d'hypothèse possible sur la forme de la relation. Le but du modèle de régression est de déterminer la variable dépendante Y dépend d'un ensemble de variable explicatives X . Supposer $X \in \mathbb{R}$. Ainsi, le problème est de déterminer pour chaque réalisation x de la variable X , la valeur de la fonction $f(x)$ appelée fonction de régression.

1.3 Notions sur l'estimateur

1.3.1 Estimateur sans biais

On dit qu'un estimateur \hat{f} de f est sans biais si : $E(\hat{f}) = f$.

1.3.2 Estimateur asymptotiquement sans biais

On dit qu'un estimateur \hat{f} de f est asymptotiquement sans biais si : $\lim_{n \rightarrow +\infty} (\hat{f}) = f$.

1.3.3 Estimateur asymptotiquement uniformément sans biais

Un estimateur \hat{f} de f est dit asymptotiquement uniformément sans biais si :

$$\lim_{n \rightarrow +\infty} \sup_x |E[\hat{f}(x) - f(x)]| = 0 ;$$

1.3.4 Espérance mathématique

Soit une variable aléatoire X absolument continue de densité de probabilité $f(x)$ définie sur $[a, b]$.

L'espérance mathématique de X est définie par :

$$E(X) = \int_a^b x f(x) dx ;$$

On a la formule suivante pour une variable aléatoire discrète X de loi de probabilité (p_k) :

$$E(X) = \sum_{i=1}^{+\infty} p_i \times x_i ;$$

où $p_i = P(X = x_i)$

1.3.5 Variance

La variance d'une variable aléatoire X absolument continue est définie par :

$$V = E(X^2) - [E(X)]^2 ;$$

Si X une variable aléatoire discrète de loi (p_i) . La variance de X se calcule ainsi :

$$V(X) = \sum_{i=1}^{+\infty} p_i \times x_i^2 - [E(X)]^2 ;$$

1.3.6 Erreur moyenne quadratique (MSE)

L'erreur moyenne quadratique MSE :

$$\begin{aligned} MSE(f(x), \hat{f}(x)) &= E(f(x) - \hat{f}(x))^2, \\ &= E(f(x))^2 + E(\hat{f}(x))^2 - 2Ef(x)\hat{f}(x), \\ &= E(f(x))^2 + E(\hat{f}(x))^2 - 2Ef(x)\hat{f}(x) + [E\hat{f}(x)]^2 - [E(\hat{f}(x))]^2, \\ &= (f(x))^2 + E(\hat{f}(x))^2 - 2f(x)E\hat{f}(x) + [E\hat{f}(x)]^2 - [E(\hat{f}(x))]^2, \\ &= [E(\hat{f}(x) - f(x))]^2 + E(\hat{f}(x))^2 - [E\hat{f}(x)]^2, \\ &= \text{Biais}(\hat{f}(x))^2 + \text{Var}(\hat{f}(x)). \end{aligned}$$

Définition (1.2) :

On dit qu'un estimateur \hat{f} de f est ponctuellement consistant en moyenne quadratique si :

$$\lim_{n \rightarrow +\infty} MSE(f(x), \hat{f}(x)) = 0 ;$$

1.3.7 Erreur moyenne quadratique intégrée (MISE)

L'erreur moyenne quadratique intégrée MISE :

$$\begin{aligned} MISE(\hat{f}, f) &= \int MSE(f(x), \hat{f}(x)) dx ; \\ &= \int \text{Biais}(\hat{f}(x))^2 dx + \int \text{Var}(\hat{f}(x)) dx ; \end{aligned}$$

Définition (1.3) :

On dit qu'un estimateur \hat{f} de f est uniformément consistant en moyenne quadratique intégrée si :

$$\lim_{n \rightarrow +\infty} MISE(f(x), \hat{f}(x)) = 0 ;$$

1.3.8 Intégrale des erreurs quadratiques (ISE)

l'intégrale des erreurs quadratiques (ISE) est défini par :

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

1.3.9 Convergence en loi

On dit qu'un estimateur \hat{f} de f est asymptotiquement normal si :

$$\hat{f} \rightarrow N(E(\hat{f}), Var(\hat{f})) ; \text{ en loi.}$$

1.4 Méthodes d'estimation de la régression non paramétrique

La particularité de la statistique non-paramétrique est que la fonction inconnu qu'on cherche à calculer, à estimer ou à classifier n'est pas supposé d'appartenir à une famille indiquée par un nombre de paramètres réels. En général, dans la théorie non-paramétrique on suppose que le nombre de paramètres qui décrivent la loi des observations est une fonction croissante du nombre d'observations, ou encore que le nombre de paramètres est infini. Dans cette partie on présente quelques méthodes d'estimation non paramétrique de la fonction de régression.

1.4.1 Méthode des k-plus proche voisin

C'est une méthode non paramétrique utilisée pour la classification et la régression. Dans les deux cas, il s'agit de classer l'entrée dans la catégorie à laquelle appartient les k-plus proche voisins dans l'espace des caractéristiques identifiées par l'apprentissage. L'estimateur k-plus proche voisins est une moyenne pondérée dans un voisinage variable. Ce voisinage est défini par ces variables X qui sont parmi les k-plus proches voisins, en trouve le X par la distance euclidienne. La suite de ces points à été introduit par [Loftsgaarden and Quesenberry] [22] dans le but de l'estimation de la densité, et [cover and hart] [8] dans le but de classification. L'estimateur des k-plus proches voisins est défini par :

$$\hat{f}_k(x) = n^{-1} \sum_{i=1}^n \omega_{ki}(x) y_i ; \quad (1.3)$$

où $\omega_{ki}(x)$ est une suite des poids des k-plus proche voisin, n est la taille d'échantillon.

Soit J_i l'ensemble d'indice ; $J_x = \{i : x_i \text{ est l'une des k observations les plus proches de } x\}$, telque :

$$\omega_{ki}(x) = \begin{cases} \frac{n}{k}, & \text{si } i \in J_i, \\ 0, & \text{sinon.} \end{cases} \quad (1.4)$$

Exemple

Soit $(x_i, y_i)_{i=1}^5 = (1, 5), (7, 12), (3, 1), (2, 0), (5, 4)$ et calculons les k -plus proches voisins $\hat{f}_k(x)$ pour $x = 4$ et $k = 3$. Les observations proches de x sont les 3 derniers points des données. Par conséquent $J_x = J_4 = 3, 4, 5$ et donc

$$w_{k1}(4) = w_{k2}(4) = 0, w_{k3}(4) = w_{k4}(4) = w_{k5}(4) = \frac{5}{3}$$

ce qui résulte

$$\hat{f}_3(4) = \frac{1 + 0 + 4}{3} = \frac{5}{3}$$

proposition (1.1) [Lai][21]

soit $k \rightarrow \infty$, $\frac{k}{n} \rightarrow \infty$, $n \rightarrow \infty$. Le biais et la variance de l'estimateur k -plus proches voisins \hat{f}_k sont donnés par :

$$E(\hat{f}(x)) - f(x) \approx \frac{1}{24f(x)^3} [(\hat{f}'' f_X + 2\hat{f}' f'_X)(x)] \left(\frac{k}{n}\right)^2;$$

$$Var(\hat{f}(x)) \approx \frac{\sigma^2(x)}{k} ;$$

Le compromis entre le *biais*² est ainsi réalisé dans un sens asymptotique en posant $k \sim n^{\frac{4}{5}}$. Une conséquence est que l'erreur moyenne quadratique elle même converge vers 0 au taux de $k \sim n^{\frac{4}{5}}$.

1.4.2 Méthode du noyau

C'est la méthode la plus connue et plus utilisée, son succès peut s'expliquer par l'expression théorique de l'estimateurs, et sa convergence dans différent sens. La méthode du noyau est une méthode simple et très pratique, lorsqu'en s'intéresse à la relation entre une variable réponse Y est une variable explicative X . L'estimation de la fonction de régression à été proposée par [Nadaraya (1964)][25] et [watson (1964)][48]. Cette approche non paramétrique par noyau a été aussi développé par [Ferraty and vieu(2002) et (2003)][13]. Cet estimation est basé sur le calcul de la moyenne pondérée des observations pour toutes les valeurs du domaine. Dans le cas multivarié continue, L'estimateur à noyau continu symétrique de la fonction de densité m est défini par [Parzen [1962)][27] :

$$\hat{m}_h(x) = \frac{1}{nh} \sum_{i=1}^n K \left(\frac{x - X_i}{h} \right) ; \quad (1.5)$$

ou K est la fonction noyau telle que $K(t) \geq 0$ et $\int_R K(t)dt = 1$ et $h > 0$ est le paramètre de lissage. Dans l'expression de l'estimateur à noyau continu (1,7), la fonction noyau K est

une densité de probabilité sur $R \rightarrow R^+$ est symétrique par rapport à zéro :

$$K(-u) = K(u) ;$$

ce qui implique l'égalité suivante :

$$\int_R uK(u)du = 0 ;$$

Construction de l'estimateur à noyau pour la régression

Nous supposons que nous avons des paires d'observation $(x_1, y_1) \dots (x_n, y_n)$ du couple (X, Y) . Nous proposons de construire un estimateur $\hat{m}(x)$ de la fonction de régression à partir des observations. L'estimateur à noyau est défini par une fonction $K(\cdot)$ et une fenêtre h similaire à l'estimation par noyau de la fonction de densité.

La suite $\{h_n : n \geq 1\}$ de nombre positif : $h_n \rightarrow 0$ et $n \rightarrow +\infty$ La fonction $K : R \rightarrow R$ mesurable et satisfait les hypothèse suivante :

1. K est bornée ;
2. $\lim_{|u| \rightarrow \infty} |u|K(u) = 0 ;$
3. $K(\cdot) \in L_1(R)$, i.e, $\int_R |K(u)|du < \infty ;$
4. $\int_R K(u)du = 1 ;$

Nous reprenons la fonction de régression

$$f(x) = E(Y|X = x) = \frac{r(x)}{m_X(x)} ;$$

Nous avons,

$$\begin{aligned} r(x) &= \int_R ym_{X,Y}(x, y)dy ; \\ &= \lim_{h \rightarrow 0} \frac{1}{2h} \int_{x-h}^{x+h} \int_R yM_{X,Y}(du, dy) ; \\ &= \lim_{h \rightarrow 0} \frac{1}{2h} E[Y1(|X_i - x| \leq h)] ; \end{aligned}$$

où $M_{X,Y}(\cdot, \cdot)$ est la fonction de répartition du (X, Y)

$$\hat{r}(x) = \frac{1}{n} \sum_{i=1}^n y_i \frac{1(|x_i - x| \leq h)}{2h} ;$$

Donc $f(x)$ est estimée par

$$\hat{f}(x) = \frac{\hat{r}(x)}{\hat{m}_X(x)} = \frac{\sum_{i=1}^n y_i 1(|x_i - x| \leq h)}{\sum_{i=1}^n 1(|x_i - x| \leq h)} ;$$

cet estimateur se présente sous forme la forme d'une moyenne locale pondérée des valeurs. Mais il présente de manière discontinue sa généralisation naturelle est définie comme suit :

$$\hat{f}(x) = \frac{\sum_{i=1}^n y_i K\left(\frac{x_i - x}{h}\right)}{\sum_{i=1}^n K\left(\frac{x_i - x}{h}\right)} ;$$

cet estimateur a été introduit par [Nadaray et watson][25]

$$\hat{f}(x) = \frac{\sum_{i=1}^n y_i K\left(\frac{x_i - x}{h}\right)}{\sum_{i=1}^n K\left(\frac{x_i - x}{h}\right)} \times 1\left\{\sum_{i=1}^n K\left(\frac{x_i - x}{h}\right) \neq 0\right\} ;$$

ou $1\{\cdot\} = 1_{\{\cdot\}}$ désigne la fonction indicatrice

$$\hat{f}(x) = \begin{cases} \frac{\sum_{i=1}^n y_i K\left(\frac{x_i - x}{h}\right)}{\sum_{i=1}^n K\left(\frac{x_i - x}{h}\right)}, & \text{si } \left\{\sum_{i=1}^n K\left(\frac{x_i - x}{h}\right) \neq 0\right\} \\ \frac{1}{n} \sum_{i=1}^n y_i, & \text{sinon} \end{cases}$$

1. Le biais de l'estimateur

$$Biais(\hat{f}(x)) = \frac{h^2}{2} f''(x) \int_{-\infty}^{+\infty} t^2 K(t) dt + o(h^2) ;$$

2. La variance

La variance de $\hat{f}(x)$ est :

$$Var(\hat{f}(x)) = \frac{1}{nh} f(x) \int_{-\infty}^{+\infty} K^2(t) dt + o\left(\frac{1}{nh}\right) ;$$

Exemples de noyaux

Voici quelques exemples de noyaux les plus communément utilisés :

1. Rectangulaire (Uniforme) :

$$k_1(u) = \begin{cases} \frac{1}{2}, & |u| \leq 1 \\ 0, & |u| > 1 \end{cases}$$

2.Triangulaire :

$$k_2(u) = \begin{cases} (1 - |u|), & |u| \leq 1 \\ 0, & |u| > 1 \end{cases}$$

3.Epanechnikov :

$$k_3(u) = \begin{cases} \frac{3}{4}(1 - u^2), & |u| \leq 1 \\ 0, & |u| > 1 \end{cases}$$

4.Biweight :

$$k_4(u) = \begin{cases} 1/\sqrt{2\pi} * \exp (u^2/2), & |u| \leq 1 \\ 0, & |u| > 1 \end{cases}$$

5.Gaussien :

$$k_5(u) = \left\{ (1/\sqrt{2\pi}) \exp (-u^2/2) \right. ; \forall u \in R$$

6. Cubique :

$$k_4(u) = \begin{cases} (3/4) * (1 - u^2)^3, & |u| \leq 1 \\ 0, & |u| > 1 \end{cases}$$

Les courbes de certains noyaux sont présentées ci-dessous :

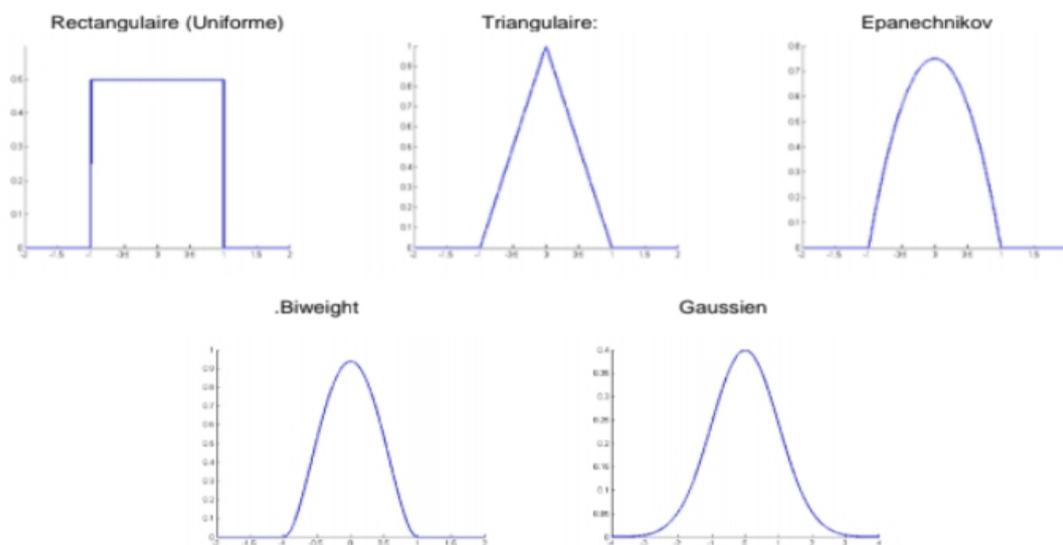


FIGURE 1.1 – Les courbes de certains noyau

1.4.3 Méthode des splines

La méthode des fonctions splines est une méthode d'analyse numérique, elle est utilisée dans le but de l'interpolation ; cette méthode est introduite pour la première fois par [Whittaker (1923)][49] et développée par [Schoenberg (1964)][34] pour servir au calcul scientifique. La méthode des spline est une technique de régression non paramétrique pouvant être vue comme une extension de régression linéaire qui modélisent automatiquement des interactions et la non-linéarités, cette méthode est énormément importante dans des différentes branches des mathématiques, est très souvent préférée à l'interpolation polynomiale, les spline sont également utilisées dans les problème de lissage des données expérimentales ou de statistique et pour représenter numériquement des contours complexes.

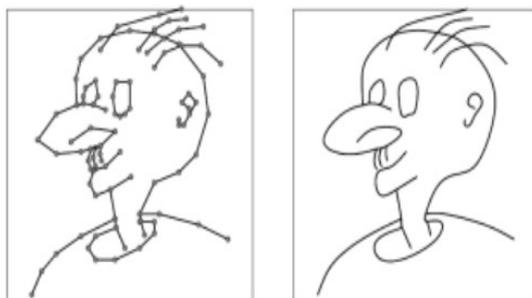


FIGURE 1.2 – Un dessin en zig-zag (à gauche) et en splines (à droite)

Une spline est une fonction définie par des polynômes par morceaux. En statistique, on l'utilise aussi pour lisser un nuage de points. Le principe des splines est de diviser l'intervalle $[a; b]$ où la fonction de spline est définie, en plusieurs sous intervalles $[a, x_1]; [x_1, x_2]; \dots; [x_n, b]$. Les points $x_1; \dots; x_n$ sont appelés les noeuds.

1- Splines de lissage

Les splines de lissage sont une façon d'utiliser les fonctions splines pour estimer la fonction de régression ; cette méthode permet de déterminer la valeur de l'estimateur en minimisant un critère bien précis.

Voir [Reinch(1967)][29], [Silverman (1985)][37], [Wahba(1990)][45] et [eubank(1999)][12]. Celui-ci combine la mesure classique de la qualité de l'ajustement, la somme des résidus au carré, et une mesure de la quantité de lissage. Les splines de lissage est une classe de fonctions non-paramétriques définie par [Whittaker (1923)][49].

Un "spline" $f : [a, b] \rightarrow R$ est une fonction polynomiale par morceaux, définie par :

$$f(x) = P_i(x), x_{i-1} \leq x \leq x_i \quad (1.6)$$

Avec $i = 1, \dots, n$ et $a = x_0 < x_1 < \dots < x_{n-2} < x_{n-1} = b$ et $P_i(x)$ est un polynôme. La

fonction f minimise la somme des carrés résidus S donnée par :

$$S(f) = \sum_{i=1}^n (y_i - f(x_i))^2 + \alpha \int_b^a (f^{(r)}(x))^2 dx ; \quad (1.7)$$

c-à-d :

$$\hat{f}(x) = \arg \min \sum_{i=1}^n (y_i - f(x_i))^2 + \alpha \int_b^a (f^{(r)}(x))^2 dx ; \quad (1.8)$$

f est une fonction différentiable telle que $f \in W_2^r[a, b]$ et α un paramètre de lissage positif.

Le paramètre α permet de réaliser un compromis entre le suivi des données exprimé par $\sum_{i=1}^n (y_i - f(x_i))^2$ et les fluctuations locales exprimé par $\int_b^a (f^{(r)}(x))^2 dx$.

Définition de l'espace de Sobolev

L'espace de sobolev $W_2^r[a, b]$ est l'ensemble des fonctions f définies sur $[a, b]$ telles que :

$$f^{(k)} \text{ est absolument continue, pour } k = 0, \dots, r - 1,$$

et

$$f^{(r)} \text{ est de carré intégrable.}$$

$W_2^r = \{f : [a, b] \rightarrow \mathbb{R} \mid f^{(k)} \text{ absolument continus pour } k = 0, \dots, r - 1, \text{ et } \int_a^b (f^{(r)}(x))^2 < \infty\}$.
i.e :

On dit que f est une fonction absolument continue, si il existe un réel a et une fonction g tels que :

$$f(x) = \int_a^b g(t) dt ;$$

2- Estimateur spline de lissage [Cao][4]

Soit \hat{f} l'estimateur spline, alors \hat{f} est défini comme suit :

$$\hat{f}(x_i) = (I + \alpha K)^{-1} Y = A_\alpha Y ; \quad (1.9)$$

avec K une matrice et A_α une matrice lisse qui ne dépend que des x_i , $i = 1, \dots, n$, l'estimateur est linéaire par rapport au vecteur Y .

3-Propriétés de l'estimateur splines de lissage

3.1 Biais de l'estimateur

$$\text{Biais}(\hat{f}, f) = E\hat{f} - f = (A_\alpha - I)f ; \quad (1.10)$$

En effet,

$$\begin{aligned} E(\hat{f}) - f &= EA_\alpha Y - f, \\ &= EA_\alpha^i Y - f, i = 1, \dots, n, \\ &= A_\alpha^i EY - f, i = 1, \dots, n, \\ &= A_\alpha^i f - f, i = 1, \dots, n, \\ &= A_\alpha f - f = (A_\alpha - I)f. \end{aligned}$$

3.2 La variance de l'estimateur

$$\text{Var}(\hat{f}) = E\|\hat{f} - E\hat{f}\|^2 = \sigma^2 \text{tr}(A_\alpha^t A_\alpha) ; \quad (1.11)$$

En effet ;

$$\begin{aligned} \text{Var}(\hat{f}) &= E\|\hat{f} - E\hat{f}\|^2, \\ &= E\|A_\alpha Y - A_\alpha f\|^2, \\ &= E \sum_{i=1}^n (A_{ij}^2)(y_i - f_{xi})^2, j = 1, \dots, n, \\ &= \sigma^2 \text{tr}(A_\alpha^t A_\alpha). \end{aligned}$$

3.3 Trace d'une matrice

La trace d'une matrice $A = (a_{ij}) 1 \leq i, j \leq n$ notée par $\text{tr}(A)$ est donnée par :

$$\text{tr}(A) = \sum_{i=1}^n (a_{ii}) ;$$

1.4.4 Méthode des séries orthogonale

Les séries orthogonale sont des utiles pour la classifications des surfaces quadratiques ; elles ont principalement appliquées dans plusieurs domaines différents (traitement du signal, statistiques fonctionnelles, . . .), le problème de l'estimation du mode de densité dans des conditions non paramétriques basées sur les séries orthogonales, à été développé par [Cencov (1962)][6] et proposé par [Kronmal et tarter(1968)][20], pour estimer des densité continues et étudiée par plusieurs auteurs [Bosq (2005)][3], [Saadi and Adjabi(2009)][32] , [Schwartz(1967)][36] et [Wahba (1981)][44] ;

Principe de la méthode

Soit $X_1 \dots X_n$ une suite des variables aléatoires indépendantes. Si on veut estimer f à partir des observations $\{X_i\}_{i=0}^n$, pour cela supposons que nous avons une fonction de régression $f \in L_2(\mathbb{R})$ puisse être représentée comme une série de Fourier.

$$f(x) = \sum_{k=0}^{\infty} a_k e_k(x) ; \quad x \in \mathbb{R} \quad (1.12)$$

où $\{e_k\}_{k=0}^{\infty}$ une base orthonormale de \mathbb{R} , et $\{a_k\}$ sont des coefficients de Fourier inconnus [Szegö][38].

Supposons que le système des fonctions $\{e_k\}, k = 1 \dots n$ forment un système orthonormal dans $[-1, 1]$ et que la variable x est limitée dans cet intervalle tel que (on choisit l'intervalle pour la simplicité de représentation) :

$$\int_{-1}^1 e_j(x) e_k(x) dx = \sigma_{jk} = \begin{cases} 0, & j \neq k \\ 1, & j = k \end{cases} \quad (1.13)$$

Alors maintenant on peut calculer les coefficients de Fourier comme suite :

$$\begin{aligned} a_j &= \sum_{k=0}^{\infty} a_k \sigma_{jk}, \\ &= a_j \int_{-1}^1 e_j(x) e_j(x) dx, \\ &= \int_{-1}^1 f(x) e_j(x) dx. \end{aligned}$$

Cette dernière intégrale ne comporte pas uniquement la base connue de fonctions, mais aussi la fonction inconnue $f(x)$. Nous donnons automatiquement une estimation de a_j si nous pouvons l'estimer de manière raisonnable. La formule du coefficient de Fourier peut s'écrire comme suite :

$$a_j = \sum_{i=1}^n \int_{-1}^1 f(x) e_j(x) dx \approx \sum_{i=1}^n f(x_i) \int_{-1}^1 e_j(x) dx. \quad (1.14)$$

Pour estimer $f(x_i)$, nous proposons de construire un estimateur sans biais de $\hat{f}(x)$. Par la méthode des moments, on peut estimer le coefficient $\{a_k; k \in N\}$ par :

$$\hat{a}_j = \sum_{i=1}^n y_i \int_{-1}^1 e_j(x) dx. \quad (1.15)$$

Si $N(n)$ termes dans la représentation (1.23) sont considérés la fonction de régression est approximée par :

$$\hat{f}_N(x) = \sum_{k=0}^N \hat{a}_k e_k(x). \quad (1.16)$$

Cet estimateur est appelé estimateur de f par les séries orthogonales. Il est une moyenne pondérée des variables Y avec le poids

$$\omega_{Ni} = n \sum_{j=1}^{N(n)} \int_{-1}^1 e_j(u) \text{ du } e_j(x).$$

proposition (1.2) [Härdle][17]

Si pour un s tel que $0 < s < 1$

$$n^{s-1} \sum_{k=0}^{N(n)} \sup_x |a_j(x)|^2 < \infty, \quad (1.17)$$

et

$$E|\epsilon_i|^{\frac{s+1}{s}} < \infty,$$

dés lors

$$N(n) \rightarrow \infty,$$

$$\hat{f}_N(x) \rightarrow f(x). \text{ en probabilité}$$

Une preuve détaillée de la consistance de \hat{m}_N peut être trouver dans [Rutkowski][31], [Szegö][38], montre que :

$$\sup_x |a_j(x)| \sim j^\rho, j = 1, 2, 3, \dots$$

avec $\rho = \frac{-1}{4}$ pour les systèmes d'Hermite et Laguerre et $\rho = 0, \frac{1}{2}$ pour les systèmes de Fourier et Legendre respectivement. L'hypothèse (1.19) apporte alors la condition sur $N(n)$ de la forme de croissance.

$$\frac{N(n)^{2\rho+1}}{n^{1-s}} \leq c < \infty \text{ quand } n \rightarrow \infty; \quad (1.20)$$

Il faut que le paramètre de lissage tend vers ∞ pour assurer la consistance, mais pas trop rapidement comme le suggère (1.20). Le taux de convergence de l'estimateur par les séries orthogonales est donné par [Härdle][15].

Conclusion

Dans ce chapitre, nous avons rappelé les notions de l'estimation non paramétrique. Ainsi nous avons donné quelques méthodes d'estimation de la régression non paramétrique, à savoir la méthode des k -plus proches voisins, la méthode du noyau, la méthode des splines, et la méthode des séries orthogonales. Enfin nous avons présenté les propriétés statistiques, biais, variance. Nous allons présenter dans le chapitre suivant en détail les deux méthodes noyau et spline et leurs principes.

2

Estimation non paramétrique par la méthode du noyau et par la méthode des splines de lissages

Introduction

Dans ce chapitre, nous allons présenter en détail l'estimateur à noyau et spline ainsi que leurs propriétés. Nous parlons aussi de son intérêt particulier sur le choix du paramètre de lissage.

2.1 Estimation non paramétrique par la méthode noyau

On suppose que lon a observé un échantillon $(X_i, Y_i); i = 1...n$ et on veut expliquer la variable aléatoire Y_i par X_i et nous considérons le modèle de régression non paramétrique donné pour $i = 1...n$

$$y_i = f(x_i) + \epsilon_i ; \quad (2.1)$$

où ϵ_i est une erreur aléatoire centré et indépendant de X_i et f est une application mesurable réelle.

La fonction de régression $f(.) = E[Y/X]$, fournit des informations sur les dépendances inconnues de Y et X ; un problème important est d'estimer f à partir de n observations. Supposons que la densité de $(X; Y)$ est $m : (x; y) \rightarrow m(x; y)$ sur R^2 ;

$m_X(x) \rightarrow m_X(x) = \int m(x; y) dy > 0$ (densité de X),

$$\forall x \in R, f(x) = E[Y/X = x] = \frac{\int y m(x, y) dy}{m_X(x)} ;$$

Les densités m et m_X sont inconnues mais on peut les estimer par :

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

$$\hat{m}_X(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{X_i - x}{h_n}\right) ;$$

par conséquent l'estimateur de la régression est donné par :

$$\forall x \in R, f(x) = E[Y/X = x] = \frac{\int y m(x, y) dy}{m_X(x)} I_{(m_X(x) \neq 0)} ; \quad (2.2)$$

Définition

Si K est un noyau d'ordre 1, l'estimateur défini par (2.2) vérifié :

$$\begin{aligned} \forall x \in R, \hat{f}(x) &= \frac{\sum_{i=1}^n Y_i K\left(\frac{X_i - x}{h_n}\right)}{K\left(\frac{X_i - x}{h_n}\right)} ; \\ &= \frac{Y_i K_{hn}(X_i - x)}{\sum_{i=1}^n K_{hn}(X_i - x)} ; \end{aligned}$$

Où $K_{hn}(\cdot) = K(\cdot/h_n)$; donc l'estimateur noyau de régression est donné

$$\hat{f}(x) = \frac{\sum_{i=1}^n Y_i K\left(\frac{X_i - x}{h_n}\right)}{K\left(\frac{X_i - x}{h_n}\right)} ; \quad (2.3)$$

$$= \frac{\hat{g}_X(x)}{\hat{m}_X(x)} ; \quad (2.4)$$

C'est l'estimateur à noyau introduit par [Nadaraya(1964)][25] et [Watson(1964)][48]

2.2 Propriétés asymptotiques de l'estimateur à noyau

Nous étudions dans cette partie deux modes de convergence, la convergence en moyenne quadratique et la convergence presque complète. nous supposons que K est un noyau vérifiant les conditions suivantes : [khadraoui][19]

1. K est bornée ,cest dire $\sup_{x \in R} |K(x)| \partial x < \infty$.

2. $\lim |x|K(x) \rightarrow 0$,quand $|x| \rightarrow +\infty$.

3. $\kappa \in L_1(R)$, c'est à dire $\int_R |K(x)| \partial x < +\infty$.

4. $\int_R |K(x)| dx = 1$.

5. $\int_R uK(x) \partial u = 0$.

6. $\int_R u^2 K(u) \partial u < +\infty$.

7. K est bornée, intégrable et a support compact.

1-Etude asymptotique du biais

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

Proposition (2.1)

i) Si $|Y| \leq C1 < \infty$ et $nh_n \rightarrow \infty$, quand $n \rightarrow \infty$, alors : $E[\hat{f}(x)] = \frac{E[\hat{g}_X(x)]}{E[\hat{m}_X(X)]} + O(\frac{1}{nh_n})$.

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

$$E[\hat{f}(x)] = \frac{E[\hat{g}_X(x)]}{E[\hat{m}_X(X)]} + O(\frac{1}{\sqrt{nh_n}}).$$

Proposition (2.2)

Si condition (4),(5) et (6) sont vérifiées si $m_X(\cdot)$ et $f(\cdot)$ sont le classe $C^2(R)$ et si $|Y|$ est borné. Alors :

$$E[\hat{f}(x) - f(x)] = \frac{h_n^2}{2} \{ \{ f''(x) + 2f'(x) \frac{m'_X(x)}{m_X(x)} \} \int_R u^2 K(u) \partial u \} (1 + O(1)).$$

Preuve

$$E[\hat{f}(x) - f(x)] = [EK(\frac{x-X}{h_n})]^{-1} \{ \int_R \frac{1}{h_n} K(\frac{x-t}{h_n}) g(t) \partial t - f(x) \int_R \frac{1}{h_n} K(\frac{x-t}{h_n}) m(t) \partial t \},$$

$$= \{ (m(x))^{-1} \} \frac{h_n^2}{2} g''(x) - \frac{h_n^2}{2} m(x) f''(x) \} \int_R u^2 K(u) \partial u \} (1 + O(1)).$$

Comme $g(x) = f(x)m(x)$, l'équation précédente peut secrire :

$$E[\hat{f} - f(x)] = \{ \frac{h_n^2}{2} \{ f''(x) + 2f'(x) \frac{m'_X(x)}{m_X(x)} \} \int_R u^2 K(u) \partial u \} (1 + O(1)).$$

D'où

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

2-Etude asymptotique de la variance

Proposition (2.3)

Sous $E(Y^2) < \infty$, alors en chaque point de continuité des fonctions $m(x)$, $f_X(x)$ et $\sigma^2(x) = \text{var}(Y|X=x)$,

on a :

$$\text{var}[\hat{f}(x)] = \frac{1}{nh_n} \left\{ \frac{\sigma^2(x)}{m_X(x)} \int K^2(u) du \right\} (O(1) + 1).$$

Preuve

soit la fonction $z(x) = \int y^2 f(x, y) \partial y$, en se basant sur le lemme de Bochner

$$\begin{aligned} \text{var}[\hat{g}, X(x)] &= \frac{1}{nh_n} \left\{ E\left[y^2 \kappa^2\left(\frac{x-X}{h_n}\right) \right] - \left[EYK\left(\frac{x-X}{h_n}\right) \right]^2 \right\}, \\ &= \frac{1}{nh_n} \left\{ \int_R K^2(u) z(x - h_n u) \partial u - h_n \left(\int_R K(u) m(x - uh_n)^2 \right) \right\}, \\ &= \frac{1}{nh_n} \left\{ \int_R K^2(u) \partial u (1 + O(1)) \right\}. \end{aligned}$$

$$E\left\{ \left\{ \hat{m}_X(x) - E(\hat{m}_X(x)) \right\} \left\{ \hat{g}_X(x) - E(\hat{g}_X(x)) \right\} \right\} = \frac{1}{nh_n} g(x) \int_R K^2(u) \partial u (1 + O(1)),$$

$$\text{var}[\hat{f}_X(x)] = \frac{1}{nh_n} f_X \int_R K^2(u) \partial u (1 + O(1)).$$

2.2.1 Choix du parametre de lissage

L'estimateur \hat{f} dépend : le noyau K et la largeur de la fenêtre h .

1-Etude de critère d'erreur quadratique moyenne de $\hat{f}(x)$

L'erreur quadratique moyenne (en anglais : mean squared error MSE) est une mesure qui permet de réduire la similarité de \hat{f} avec la fonction de régression inconnue f au point x .

Notre but est de minimiser

$$MSE(\hat{f}(x)) = E[\hat{f}(x) - f(x)]^2.$$

Le développement de cette expression faite précédemment , nous donne

$$MSE(\hat{f}(x)) = var[\hat{f}(x)] + [biais(\hat{f}(x))]^2.$$

Nous constatons d'une part que les expressions du biais de $\hat{f}(x)$ et de la variance de $\hat{f}(x)$

$$MSE(\hat{f}(x)) = \frac{h_n^4}{4} [(f''(x) + 2f'(x) \frac{m'_X(x)}{m_X(x)} (u^2 K(u)) + O(1))]^2 + \frac{1}{nh_n} (\frac{\sigma^2(x)}{m_X(x)} [K^2(u)] + (1 + O(1)))$$

où

$$[u^p K^q(u)] = \int t^p K^q(t) dt.$$

2-Fenêtre optimal

Pour trouver un compromis entre biais et variance, nous minimisons l'expression de l'erreur quadratique moyenne asymptotique AMSE (asymptotic mean square error) par rapport à h_n , donnée par :

$$MSE(\hat{f}(x)) = \left\{ \frac{h_n^4}{4} [(f''(x) + 2f'(x) \frac{m'_X(x)}{m_X(x)})]^2 [u^2 K(u)]^2 + \frac{1}{nh_n} (\frac{\sigma^2(x)}{m_X(x)} [K^2(u)]) \right\}.$$

Comme AMSE est fonction convexe. La fenêtre $h_{opt(\hat{f}(x))}^{MSE}(\hat{f}(x)) = (AMSE \hat{f}(x))$ est solution de l'équation suivante :

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

Lorsque :

$$\frac{h_n^4}{4} [(f''(x) + 2f'(x) \frac{m'_X(x)}{m_X(x)})]^2 [u^2 K(u)]^2 \neq 0;$$

d'où

$$h_{opt(\hat{f}(x))}^{MSE} = n^{-1/5} \left\{ \frac{\frac{\sigma^2(x)}{m_X(x)} [K^2(u)]}{\frac{m'_X(x)}{m_X(x)} [u^2 K(u)]^2} \right\}^{1/5}.$$

2-MISE (Mean Integrated Squared Error)

$$MISE[f_n(x)] = E\left[\int R(F_n(x) - f(x))^2 dx\right].$$

En appliquant le théorème de Fubini, on a :

$$MISE[\hat{f}(x)] = E\left[\int R(\hat{F}(x) - f(x))^2 dx\right].$$

Sous les mêmes hypothèses que les propositions (biais) et (variance), on a

$$MSE(\hat{f}(x)) = \left\{ \frac{h_n^4}{4} \left[(f''(x) + 2f'(x) \frac{m'_X(x)}{m_X(x)}) \right]^2 [u^2 K(u)]^2 + \frac{1}{nh_n} \left(\frac{\sigma^2(x)}{m_X(x)} \right) [K^2(u)] \right\}.$$

3-AMISE sous condition de continuité

Supposons que $\exists \beta > 0$ telle que $\inf_{x \in c} m(x) > 0$;

$$AMISE[\hat{f}(x)] \rightarrow 0$$

la fenêtre $hMISE_{opt(\hat{f}(x))}^{MISE}$ minimisant l'AMISE du critère global, et $h \rightarrow 0$, $nh \rightarrow \infty$ et K est borné, intégrable, positif, symétrique et a support compact On a :

$$h_{opt(\hat{f}(x))}^{MSE} = n^{-1/5} \left\{ \frac{\frac{\sigma^2(x)}{m_X(x)} [K^2(u)]}{\frac{m'_X(x)}{m_X(x)} \left\}^2 [u^2 K(u)]^2 \right\}^{1/5}.$$

Un travail similaire se fait pour le choix optimum du paramètre de lissage dans le cas de l'estimateur de [Parzen][27] et [Roseblatt][30], nous obtenons :

$$h_{MSE}^{\hat{f}(x)} = n^{-1/5} \left\{ \frac{m_X(x) [K^2]}{(m''_X(x))^2 dx [t^2 K]^2} \right\}^{1/5}$$

$$h_{MSE}^{\hat{f}(x)} = n^{-1/5} \left\{ \frac{[K^2]}{\int_R (m''(x))^2 dx [t^2 K]^2} \right\}^{1/5}.$$

Nous notons que l'expression de h_n optimal, minimisant asymptotiquement les quatre critères de erreurs la forme

$$h_{opt} = C_n^{-1/5}.$$

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

3-Validtaion croisé

La cross validation ou validation croisée est une méthode statistique qui permet d'évaluer les performances des modèles d'apprentissage automatique. Lorsqu'on entraîne un modèle sur des données étiquetées, l'hypothèse qu'il doit également fonctionner sur de nouvelles données.

La fonction validation croisée est définie par la quantité :

$$CV(h) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}^{-i}(x_i, h))^2.$$

2.2.2 Normalisation asymptotique

La première démonstration de la normalité asymptotique de l'estimateur [Nadaraya][25] et [Watson][48] est due à [Schuster][35]. On se réfère également aux théorèmes

$$\begin{aligned} f(x) &= \arg \min_a E[(Y - a)^2 | X = x]; \\ &= E(Y | X = x). \end{aligned}$$

de [Nadaraya][25] et de [Härdle] [16] qui proposent d'autres méthodes de démonstration. Le noyau K est supposé borné, à support compact et d'ordre 2. La fenêtre h est choisie égale à $cn^{-\frac{1}{5}}$.

2.3 Estimation non paramétrique par la méthode des fonctions splines

2.3.1 Fonction Spline

Soient $I = [a, b]$ un intervalle de \mathbb{R} , et soit $a < x_1 < \dots < x_n < b$ une partition de l'intervalle I et soit $f : I \rightarrow \mathbb{R}$ une fonction.

La fonction f est dit spline polynomiale de degré $(r - 1)$ (où d'ordre r) ayant pour noeuds x_1, \dots, x_n si elle vérifié les propriétés suivantes :

- f est continûment dérivable jusqu'à l'ordre r .
- Dans chaque sous-intervalle la fonction f est un polynôme de degré $\leq r, r \in \mathbb{N}$.

2.3.2 Fonction Spline Naturelle

soit $a < x_1 < \dots < x_n < b$ une partition de l'intervalle I et soit f une fonction spline de degré $\leq r$.

La fonction f est dite naturelle, si elle coïncide avec un polynôme de degré inférieur ou égale à $(r - 1)$ en dehors de l'intervalle $[x_1, x_n]$.

On note par $S_{2r}(x_1, \dots, x_n)$ l'espace des fonctions splines naturelles d'ordre $2r$, telque $\dim S_{2r}(x_1, \dots, x_n) = n$

Alors nous avons

- f est une spline de degré $2r - 1$,
- $f \in P_{r-1}$, pour $x \in [x_1, a] \cup [x_n, b]$. P_{r-1} où $P_{r-1} = p : \mathbb{R} \rightarrow \mathbb{R}$, P polynôme de degré $\leq (r - 1)$, $r \in \mathbb{N}$.

Théorème (2.1) : [Thomas-Agnan][40]

$S_r(x_1, \dots, x_n)$ est un sous espace vectoriel de l'espace des fonctions dérivables jusqu'à l'ordre $r - 2$ dont une base est donnée par les fonctions $1, x, \dots, x^{r-1}$ et les fonctions $(x - x_1)_+^{r-1}, \dots, (x - x_n)_+^{r-1}$, avec

$$\dim S_r(x_1, \dots, x_n) = r + n$$

tel que

$$u_+^{r-1} = \begin{cases} u^{r-1}, & \text{si } u > 0; \\ 0, & \text{sinon.} \end{cases}$$

2.3.3 Spline Cubique Naturelle

Soient $[a, x_1], [x_1, x_2], \dots, [x_n, b]$ sous-intervalles dans $[a, b]$, et soit f une fonction définie sur I .

f est appelée une spline cubique si les conditions suivantes sont satisfaites :

- f est un polynome cubique sur chaque sous-intervalle dans $[a, b]$,
- La fonction f est deux fois continument differentiable sur $[a, b]$, et donc f et ses dérivées d'ordre 1 et 2 sont continues aux points x_i .

2.3.4 Interpolation par les fonctions splines

Rappelons que l'interpolation consiste à trouver une fonction passant exactement par les point à notre disposition.

soient x_0, \dots, x_n les points d'interpolation. On cherche un polynôme P sur chaque segment $[x_i, x_{i+1}]$ tel que :

$$P_i : [x_i, x_{i+1}[\longrightarrow R, P_i(x) = a_i(x - x_i)^3 + b_i(x - x_i)^2 + c_i(x - x_i) + d_i,$$

Soit f une fonction spline cubique peût s'écrire sous la forme polynomiale suivante :

$$f(x) = \sum_{i=1}^n P_i(x) 1_{[x_i, x_{i+1}[}(x). \tag{2.5}$$

— La 1^{er} dérivée de f est :

$$\begin{aligned} f'(x) &= \sum_{i=1}^n P'_i(x) 1_{[x_i, x_{i+1}[}(x), \\ &= \sum_{i=1}^n 3a_i(x - x_i)^2 + 2b_i(x - x_i) + c_i, \quad x_i \leq x < x_{i+1}. \end{aligned}$$

— La 2^{me} dérivée de f est :

$$\begin{aligned} f''(x) &= \sum_{i=1}^n P''_i(x) 1_{[x_i, x_{i+1}[}(x), \\ &= \sum_{i=1}^n 6a_i(x - x_i) + 2b_i, \quad x_i \leq x < x_{i+1}. \end{aligned}$$

les coefficients a_i , b_i , c_i et d_i sont déterminés à partir des conditions d'interpolation ($P_i(x_i) = f(x_i)$ pour $i = 1, \dots, n$) et les conditions de continuité pour la première et la deuxième dérivées à chaque tous noeuds et sur chaque sous-intervalle.

Notons $h_i = x_{i+1} - x_i$, $i = 1, \dots, n$;

1- La condition (1) :

l'égalité des polynômes adjacents au point (x_i, y_i) (i.e f est continue)

$$P_{i-1}(x_i) = P_i(x_i) \Leftrightarrow a_{i-1}h_{i-1}^3 + 2b_{i-1}h_{i-1}^2 + c_{i-1}h_{i-1} + d_{i-1} = d_i = y_i$$

alors

$$d_i = y_i$$

2- La condition (4) :

l'égalité des dérivées premières au point (x_i, y_i)

$$P'_{i-1}(x_i) = P'_i(x_i) \Leftrightarrow 3a_{i-1}h_{i-1}^2 + b_{i-1}h_{i-1} + c_{i-1} = c_i$$

3- La condition (3) :

l'égalité des dérivées seconde au point (x_i, y_i)

$$P''_{i-1}(x_i) = P''_i(x_i) \Leftrightarrow 6a_{i-1}h_{i-1} + 2b_{i-1} = 2b_i$$

4- La condition (4) :

les dérivées seconde et troisième doivent être égales à zéro aux bords :

$$P_1''(x_1) = P_n''(x_n) = 0 \Leftrightarrow d_1 = d_n = 0,$$

et

$$P_1'''(x_1) = P_n'''(x_n) = 0 \Leftrightarrow a_1 = a_n = 0.$$

D'après la continuité de f et ses deux dérivées on obtient les différentes relations entre les coefficients.

Dans la première condition on obtient :

$$c_{i-1} = \frac{y_i - y_{i-1}}{h_{i-1}} - a_{i-1}h_{i-1}^2 + 2b_{i-1}h_{i-1}. \quad (1)$$

La seconde condition donne :

$$a_{i-1} = \frac{b_i - b_{i-1}}{3h_{i-1}}. \quad (2)$$

Remplaçons (2) dans (1), on trouve :

$$c_{i-1} = \frac{y_i - y_{i-1}}{h_{i-1}} - \frac{1}{3}(b_i + 2b_{i-1}h_{i-1}). \quad (3)$$

On remplace (2) et (3) dans la seconde condition, on obtient :

$$\frac{1}{3}[b_{i-1}h_{i-1} + 2b_i(h_{i-1} + h_i) + b_{i+1}h_i] = y_{i-1}\frac{1}{h_{i-1}} - y_i\left(\frac{1}{h_i} + \frac{1}{h_{i-1}}\right) + y_{i+1}\frac{1}{h_i}.$$

Nous avons

$$P_i(x_i) = y_i, \text{ et } P_i''(x_i) = 2b_i.$$

Supposons que f est une fonction spline cubique naturelle (SCN) de nœuds X_1, \dots, X_n et soient $F = (F_1, \dots, F_n)^t$ et $\Gamma = (\gamma_2, \dots, \gamma_{n-1})$ telque

$$F_i = f(x_i), \gamma_i = f''(x_i), i = 1, \dots, n, \text{ et } \gamma_1 = \gamma_n = 0.$$

Théorème (2.2) : [Green et Silverman(1994)][14]

Les vecteurs F et Γ définissent une spline cubique naturelle si et seulement s'il satisfait :

$$Q^t F = R\Gamma$$

Si cette relation est vérifiée, alors :

$$\int_a^b f''(x)^2 dx = \Gamma^2 R\Gamma = F^t K F. \quad (2.6)$$

Où Q , R , et K des matrices définies comme suit :

— Q est une matrice de taille $n \times (n - 2)$ de composantes q_{ij} :

$$\begin{aligned} q_{j-1,j} &= \frac{1}{h_{j-1}} \\ q_{jj} &= -\left(\frac{1}{h_{j-1}} + \frac{1}{h_j}\right), \\ q_{j+1,j} &= \frac{1}{h_j} \\ q_{ij} &= 0, \text{ si } |i - j| \geq 2. \end{aligned}$$

— R est une matrice symétrique de taille $(n - 2) \times (n - 2)$ de composantes r_{ij} :

$$\begin{aligned} r_{ii} &= \frac{1}{3}(h_{i-1} + h_i), i = 1, \dots, n - 1, \\ r_{i,i+1} &= r_{i+1,i} = \frac{1}{6}h_i, i = 1, \dots, n - 1, \\ r_{ij} &= 0, \text{ si } |i - j| \geq 2. \end{aligned}$$

— K est définie par :

$$K = QR^{-1}Q^t.$$

Interpolation polynomiale : Polynômes de Lagrange

Soit $f : R \rightarrow R$ et x_1, \dots, x_n n réels donnés distincts.

Interpoler la fonction f par un polynôme de degré $n - 1$ aux points x_1, \dots, x_n consiste à trouver un polynôme de degré $\leq (n - 1)$ tel que :

$$p(x_i) = f(x_i), 1 \leq i \leq n.$$

Si un polynôme existe, il s'écrit avec une forme unique suivante

$$p(x_i) = \sum_{i=1}^n \alpha_i l_i(x). \tag{2.7}$$

Avec $l_i(x)$ sont des polynômes de Lagrange associé aux points (x_i, y_i) est :

$$l_i(x) = \prod_{j=1, j \neq i}^n \frac{x - x_j}{x_i - x_j},$$

$l_i(x)$ est de degré $n - 1$ et vérifié

$$l_i(x_i) = 1, l_i(x_j) = 0. \tag{2.8}$$

En utilisant (2.8) et en prenant $j = 1, \dots, n$, on obtient

$$\alpha_j = p(x_j) = f(x_j)$$

Exemple 1

Pour $n = 1$ le polynôme de Lagrange s'écrit

$$P_1(x) = \alpha_0 \frac{x - x_1}{x_0 - x_1} + \alpha_1 \frac{x - x_0}{x_1 - x_0}.$$

C'est l'équation de la droite qui passe par les points (x_0, α_0) et (x_1, α_1) .

Exemple 2

Pour $n = 2$ le polynôme de Lagrange s'écrit

$$P_2(x) = \alpha_0 \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} + \alpha_1 \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} + \alpha_2 \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)}.$$

C'est l'équation de parabole qui passe par les points (x_0, α_0) , (x_1, α_1) et (x_2, α_2) .

Théorème (2.3) [Barbillon][2]

Le polynôme

$$L(x) = \sum_{i=1}^n l_i(x) y_i,$$

est l'unique polynôme de degré $n - 1$ vérifiant

$$L(x_i) = y_i.$$

2.3.5 Existence et unicité des spline d'interpolation

Théorème (2.4) [Thomas-Agnan] [40]

Etant donné n points (x_i, y_i) d'abscisses distinctes dans l'intervalle $[a, b]$ et $n \geq r$. Il existe une fonction et une seule f de l'espace de Sobolev $W_r^2[a, b]$ tel que :

1. f satisfait les conditions d'interpolation

$$f(x_i) = y_i, i = 1, \dots, n.$$

2. f minimise la quantité $\int_a^b f^{(r)}(x)^2 dx$ dans l'ensemble des fonctions $W_r^2[a, b]$ qui satisfait les conditions d'interpolation.

De plus, cette fonction est une spline polynomiale naturelle d'ordre $(2r)$ ayant pour noeuds les points x_1, \dots, x_n .

Le théorème suivant montre que l'interpolation par spline cubique naturelle est l'unique qui minimise $\int_a^b f^{(r)}(x)^2 dx$ par rapport à toutes les fonctions dans $W_2^2[a, b]$.

Théorème (2.5) [Cao and Golubev(2006)] [5]

Supposons que $n \geq 2$ et f_λ est la spline cubique naturelle pour les valeurs y_1, \dots, y_n en points x_1, \dots, x_n , ou $a < x_1 < \dots < x_n < b$. Soit \tilde{f} fonction dans $W_2^2[a, b]$ telle que :

$$\tilde{f}(x_i) = y_i, i = 1, \dots, n.$$

Alors

$$\int_a^b \tilde{f}''(x)^2 dx \geq \int_a^b \hat{f}''(x)^2 dx. \quad (2.9)$$

On a l'égalité si et seulement si \tilde{f} et \hat{f} sont identiques.

2.3.6 Splines de lissage

Supposons que $a < x_1 < \dots < x_n < b$, et considérons le modèle non paramétrique suivant :

$$y_i = f(x_i) + \epsilon_i; \quad (2.10)$$

où

$y_i, i = 1, \dots, n$ sont des valeurs observées de la variable réponse Y ,

f est une fonction lisse inconnue dans $W_2^2[a, b]$,

$x_i, i = 1, \dots, n$ sont des valeurs observées de la variable X ,

$\epsilon_i, i = 1, \dots, n$ sont des erreurs normalement distribuées, de moyenne nulle et de variance σ^2 .

L'estimateur de f est donné par la minimisation de la somme convexe pour $q \in]0, 1[$:

$$(1 - q) \sum_{i=1}^n (y_i - f(x_i))^2 + q \int_a^b f^{(r)}(x)^2 dx, \alpha > 0. \quad (2.11)$$

On pose $\alpha = \frac{q}{1-q}$. La formule précédents devient :

$$\sum_{i=1}^n (y_i - f(x_i))^2 + \alpha \int_a^b f^{(r)}(x)^2 dx, \alpha > 0,$$

Il est équivalent d'estimer f par la fonction \hat{f} qui minimise la fonctionnelle suivante par rapport à toute fonctions f dans $W_2^r[a, b]$.

Théorème (2.6) : [Schoenberg][34]

Le minimum du problème (2.11) admet une solution unique $\hat{f}(x)$ qui est une fonction spline dans l'ensemble $S_{2r}(x_1, \dots, x_n)$. Nous allons traiter le cas où $r = 2$, déjà vue précédemment, ce cas est souvent utilisé car il donne un algorithme très simple pour la construction de la fonction spline. En plus, les spline cubique sont facilement évaluées et

donnent "en général" des résultats satisfaisants. Poser $r = 2$ revient à résoudre le problème de minimisation suivant : Trouver la fonction f qui minimise

$$R(\hat{f}) = \sum_{i=1}^n (y_i - f(x_i))^2 + \alpha \int_a^b (f''(x))^2 dx, \quad (2.12)$$

par rapport à toute les fonctions f dans $W_2^2[a, b]$.
la solution du problème (2.12) est une spline de lissage cubique naturelle aux noeuds x_i , $i = 1, \dots, n$ (voir [Reinsh][29]).

2.3.7 Existence et unicité de la spline de lissage minimisante

Théorème (2.7) : [Thomas-Agnan][40]

Etant donnés n points (x_i, y_i) , d'abscisses distinctes dans l'intervalle $[a, b]$ et un réel $\alpha > 0$, il existe une fonction et une seule \hat{f} de l'espace de Sobolev $W_2^2[a, b]$ qui minimise la quantité $R(\hat{f})$ dans l'ensemble des fonctions $W_2^2[a, b]$. De plus, cette fonction est une spline polynomiale d'ordre 4 ayant pour noeuds les points x_1, \dots, x_n .

Théorème (2.8) : [Cao (2008)][4]

Supposons que $n \geq 0$ et que le paramètre de lissage α est positif, alors \hat{f} est une spline cubique naturelle telle que :

$$F = (I + \alpha K)^{-1} Y,$$

et pour toute f dans $W_2^2[a, b]$, on a

$$R(\hat{f}) \leq R(f). \quad (2.13)$$

2.3.8 Propriétés de l'estimateur splines de lissage

1- Biais de l'estimateur

$$Biais(\hat{f}, f) = E\hat{f} - f = (A_\alpha I)f. \quad (2.14)$$

2-Variance de l'estimateur

$$Var(\hat{f}) = E\|\hat{f} - E\hat{f}\|^2 = \sigma^2 tr(A_\alpha^t A_\alpha). \quad (2.15)$$

2.3.9 Propriétés asymptotiques de l'estimateur

1- Convergence en moyenne quadratique

Théorème (2.9) : [Wahba][41]

Supposons que $\sigma^2 \neq 0$.

Pour avoir la convergence il faut que α décroît avec n . Si cette condition est vérifiée, on a alors ;

$$- \hat{f}_\alpha(x) \longrightarrow f(x), \text{ en moyenne quadratique.}$$

Dans le modèle (2.4), [Craven and Wahba][9] ont majoré l'erreur moyenne quadratique de l'estimateur \hat{f}_α aux points d'observations x_1, \dots, x_n . Avec les notations

$$f = (f(x_1), \dots, f(x_n))^t \text{ et } \hat{f}_\alpha = (\hat{f}_\alpha(x_1), \dots, \hat{f}_\alpha(x_n))^t.$$

cette erreur est appelée MASE s'écrit

$$n^{-1} E \|f - \hat{f}_\alpha\|^2.$$

Théorème (2.10) : [Craven and Wahba][9]

Dans le modèle (2.4), on suppose que

- Les points d'observations sont tel que $\int_0^{x_j} \omega(x) dx = \frac{j}{n}$ pour $1 \leq j \leq n$ où ω est une fonction strictement positive et continue sur $[0, 1]$;
- La fonction f est dans l'espace de Sobolev $W_2^r[0, 1]$ pour un $r \geq 1$ donné ;
- Les erreurs $\epsilon_1, \dots, \epsilon_n$ sont décorrelées et de variance commune $\sigma^2 > 0$.

En supposant que m n'est pas un polynôme de degré $< r$, l'estimateur spline de lissage \hat{f}_α d'ordre $2r$ est consistant si et seulement si

$$\alpha = \alpha(n) \rightarrow 0 \text{ et } \alpha n^{2r} \rightarrow \infty \text{ lorsque } n \rightarrow \infty.$$

D'autre part, la MASE de cet estimateur vérifie

$$\frac{1}{n} E \|f - \hat{f}_\alpha\|^2 = O(\alpha) + O\left(\frac{1}{\alpha^{\frac{1}{2r}} n}\right) \quad (2.16)$$

Considérons maintenant l'erreur moyenne quadratique intégrée, notée MISE. Dans le cas du modèle (2.4), sous l'hypothèse du bruit blanc, [Ragozin(1983)][28] a majoré la MISE des dérivées de la fonction par les les dérivées correspondantes de la spline de lissage \hat{f}_α .

Théorème (2.11) : [Ragozin][28]

Dans le modèle (2.4), on suppose que

- Les points d'observations sont équidistants : $x_j = j/n, 1 \leq i \leq n$;
- La fonction f est dans l'espace de Sobolev $W_2^k[0, 1]$ avec $0 < kr$;
- Les erreurs ϵ_i sont centrées, décorrelées et de variance commune σ^2

Soit \hat{f}_α l'estimateur spline de lissage d'ordre $2r$ de f . Si le paramètre de lissage $\alpha = \alpha(n)$ est choisi de sorte que $(n\alpha^{\frac{1}{2r}}) - 1 \leq c$ pour une constante c et pour tout $n \geq r$, alors la MISE de la j^{me} dérivé ($j < k$) de \hat{f}_α vérifie

$$E|f - \hat{f}_\alpha|_j^2 = E \int_0^1 (f^j(x) - \hat{f}_\alpha^j(x))^2 dx \leq P(\alpha + n^{-2r})^{\frac{k-j}{r}} |f|_k^2 + \frac{Q\sigma^2}{n^{\frac{\alpha(2j+1)}{2r}}} \quad (2.17)$$

Pour des constantes, P et Q ne dépendent que de r , k et C . En particulier si $\alpha(n) \sim n^{\frac{-2r}{2k+1}}$ lorsque $n \rightarrow \infty$, alors

$$E|f - \hat{f}_\alpha|_j^2 = O(n^{\frac{-2r}{2k+1}}).$$

2.3.10 Choix du paramètre de lissage

1-Méthode de la validation croisée(CV)

soit $(A)_{ii}$ le i^{me} élément de la diagonale associée à la matrice de lissage A_α .

Théorème(2.12) [Moulines][24]

Le score de la validation croisée vérifie :

$$CV(\alpha) = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{f}(x_i)}{1 - (A_{\alpha ii})} \right)^2, \quad (2.18)$$

α est choisi de telle manière qu'il minimise $CV(\alpha)$.

2-Méthode de la validation croisée généralisée (GCV)

Il suffit de remplacer les dénominateurs $1 - (A_{\alpha ii})$ dans la validation croisée $CV(\alpha)$ par leurs moyenne $1 - \frac{1}{n}tr(A_\alpha)$ ainsi le score de la validation croisée généralisée :

$$GCV(\alpha) = \frac{1}{n} \frac{\sum_{i=1}^n (y_i - \hat{f}(x_i))^2}{(1 - \frac{1}{n}tr(A_\alpha))^2} \quad (2.19)$$

α est choisi de telle manière qu'il minimise $GCV(\alpha)$.

$$\alpha_{GCV} = \arg \min GCV(\alpha)$$

3-Estimation du risque sans biais

Le risque associé à l'estimateur \hat{f} de f est donnée par :

$$R(f, \hat{f}) = \frac{1}{n} E\|f - A_\alpha Y\|^2 = \frac{1}{n} \|(I - A_\alpha)f\|^2 + \frac{\sigma^2}{n} tr(A_\alpha A_\alpha^t).$$

Le principe de l'estimation du risque sans biais est de minimiser le membre droit de l'équation précédente en se basant sur les observations. Notons que le risque dépend du biais $\|(I - A_\alpha)f\|^2$. Le principe c'est le remplacer $\|(I - A_\alpha)f\|^2$ par son estimateur sans biais, qui peut être calculé par

$$E\|(I - A_\alpha)Y\|^2 = \|(I - A_\alpha)f\|^2 + \sigma^2n - 2\sigma^2tr(A_\alpha) + \sigma^2tr(A_\alpha A_\alpha^t).$$

Par conséquent

$$\|(I - A_\alpha)f\|^2 = E\|(I - A_\alpha)Y\|^2 - \sigma^2n + 2\sigma^2tr(A_\alpha) - \sigma^2tr(A_\alpha A_\alpha^t).$$

Le risque devient

$$R(f, \hat{f}) = \frac{1}{n}E\|(I - A_\alpha)Y\|^2 - \sigma^2 + 2\frac{\sigma^2}{n}tr(A_\alpha)$$

la seule chose qu'on puisse faire est de prendre $\|(I - A_\alpha)Y\|^2 - \sigma^2 + 2\frac{\sigma^2}{n}tr(A_\alpha)$ comme estimateur du biais car l'espérance n'est pas estimable à partir des observations. Ainsi, le choix de α basé sur l'estimation du risque sans biais est donné par la formule suivante :

$$\hat{\alpha} = \arg \min_{\alpha > 0} \|Y - A_\alpha Y\|^2 + 2tr(A_\alpha) \quad (2.20)$$

Dans la pratique, cette méthode a un petit inconvénient, car $\hat{\alpha}$ dépend de σ^2 qui est difficilement connu en pratique, par conséquent σ^2 sera estimé

$$\hat{\sigma}^2 = \frac{\|Y - A_\alpha Y\|^2}{n}$$

On insère cet estimateur dans la formule de $\hat{\alpha}$ et on obtient

$$\hat{\alpha} = \arg \min_{\alpha > 0} \|Y - A_\alpha Y\|^2 \left[1 + \frac{2}{n}tr(A_\alpha)\right].$$

Conclusion

Dans ce chapitre nous avons présenter l'estimation non-paramétrique de la fonction de régression par la méthode de noyau et spline. Ainsi nous avons donné leurs propriétés (l'espérance, la variance). Ensuite, nous avons utilisé la méthode de validation croisée pour le choix du paramètre de lissage pour chaque méthode.

Dans le chapitre suivant nous allons comparer les deux méthodes d'estimation de la courbe de régression de la moyenne à savoir la méthode du noyau et la méthode des splines par la simulation et un cas réel.

3

Application

Introduction

Dans ce chapitre, nous allons comparer deux méthodes d'estimation de la courbe de régression $f(x)$ à travers les différentes simulations et un cas réel : Méthodes noyaux et méthodes splines. La comparaison est basée sur les critères suivants ASE (erreur moyenne quadratique)

3.1 Application sur R

Dans cette partie, nous allons utiliser le logiciel R, c'est un langage de programmation, un logiciel libre destiné aux statistiques et à la science des données soutenu par la R Foundation for Statistical Computing. Il fait partie de la liste des paquets GNU et est écrit en C.

On a utilisé le R pour calculer et représenter graphiquement les fonctions de régression et leurs estimateurs afin de pouvoir les comparer dans des situations simulées.

Le ASE est calculé par la formule suivante :

$$ASE(\hat{f}, f) = \frac{1}{n} \sum ((\hat{f}(x) - f(x))^2)$$

3.2 Etude comparatif

3.2.1 Méthode de noyau

Rappelons qu'on suppose que lon 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 :

$$y_i = f(x_i) + \epsilon_i ;$$

où ϵ_i est l'aléatoire centré et indépendante de X_i . Aussi la fonction de régression :

$$f(x) = E[Y/X = x] = \frac{\int ym(x, y)\partial y}{m_X(x)}$$

où $m_X(x)$ est la densité de la variable X , et $g_X(x) = \int ym(x, y)\partial y$.

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

$$\hat{f}(x) = \frac{\sum Y_i K\left(\frac{X_i - x}{h_n}\right)}{\sum K\left(\frac{X_i - x}{h_n}\right)} = \frac{\hat{g}_X(x)}{\hat{m}_X(x)};$$

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 $\hat{f}(x)$: avec $\hat{g}_X(x)$ est l'estimateur naturel de $g_X(x)$:

$$\hat{g}_X(x) = \frac{1}{nh_n} \sum Y_i K_{h_n}(X_i - x);$$

et $\hat{m}_X(x)$ l'estimateur du noyau de la densité de la fonction m

$$\hat{m}_X(x) = \frac{1}{nh_n} \sum K_{h_n}(X_i - x).$$

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

$$y = f(x) + \epsilon, \quad \text{où } \epsilon \rightarrow N(0; \sigma^2)$$

3.2.2 Méthode des splines

L'estimateur de la fonction de régression spline :

$$\hat{f}_\alpha(x_i) = (I + \alpha K)^{-1} Y = A_\alpha Y.$$

avec A_α une matrice lisse qui ne dépend que des $x_i, i = 1, \dots, n$, l'estimateur est linéaire par rapport au vecteur Y .

Le paramètre de lissage est choisi par la méthode de la validation croisée généralisée :

$$\hat{\alpha} = \arg \min \frac{1}{n} \frac{\sum_{i=1}^n (y_i - \hat{f}(x_i))^2}{(1 - \frac{1}{n} \text{tr}(A_\alpha))^2} ;$$

Les simulations sont effectuées pour différentes tailles d'échantillon de plus en plus grandes :

$$n \in \{50, 100, 200, 500, 1000, 2000\}$$

3.3 Modèles de fonction

Les modèles considérés sont :

1. $f_1(x) = 4.26(\exp(-3.25x_i) - 4\exp(-6.5x_i) + 3\exp(-9.75x_i)) + e_i$, [Eubank][12]
2. $f_2(x) = x_i + 0.5\exp(-50(x_i - 0.5)^2) + e_i$. [Eubank][12]

3.3.1 Modèle $f_1(x)$

1. x est uniforme sur $[0, 5]$;
2. e sont les résidus qui sont normalement distribués de moyenne 0 et de variance $\sigma^2 = 0.05$;

ASE associé à f_1

n	ASE noyau	ASE spline
50	0.1480599	0.0898976
100	0.0975050	0.0654594
200	0.0750013	0.0497245
500	0.0655905	0.0374850
1000	0.0674577	0.0216823
2000	0.0414707	0.0160180

TABLE 3.1 – ASE donnée par les deux méthodes associée au modèle $f_1(x)$ en fonction de la taille de l'échantillon n .

Pour $n = 50$

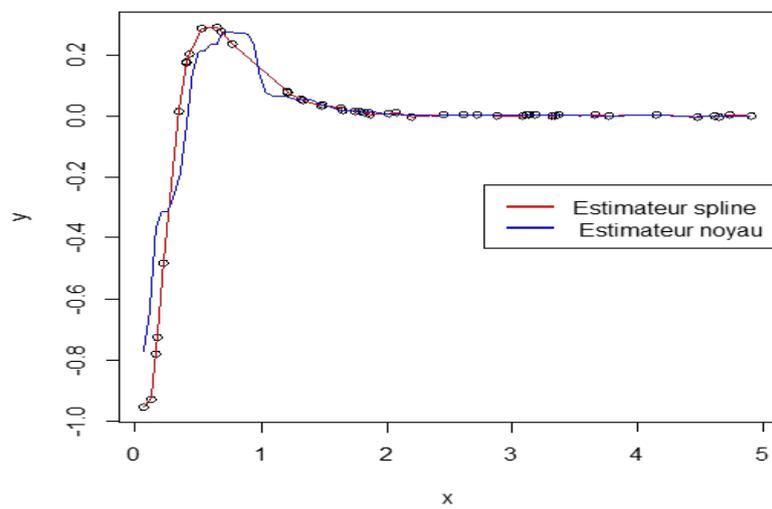


FIGURE 3.1 – Estimation de f_1 , $n = 50$.

Pour $n = 100$

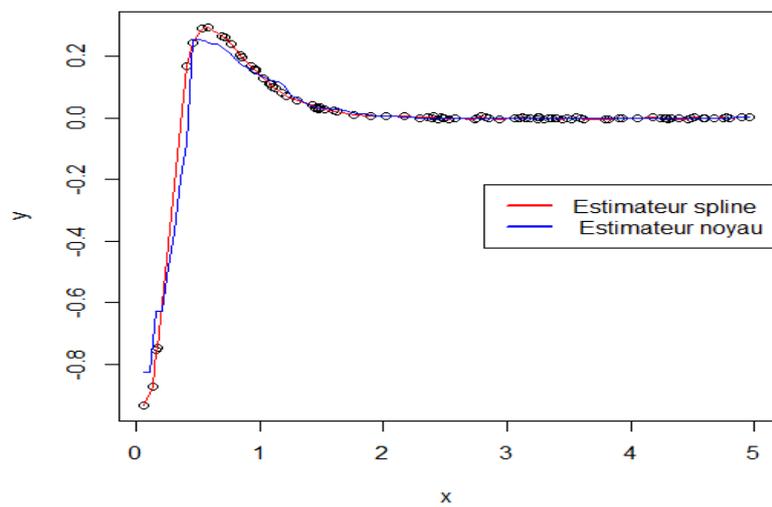


FIGURE 3.2 – Estimation de f_1 , $n = 100$.

Pour $n = 200$

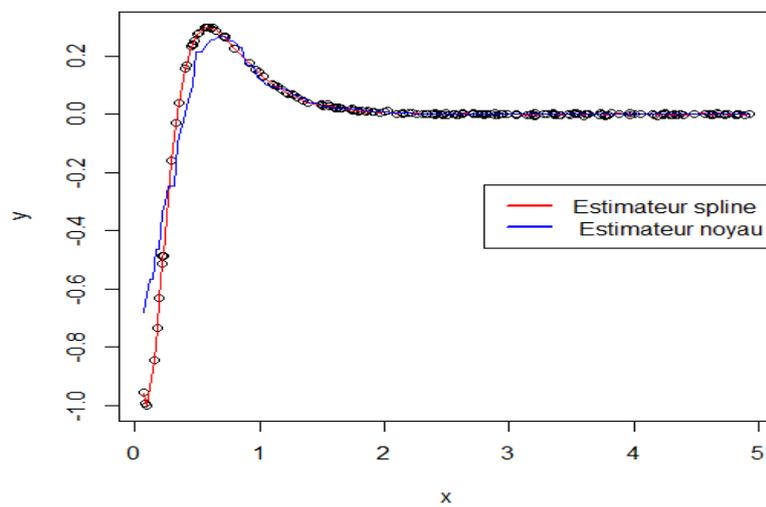


FIGURE 3.3 – Estimation de f_1 , $n = 200$.

Pour $n = 500$

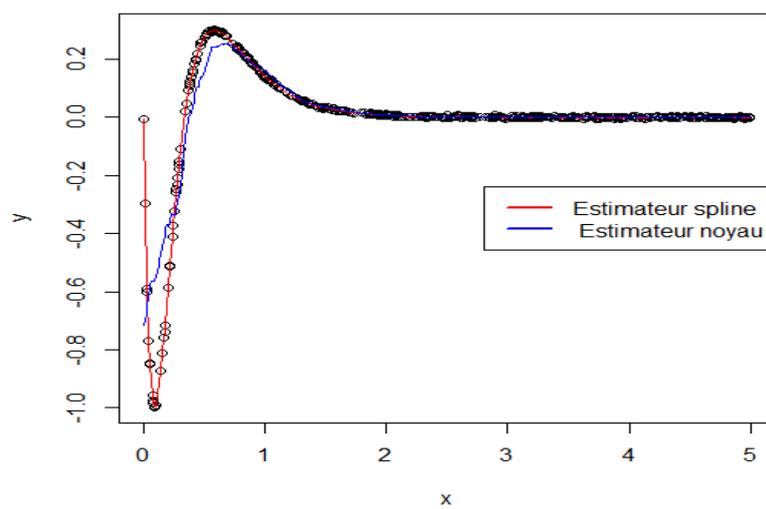


FIGURE 3.4 – Estimation de f_1 , $n = 500$.

Pour $n = 1000$

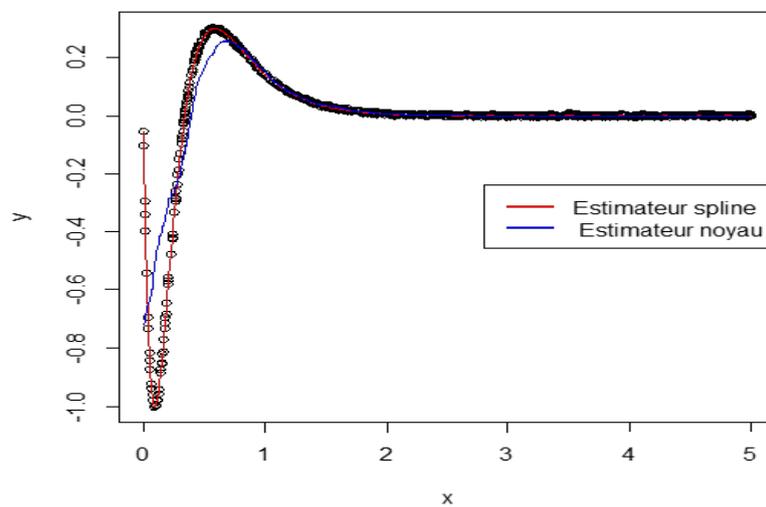


FIGURE 3.5 – Estimation de f_1 , $n = 1000$.

Pour $n = 2000$

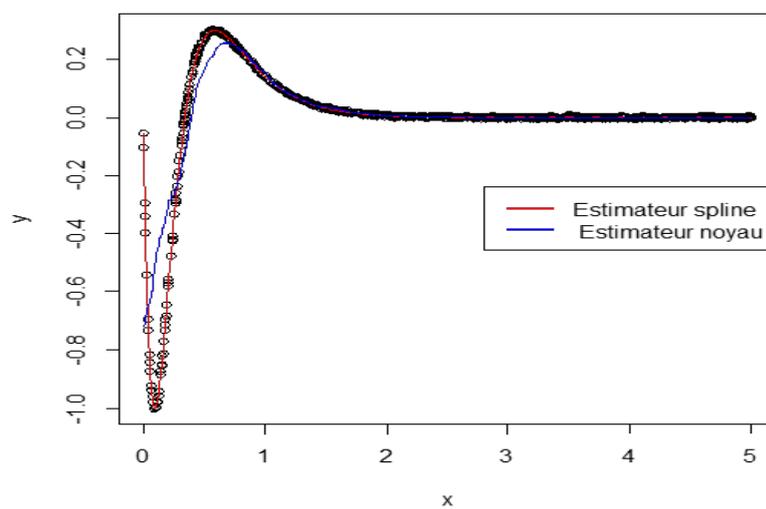


FIGURE 3.6 – Estimation de f_1 , $n = 2000$.

Discussion des résultats du modèle f_1 :

D'après les résultats du tableau(3.1) on peut observée que :

- La valeur de ASE diminue lorsque la taille d'échantillon n augmente et ce résultat est vrai pour les deus méthode.

On remarque que quelque soit la taille n de l'échantillon l'erreur associée à la méthode des fonction spline est la plus petite par rapport a la méthode du noyau (ASE spline $<$ ASE noyau),

Graphiquement :

On remarque que la meilleur méthode qui ajuste bien les données est la méthode des splines.

On conclue que la méthode des splines est meilleure que celle du noyau.

3.3.2 Modèle $f_2(x)$

1. x est uniforme sur $[0, 3]$;

2. e sont les résidus qui sont normalement distribués de moyenne 0 et de variance $\sigma^2 = 0.0225$;

ASE associé à f_2

n	ASE noyau	ASE spline
50	0.516535	0.08642
100	0.405451	0.07123
200	0.065675	0.00684
500	0.045677	0.00478
1000	0.033895	0.00474
2000	0.017639	0.00356

TABLE 3.2 – ASE donnée par les deux méthodes associée au modèle $f_2(x)$ en fonction de la taille de l'échantillon n .

Pour $n = 50$

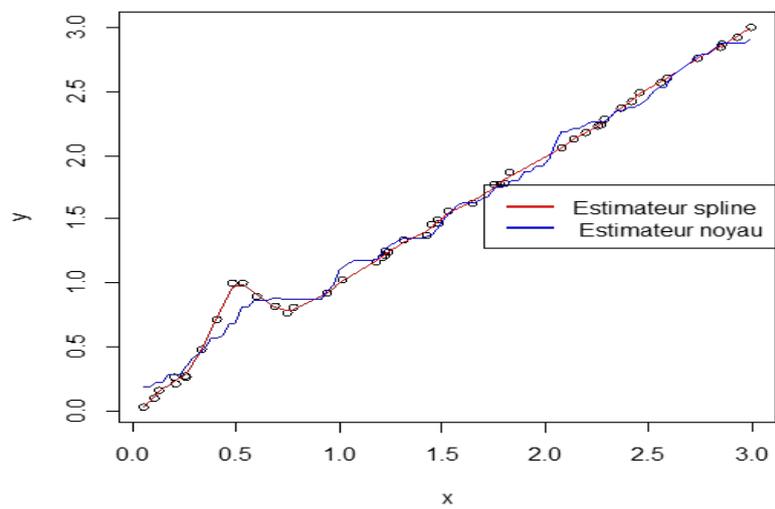


FIGURE 3.7 – Estimation de f_2 , $n = 50$.

Pour $n = 100$

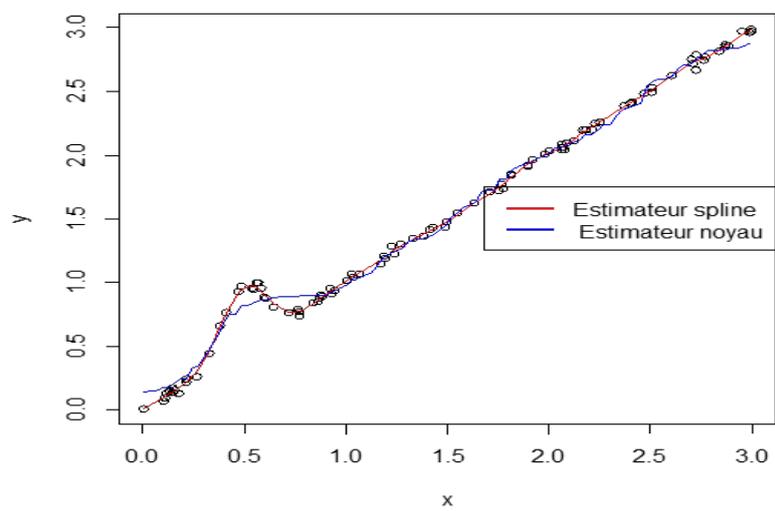


FIGURE 3.8 – Estimation de f_2 , $n = 100$.

Pour $n = 200$

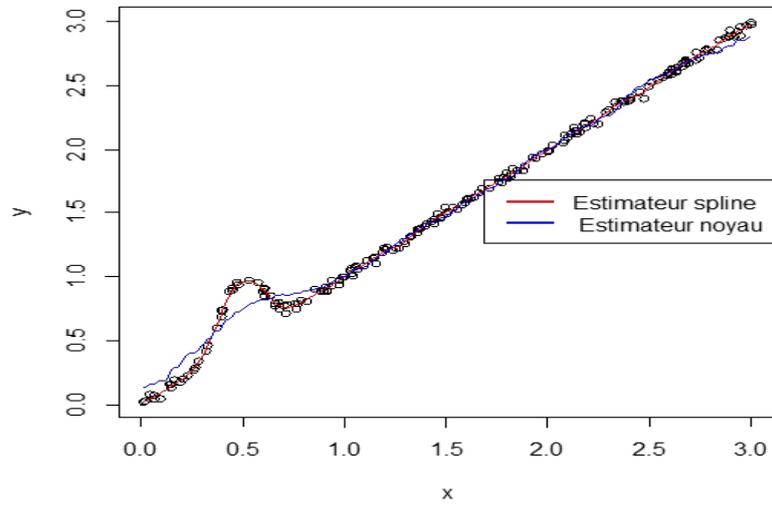


FIGURE 3.9 – Estimation de f_2 , $n = 200$.

Pour $n = 500$

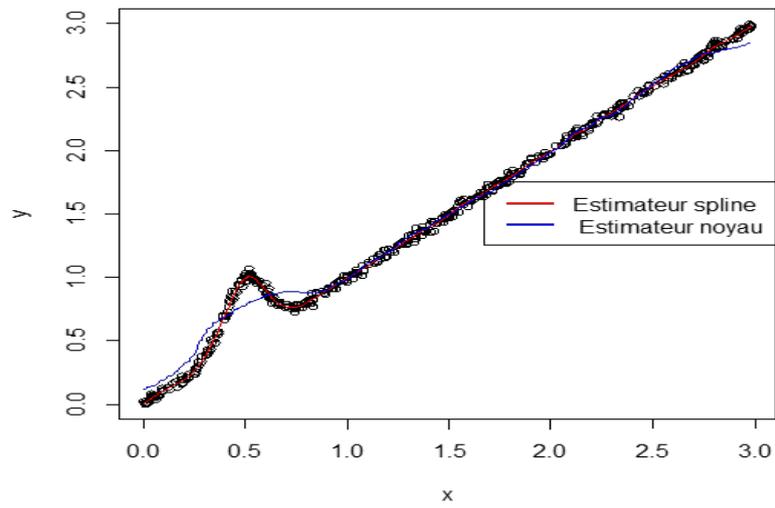


FIGURE 3.10 – Estimation de f_2 , $n = 500$

Pour $n = 1000$

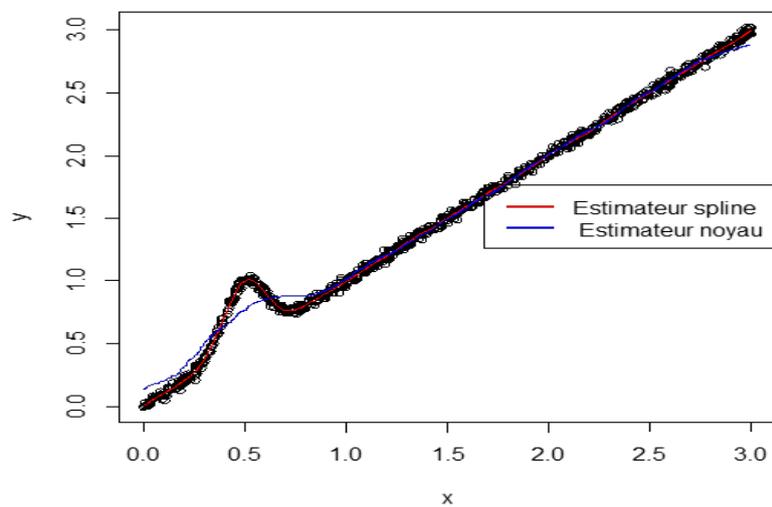


FIGURE 3.11 – Estimation de f_2 , $n = 1000$.

Pour $n = 2000$

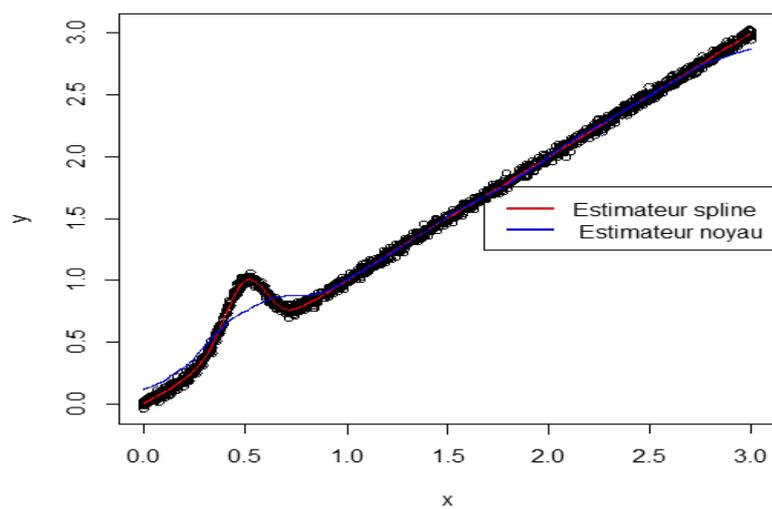


FIGURE 3.12 – Estimation de f_2 , $n = 2000$.

Discutez les résultats du modèle f_2 :

D'après les résultats du tableau(3.2) on peut observée que :
Les valeur de ASE diminue lorsque la taille d'échantillon n augmente,ce résultat est vrai pour les deux méthodes .
On trouve les mêmes résolution que le modèle f_1 .Donc la méthode optimale pour nos données est toujours la méthode spline .Et cette constatation est vérifiée graphiquement.

3.3.3 Cas réel

Soit l'ensemble données réelles concerne les données de l'éolienne, voir [Ebank(1999)][\[12\]](#) :

- La variable dépendante Y est la sortie du courant continu.
- La variable indépendante X est la vitesse du vent (miles/h).

x	y	x	y	x	y
2.45	0.123	2.70	0.500	2.90	0.653
3.05	0.558	3.40	1.057	3.60	1.137
3.95	1.144	4.10	1.194	4.60	1.562
5.00	1.582	5.45	1.501	5.80	1.737
6.00	1.822	6.20	1.866	6.35	1.930
7.00	1.800	7.40	2.088	7.85	2.179
8.15	2.166	8.80	2.112	9.10	2.303
9.55	2.294	9.40	2.386	10.00	2.236
10.20	2.310				

TABLE 3.3 – Données sur les éoliennes.

ASE associée au cas réel

ASE noyau	ASE spline
0.751359	0.386751

TABLE 3.4 – ASE associée au cas réel.

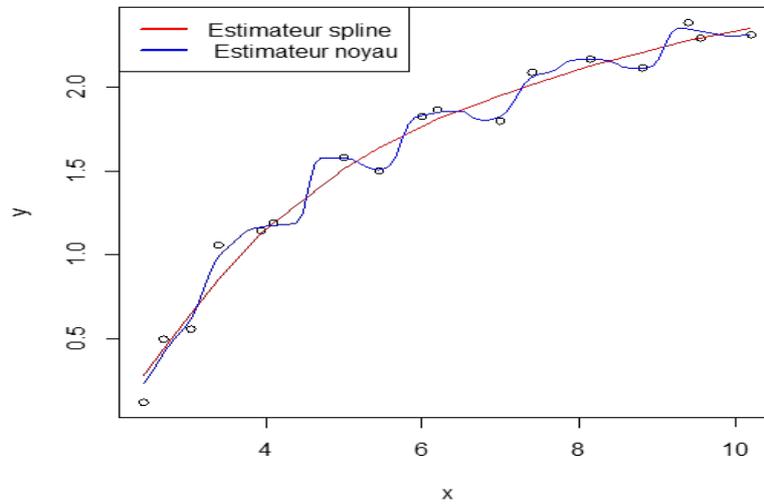


FIGURE 3.13 – Estimation de la courbe d'éoliennes par la méthode du noyau et la méthode des fonctions splines.

Discussion des résultats :

Pour les résultats du tableau (3.3) en conclut que l'erreur associée à la méthode des fonctions splines est petite par rapport à celle donnée par la méthode du noyau. Ce qui signifie que la méthode des spline estime bien les données même dans le cas réel.

Conclusion

Dans ce dernier chapitre, nous avons comparé la méthode du noyau et la méthode des splines pour estimer la fonction de régression. Les résultats obtenus montrent que la méthode des splines donne de meilleurs résultats que celle du noyau et ces résultats sont confirmés graphiquement.

Conclusion générale

Dans ce travail nous avons présenté deux méthodes d'estimation non paramétrique de la fonction de régression, les méthodes étudiées sont la méthode du noyau et la méthode des fonctions des splines . Le but principal est de comparer ces deux méthodes en s'appuyant sur des exemples de fonctions de régressions simulé et sur un cas réel. Ce travail peut être résumer en deux parties principales :

La première partie est théorique, elle comporte les deux premiers chapitres. Dans le premier chapitre nous avons défini quelques méthodes d'estimation non paramétrique de la fonction de régression, la méthode de noyau, les k plus proche voisin, la méthode des séries orthogonales et la méthode des splines .

Dans le deuxième chapitre, nous nous sommes intéressé à la méthode du noyau et la méthode spline, nous avons défini les estimateurs et leurs différentes propriétés statistiques et asymptotiques.

La deuxième partie de ce travail, est la partie simulation ,l'utilisent le logiciel R qui nous a permis de comparer par simulation, sur des fonctions cibles et un cas réel les deux méthode. Cette comparaison basé sur le critère ASE .

On conclut :

. Le meilleur résultat est obtenu avec la méthode de spline.

. l'erreur associer à la méthode des spline est la plus petite à celle associer a la méthode du noyaux, ($ASE_{spline} < ASE_{noyau}$). Donc le meilleur résultat est obtenu par la méthode des spline.

Bibliographie

- [1] ADJABI, S., AMROUN, S., ET AL. *L'estimation de la courbe de régression de la moyenne*. PhD thesis, Université de Bejaia, 2011.
- [2] BARBILLON, P. Méthodes d'interpolation à noyau pour l'approximation de fonction type boîte noire coûteuses. *Université Paris Sud-Paris* (2010).
- [3] BOSQ, D. Estimation optimale de la densité par projection. *Canadian Journal of Statistics* 33 (2005), 21–37.
- [4] CAO, Y. Inégalités d'oracle pour l'estimation de la régression. *Thèse Doctorat, Université de Provence* (2008).
- [5] CAO, Y., AND GOLUBEV, Y. On oracle inequalities related to smoothing spline. *Mathematical Methods of Statistics* 15 (2006), 398–414.
- [6] CENCOV, N. Evaluation of an unknown distribution density from observations. *Journal of the Institute of Mathematics and Its Applications* 18 (1962), 1559–1962.
- [7] COLLOMB, G. Estimation non paramétrique de la régression par la méthode du noyau : Propriétés de convergence asymptotiquement normale indépendante. *Annales scientifiques de l'Université de Clermont-Ferrand 2, tome 65, série mathématiques*, 15 (1977), 24–46.
- [8] COVER, T.M. ET HART, P. Nearest neighbor pattern classification. *IEEE Transactions on Information Theory* 13 (1967), 21–27.
- [9] CRAVEN, P., AND WAHBA, G. Smoothing noisy data with spline functions. Estimating the correct degree of smoothing by the method of generalized cross-validation. *Numer. Math* 31(4).
- [10] EPANECHNIKOV, V. A. Nonparametric estimation of a multidimensional probability density. *Theory Probab. Appl* (1969).
- [11] EUBANK, R. Spline smoothing and nonparametric regression. *Dekker, New York* (1988).
- [12] EUBANK, R. Nonparametric regression and spline smoothing. (2nd ed.) *New York, Dekker* (1999).
- [13] FERRATY, F., AND VIEU, P. Statistique fonctionnelle : Modèle non paramétrique de régression. *Note de cours de DEA* (2002).
- [14] GREEN, P. J., AND SILVERMAN, B. W. Nonparametric regression and generalized linear models : A roughness penalty approach. *of Monographs on Statistics and Applied Probability. Chapman et Hall, London* 58 (1994).

- [15] HARDLE, W. Robust regression function estimation. *Journal of Multivariate Analysis* 14 (1984), 169–80.
- [16] HARDLE, W. Applied nonparametric regression. *Cambridge University Press, Cambridge* (1990).
- [17] HARDLE, W. Applied nonparametric regression. *Humboldt-Universität zu Berlin*. (1994).
- [18] HURLIN, C. Régressions non paramétriques univariées. *Master Econométrie et Statistique Appliquée (ESA), Université d'Orléans*. (2007-2008).
- [19] KHADRAOUI, R. Estimation à noyau de la fonction de régression.
- [20] KRONMAL, R., AND M.TARTER. the estimation of probability densities and cumulatives by fourier series methodes. *J.Amer. Statist. Assoc* 63 (1968), 925–952.
- [21] LAI, S. L. Large sample properties of k-nearest neighbor procedures. *Ph.d. dissertation, Dept. Mathematics, UCLA, Los Angeles*. (1977).
- [22] LOFTSGAARDEN, D. O., AND QUESENBERRY, C. P. A nonparametric estimate of a multivariate density function. *Annals of Mathematical Statistics* 36 (2005), 1049–51.
- [23] M.H, S. Spline analysis. *Prentice Hall, Englewood Cliffs, N J* (1973).
- [24] MOULINES, M. C. Classification et sélection de caractéristiques de textures.
- [25] NADARAYA, E. A. On estimating regression. *Theor. Prob. Appl.* 9 (1964), 141–142.
- [26] OLIVEIRA, J. T. D. Estatística de densidades : resultados assintoticos. *Universidade de Lisboa. Faculdade de Ciências de Lisboa*, (1963).
- [27] PARZEN, E. On estimation of a probability density function and mode. *The annals of mathematical statistics* 33, 3 (1962), 1065–1076.
- [28] RAGOZIN., D. L. Error bounds for derivative estimates based on spline smoothing of exact or noisy data. *J. Approx. Theory*, 37(4) (1983), 335–355.
- [29] REINSCH, C. Smoothing by spline function. *Numererisch Matheematik* 10 (1967), 177–183.
- [30] ROSENBLATT, M. Remarks on some nonparametric estimates of a density function. *Annals of Mathematical Statistic* 27 (1956), 832–837.
- [31] RUTKOWSKI, L. Sequential estimates of a regression function by orthogonal series with applications in discrimination, in revesz, schmetterer and zolotarev (eds). *The First Pannonian Symposium on Mathematical Statistics, Springer-Verlag, pp* (1981), 263–44.
- [32] SAADI, N., AND ADJABI, S. On the estimation of the probability density by trigonometric series. *communications in Statistics-Theory and Methods* 38 (2009), 3583–3595.
- [33] SARDA, P., AND VIEU, P. Kernal regression. smoothing and regression : Approches, comutation, and application. *Ed. M.G. Schimek, Wiley series in probability and statistics* (2000), 43–47.

- [34] SCHOENBERG, I. J. Spline functions and the problem of graduation. *Mathematics* 52 (1964), 974–50.
- [35] SCHUSTER., E. F. Joint asymptotic distribution of the estimated regression function at a finite number of points. *Annals of Mathematical Statistics* 43, 1 (1972), 84–88.
- [36] SCHWARTZ, S. C. Estimation of probability densities by an orthogonal series. *Annals of Mathematical Statistics* 38 (1967), 1261–1265.
- [37] SILVERMAN, B. W. Some aspects of the spline smoothing approach to non- parametric regression curve fitting (with discussion). *of the Royal Statistical Society Series B* 47 (1985), 1–52.
- [38] SZEGO, G. Orthogonal polynomials. *Amer. Math. Soc. Coll. Publ.* (1959).
- [39] T P. VIEU F. FERRATY. Statistique fonctionnelle : Modèle non-paraamétrique de régression. *Notes de cours de DEA.* (2002/2003).
- [40] THOMAS-AGNAN, C. Spline function and filtering. *The Annals of Statistics*, 19 : (1999), 1512–152.
- [41] WAHBA, G. Optimal convergence properties of variable knot, kernel, and orthogonal series methods for density estimation. *Annals of Statistics* 3.
- [42] WAHBA, G. Convergence properties of the method of regularization for noisy linear operator equations. *TSR No. 1132, Math. Res. Center, Univ. of Wisconsin-Madison* (1973).
- [43] WAHBA, G. Smoothing noisy data by spline functions. *II.Tech. Report No. 380, Dept. of Statist. Univ. of Wisconsin- Madison.* (1974).
- [44] WAHBA, G. Data-based optimal smoothing of orthogonal series density estimates. *The Annals of Statistics* 9(1) (1981), 146–156.
- [45] WAHBA, G. Spline models for observational data. *CBMS-NSF series. SIAM, Philadelphia* (1990).
- [46] WAHBA, G., AND WOLD., S. A completely automatic french curve : Fitting spline fuctions by cross validation. *Simulation and Computation* (1975), 1–17.
- [47] WALTER, G. Properties of hermite series estimation of probability density. *Annals of Statistics* 5 (19), 1258–64.
- [48] WATSON, G. S. Smooth regression analysis. *Sankhyà Ser. A* 26 (1964), 359–372.
- [49] WHITTAKER, E. On a new method of graduation. *Proceedings of the Edinburgh Mathematical Society* 41 (1923), 63–75.

Résumé

Dans ce travail nous étudions l'estimateur non paramétrique de la fonction de régression pour les données réelles en utilisant la méthode de noyau et spline. Puis nous allons la comparer avec la méthode de ASE (average squared error) Finalement un travail de simulation on a donné des explications graphiques des résultats théoriques appliqués sur des exemples de régression non linéaire à l'aide du logiciel R. Pour vérifier la bonne estimateurs entre noyaux et spline.

Mots clés : régression non-paramétrique ; Fonction ; Spline de lissage ; noyau ; VCG

Abstract

In this work we study the nonparametric estimator of the regression function. sion for real data using the kernel and spline method. Then we compare it with the ASE (average squared error) method Finally, a simulation work graphical explanations of the theoretical results applied to nonlinear examples of non-linear regression using R.software to check that the estimators between kernels and splin.

words :non-parametric regression ;Function ;Smoothing spline ;kernels ;VCG