

Estimation non-paramétrique
Compte-rendu d'article, 1er trimestre 2021–2022

NOM: DURAND

Prénom: Homer

Titre de l'article : *Nonparametric density and survival function estimation in the multiplicative censoring model*

Table des matières

1	Rappels et Notations	2
1.1	Rappels	2
1.2	Notations	3
2	Contexte	4
3	Propriétés théoriques	6
4	Discussion	16

1 Rappels et Notations

1.1 Rappels

Définition 1 (Noyau). *Un noyau est une fonction $K : \mathbb{R} \rightarrow \mathbb{R}$ intégrable ($K \in L^1(\mathbb{R})$), d'intégrale égale à 1 : $\int_{\mathbb{R}} K(u)du = 1$.*

Définition 2 (Noyau d'ordre l). *Soit $l \in \mathbb{N}^*$. On dit que $K : \mathbb{R} \rightarrow \mathbb{R}$ est un noyau d'ordre l si pour tout $j \in \{0, \dots, l\}$ les fonctions $u \mapsto u^j K(u)$ sont intégrables telles que*

$$\int_{\mathbb{R}} K(u)du = 1$$

$$\int_{\mathbb{R}} u^j K(u)du = 0, \quad \forall j \in \{1, \dots, l\}$$

Définition 3 (Classe de Hölder). $\Sigma_I(\beta, R) = \{f : I \rightarrow \mathbb{R}, f^{(l)}$ existe pour $l = \lfloor \beta \rfloor, |f^{(l)}(x) - f^{(l)}(y)| \leq R|x - y|^{\beta-l}, \forall x, y \in I\}$

Définition 4 (Classe de Nikol'ski). *Soit $\beta, R \in \mathbb{R}^{+*}$, on définit la classe de Nikol'ski $\mathcal{N}(\beta, R)$ comme :*

$$\mathcal{N}(\beta, R) = \left\{ f \in L_2(\mathbb{R}), \forall y \in \mathbb{R}, \left(\int |f^{\lfloor \beta \rfloor}(x+y) - f^{\lfloor \beta \rfloor}(x)|^2 \right)^{1/2} \leq R|y|^{\beta-\lfloor \beta \rfloor}, \|f^{\lfloor \beta \rfloor}\| \leq R \right\}$$

On peut voir cette classe une version intégrée de la classe de Hölder.

Définition 5. *On note $\hat{F}_Y(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{Y_i \geq x\}}$*

Définition 6 (Inégalité de Young). *Soit f une fonction de $L^p(\mathbb{R})$ et g une fonction de $L^q(\mathbb{R})$. Soit $p, q, r \in [1, +\infty]$ tels que*

$$\frac{1}{p} + \frac{1}{q} = 1 + \frac{1}{r}$$

Alors on a

$$\|f \star g\| \leq \|f\|_p \|g\|_q$$

En particulier, si $p = 1, r = q = 2$, on a

$$\|f \star g\|_2 \leq \|f\|_1 \|g\|_2 \tag{1}$$

Théorème 1 (Talagrand). *Soit X_1, \dots, X_n des variables aléatoires indépendantes à valeur dans un espace Ξ mesurable. Soit \mathcal{F} un espace dénombrable de fonctions mesurables de Ξ dans \mathbb{R} . On définit pour tout $\psi \in \mathcal{F}$,*

$$\nu_n(\psi) = \frac{1}{n} \sum_{i=1}^n (\psi(X_i) - \mathbb{E}[\psi(X_i)]).$$

En supposant qu'il existe trois constantes positives M, H et v telles que

$$\sup_{\psi \in \mathcal{F}} \|\psi\|_{\infty} \leq M, \quad \mathbb{E}[\sup_{\psi \in \mathcal{F}} |\nu_n \psi|] \leq H, \quad \sup_{\psi \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \text{Var}[\psi(X_i)] \leq v$$

alors, pour tout $\alpha > 0$,

$$\mathbb{E} \left[\left(\sup_{\psi \in \mathcal{F}} |\nu_n(\psi)|^2 - 2(1 + 2\alpha)H^2 \right)_+ \right] \leq \frac{4}{a} \left(\frac{v}{n} \exp \left(-a\alpha \frac{nH^2}{v} \right) + \frac{49}{aC^2(\alpha)} \frac{M^2}{n^2} \exp \left(-\frac{\sqrt{2}aC(\alpha)\sqrt{\alpha} nH}{7M} \right) \right)$$

avec $C(\alpha) = (\sqrt{1 + \alpha} - 1) \wedge 1$ et $a = 1/6$.

1.2 Notations

- **Norme p** : $\|f\|_p = \int |f(x)|^p dx$. On notera $\|\cdot\|$ la norme euclidienne.
- **Produit de convolution** : $f \star g(x) = \int f(x - u)g(u)du$
- Pour tout $h > 0$ et K , un noyau, on note $K_h(u) = \frac{1}{h}K\left(\frac{u}{h}\right)$

2 Contexte

L'article [2] propose d'étudier le modèle dit de censure multiplicative ("*multiplicative censoring model*") défini par

$$Y_i = X_i U_i, i = 1, \dots, n \quad (2)$$

où $X_i, i = 1, \dots, n$ sont des réels indépendants et identiquement distribués dont la densité f et la fonction de répartition F sont inconnues et où les $U_i, i = 1, \dots, n$ sont des variables aléatoires *i.i.d* suivant une loi uniforme sur $[0, 1]$. Ce modèle est dit modèle de censure multiplicative car les X sont "censurés" par une variable aléatoire U indépendante. On parle également de *shrinkage censoring model*. Les auteurs proposent donc des estimateurs des fonctions de densité f et de survie $\bar{F}(x) = 1 - F(x)$ basés sur les méthodes d'estimation à noyau largement étudiées dans [4].

L'ensemble des méthodes étudiées repose sur la relation entre les fonctions de densité f_Y et de survie $\bar{F}_Y = 1 - F_Y$ des Y_i et celles de X_i :

$$\bar{F}_Y(y) + y f_Y(y) = \bar{F}(y), \forall y \in \mathbb{R} \quad (3)$$

Cela se montre facilement en utilisant la proposition 1, en particulier le fait que

$$\forall y \in \mathbb{R}, f_Y(y) = \int_y^{+\infty} \frac{f(x)}{x} \mathbb{1}_{y \geq 0} dx + \int_{-\infty}^y \frac{f(x)}{|x|} \mathbb{1}_{y < 0} dx$$

Ainsi on a pour $y \geq 0$

$$\begin{aligned} \bar{F}_Y(y) &= \int_y^{+\infty} f_Y(z) dz = \int_y^{+\infty} \int_z^{+\infty} \frac{f(x)}{x} dx dz = \int \left(\int_x^y dz \right) \frac{f(x)}{x} \mathbb{1}_{y \leq x} dx \\ &= \int_y^{+\infty} f(x) dx - y \int_y^{+\infty} \frac{f(x)}{x} dx = \bar{F}(y) - y f_Y(y) \quad \text{par (1)} \end{aligned}$$

et pour $y < 0$

$$\begin{aligned} F_Y(y) &= \int_{-\infty}^y f_Y(z) dz = \int \left(\int_x^y dz \right) \frac{f(x)}{|x|} \mathbb{1}_{x \leq y} dx \\ &= y f_Y(y) + F(y) \end{aligned}$$

On retrouve donc bien (3). La propriété suivante découle immédiatement :

Propriété 1. Soit $t : \mathbb{R} \rightarrow \mathbb{R}$ une fonction bornée et dérivable telle que $t' \in \mathbb{L}^2(\mathbb{R})$. Supposons que $\mathbb{E}[|X|] < +\infty$, alors

$$\mathbb{E}[t(Y) + Y t'(Y)] = \mathbb{E}[t(X)] \quad (4)$$

Cette relation nous permet d'obtenir de simples estimateurs des fonctions de survie et de densité \bar{F} et f en apportant une correction aux données Y .

Définition 7 (Estimateur à noyau). Soit $K : \mathbb{R} \mapsto \mathbb{R}$ un noyau. L'estimateur de la fonction de survie $\bar{F}(x)$ est défini par :

$$\begin{aligned}\hat{F}(x) &= \frac{1}{nh} \sum_{i=1}^n \left(\int K\left(\frac{u-x}{h}\right) \mathbb{1}_{\{Y_i \geq u\}} du + Y_i K\left(\frac{Y_i-x}{h}\right) \right) \\ &= K_h \star \hat{F}_Y(x) + \frac{1}{n} \sum_{i=1}^n Y_i K_h(Y_i - x)\end{aligned}$$

Celui de la densité $f(x)$ est défini par :

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n (Y_i K_h'(Y_i - x) + K_h(Y_i - x))$$

On remarquera que l'estimateur de la densité de X découle assez directement de la propriété (1).

Les auteurs de [2] proposent dans un second temps des estimateurs spécifiques pour le cas où les X_i sont non-négatifs permettant d'éviter les effets de bord proche de 0, à savoir les *convolution power kernel* présentés dans [1]. On définit pour cela pour k une densité sur \mathbb{R}^+ d'espérance 1 et K_m la densité de $\frac{1}{m} \sum_{i=1}^m E_i$ avec E_i *i.i.d* de densité k telle que

$$k_m(u) = mk^{*m}(mu)$$

où $k^{*m} = k \star \dots \star k$ m fois. On note de plus la transformée de fourrier de k_m la fonction k_m^* telle que

$$k_m^*(t) = \left(k^*\left(\frac{t}{m}\right)\right)^m, t \in \mathbb{R}.$$

Pour $\alpha_1, \dots, \alpha_L$ des réels tels que $\sum_{j=1}^L \alpha_j = 1$, $k^{(1)}, \dots, k^{(L)}$ des densités sur \mathbb{R}^+ d'espérance 1, on définit le noyau par puissance de convolutive (*convolution power kernel* ou *CPK*) par

$$K_m = \sum_{j=1}^L \alpha_j k_m^{(j)}$$

On peut donc désormais définir les estimateurs de la fonction de densité

$$\tilde{f}_m(x) = \frac{1}{nx} \sum_{i=1}^n \left[K_m\left(\frac{Y_i}{x}\right) + \frac{Y_i}{x} K_m'\left(\frac{Y_i}{x}\right) \right]$$

et de la fonction de survie

$$\tilde{\bar{F}}_m(x) = \frac{1}{x} \int_{\mathbb{R}^+} K_m\left(\frac{u}{x}\right) \hat{F}_Y(u) du + \frac{1}{nx} \sum_{i=1}^n Y_i K_m\left(\frac{Y_i}{x}\right).$$

Les auteurs fournissent des propriétés théoriques de l'estimateur de la fonction de survie que nous n'étudierons pas dans ce rapport.

3 Propriétés théoriques

Nous nous proposons dans cette section d'étudier les principaux résultats théoriques pour les estimateurs définis dans (7), i.e. les erreurs quadratique moyenne ponctuelles (voir Proposition (2)) et intégrée (voir Proposition (3)) ainsi que l'estimation des risques pour l'estimation adaptative de h par la méthode de Lepski (voir Théorème (2)).

Nous montrons dans un premier temps un résultats théorique préliminaire nécessaire pour la relation (3).

Proposition 1 (Densité du produit de deux variables indépendantes). *Soit $Z = XY$ une variable aléatoire avec X et Y deux variables aléatoires indépendantes de densités de probabilité respectives f_X et f_Y . Alors la densité de probabilité de Z est donnée par :*

$$f_Z(z) = \int_{\mathbb{R}} \frac{1}{|x|} f_X(x) f_Y(z/x) dx \quad (5)$$

Démonstration. Soit F_Z la fonction de répartition de Z . Par définition, on a

$$\begin{aligned} F_Z(z) &= \mathbb{P}(Z \leq z) \\ &= \mathbb{P}(XY \leq z) \\ &= \mathbb{P}(XY \leq z | X \geq 0) + \mathbb{P}(XY \leq z | X < 0) && \text{car } X \text{ indépendant de } Y \\ &= \mathbb{P}(Y \leq z/X | X \geq 0) + \mathbb{P}(Y \leq z/X | X < 0) \\ &= \int_{\mathbb{R}_+} f_X(x) \int_{-\infty}^{z/x} f_Y(y) dy dx + \int_{\mathbb{R}_-} f_X(x) \int_{z/x}^{+\infty} f_Y(y) dy dx \end{aligned}$$

En dérivant par rapport à z on obtient

$$\begin{aligned} f_Z(z) &= \int_{\mathbb{R}_+} f_X(x) f_Y(z/x) \frac{1}{x} dx - \int_{\mathbb{R}_-} f_X(x) f_Y(z/x) \frac{1}{x} dx \\ &= \int_{\mathbb{R}} \frac{1}{|x|} f_X(x) f_Y(z/x) dx \end{aligned}$$

□

Nous pouvons désormais étudier le risque ponctuel des estimateurs de la densité et de la fonction de survie. Considérons les hypothèse suivante

1. K est bornée
2. K est pair avec $\lim_{u \rightarrow +\infty} K(u) = 0$, K dérivable, K' intégrable
3. $\int (K'(u))^2 du < +\infty$
4. $\int |u| (K'(u))^2 du < +\infty$

Proposition 2 (Risque carré moyen ponctuel). *Supposons $\mathbb{E}[|X_1|] < +\infty$. Soit $x_0 \in \mathbb{R}$. Supposons que f est Hölder pour I un voisinage de x_0 , i.e. $f \in \Sigma_I(\beta, R)$. Si le noyau*

K est d'ordre $l + 1$ avec $l = \lfloor \beta \rfloor$ et $\int |u|^{\beta+1} |K(u)| du < +\infty$ alors si l'hypothèse (1) est satisfaite, le risque quadratique ponctuel de l'estimateur de la fonction de survie satisfait

$$\mathbb{E}[(\hat{F}_h(x_0) - \bar{F}(x_0))^2] \leq C_1^2 h^{2(\beta+1)} + \frac{C_2}{nh} + \frac{C_3}{n} \quad (6)$$

Plus particulièrement, en $x_0 = 0$, on a

$$\mathbb{E}[(\hat{F}_h(0) - \bar{F}(0))^2] \leq C_1^2 h^{2(\beta+1)} + \frac{C_4}{n} \quad (7)$$

avec $C_1 = \frac{R}{(\beta+1)!} \int |u|^{\beta+1} |K(u)| du$, $C_2 = 2\mathbb{E}[|X_1|] \|K\|^2$, $C_3 = 2\|K\|^2$ et $C_4 = 2\|K\|^2 + \int |u| K^2(u) du$. Si le noyau K est d'ordre l avec $l = \lfloor \beta \rfloor$ et $\int |u|^\beta |K(u)| du < +\infty$ et que de plus les hypothèses (2) et (3) sont vérifiées, le risque quadratique ponctuel de l'estimateur de la densité satisfait

$$\mathbb{E}[(\hat{f}_h(x_0) - f(x_0))^2] \leq C_5^2 h^{2\beta} + \frac{C_6}{nh^3} \quad (8)$$

Plus particulièrement, en $x_0 = 0$, sous les hypothèses (2), (3) et (4) on a

$$\mathbb{E}[(\hat{f}_h(0) - f(0))^2] \leq C_5^2 h^{2\beta} + \frac{C_7}{nh^2} \quad (9)$$

Si $\mathbb{E}[1/X] = \|f_Y\|_\infty < +\infty$, sous les hypothèses (2) et (3) on a

$$\mathbb{E}[(\hat{f}_h(0) - f(0))^2] \leq C_5^2 h^{2\beta} + \frac{C_8}{nh} \quad (10)$$

avec $C_5 = \frac{R}{\beta!} \int |u|^\beta |K(u)| du$, $C_6 = 2(\mathbb{E}[|X_1|] \|K'\|^2 + \|K\|_\infty^2)$, $C_7 = \|K\|_\infty + \int |u| K'(u)^2 du$ et $C_8 = \|f\|_\infty \|K\|^2 + \|f_Y\|_\infty \int u^2 K'(u)^2 du$.

Démonstration. Commençons par étudier le risque quadratique ponctuel de \hat{f}_h . Pour tout $x_0 \in \mathbb{R}$, on a

$$\begin{aligned} \mathbb{E}[\hat{f}_h(x_0)] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n (Y_i K'_h(Y_i - x_0) + K_h(Y_i - x_0))\right] \\ &= \mathbb{E}[Y_1 K'_h(Y_1 - x_0) + K_h(Y_1 - x_0)] && \text{car les } Y_i \text{ sont i.i.d} \\ &= \mathbb{E}[K_h(X_1 - x_0)] && \text{par (4)} \\ &= \int K_h(y - x_0) f(y) dy \\ &= K_h \star f(x_0) =: f_h(x_0) \end{aligned}$$

Ainsi \hat{f}_h est un estimateur sans biais de f_h . On a donc

$$\mathbb{E}[(\hat{f}_h(x_0) - f(x_0))^2] = (f(x_0) - f_h(x_0))^2 + \text{Var}[\hat{f}_h(x_0)]$$

Estimation du terme de biais de \hat{f}_h :

$$\begin{aligned}
 f_h(x_0) - f(x_0) &= \frac{1}{h} \int K\left(\frac{y-x_0}{h}\right) f(y) dy - f(x_0) \\
 &= \frac{1}{h} \int K\left(\frac{y-x_0}{h}\right) (f(y) - f(x_0)) dy \\
 &= \int K(u) (f(x_0 + uh) - f(x_0)) dy \quad \text{en posant } u = \frac{y-x_0}{h}
 \end{aligned}$$

Par un développement de Taylor-Young de $f(x_0 + uh)$ on obtient

$$f(x_0 + uh) = f(x_0) + uh f'(x_0) + \dots + \frac{(uh)^l}{l!} f^{(l)}(x_0 + \tau uh) \quad (11)$$

avec $0 \leq \tau \leq 1$. Or K est d'ordre $l = \lfloor \beta \rfloor$ et on a donc

$$\begin{aligned}
 f_h(x_0) - f(x_0) &= \int K(u) \frac{(uh)^l}{l!} f^{(l)}(x_0 + \tau uh) du \\
 &= \int K(u) \frac{(uh)^l}{l!} f^{(l)}(x_0 + \tau uh) du - \int K(u) \frac{(uh)^l}{l!} f^{(l)}(x_0) du \quad \text{car } K \text{ d'ordre } l \\
 &= \frac{h^l}{l!} \int K(u) u^l (f^{(l)}(x_0 + \tau uh) - f^{(l)}(x_0)) du
 \end{aligned}$$

En passant à la valeur absolue on obtient,

$$\begin{aligned}
 |\hat{f}_h(x_0) - f(x_0)| &\leq \frac{h^l}{l!} \int |K(u)| |u|^l |f^{(l)}(x_0 + \tau uh) - f^{(l)}(x_0)| du \\
 &\leq \frac{Rh^l}{l!} \int |u|^l |K(u)| |\tau uh|^{\beta-l} du \\
 &= \frac{Rh^\beta}{l!} \int |u|^\beta |K(u)| |\tau|^{\beta-l} du
 \end{aligned}$$

Et on a finalement,

$$(\hat{f}_h(x_0) - f(x_0))^2 \leq h^{2\beta} \left(\frac{R}{l!} \int |u|^\beta |K(u)| du \right)^2$$

Etudions désormais le **terme de variance ponctuelle de \hat{f}_h** :

$$\begin{aligned}
\text{Var}[\hat{f}_h(x_0)] &= \text{Var}\left[\frac{1}{nh} \sum_{i=1}^n \left(\frac{Y_i}{h} K'\left(\frac{Y_i - x_0}{h}\right) + K\left(\frac{Y_i - x_0}{h}\right)\right)\right] \\
&\leq \frac{1}{nh^2} \text{Var}\left[\frac{Y_1}{h} K'\left(\frac{Y_1 - x_0}{h}\right) + K\left(\frac{Y_1 - x_0}{h}\right)\right] && \text{car les } Y_i \text{ sont i.i.d} \\
&= \frac{1}{nh^2} (\mathbb{E}\left[\frac{Y_1^2}{h^2} \left(K'\left(\frac{Y_1 - x_0}{h}\right)\right)^2 + \frac{Y_1^2}{h} K'\left(\frac{Y_1 - x_0}{h}\right) K\left(\frac{Y_1 - x_0}{h}\right) + K^2\left(\frac{Y_1 - x_0}{h}\right)\right]) \\
&= \frac{1}{nh^2} (\mathbb{E}\left[\frac{Y_1^2}{h^2} \left(K'\left(\frac{Y_1 - x_0}{h}\right)\right)^2\right] + \mathbb{E}\left[K^2\left(\frac{X_1 - x_0}{h}\right)\right]) && \text{Par (4) (12)}
\end{aligned}$$

Or on peut majorer chacun des deux termes comme suit :

$$\begin{aligned}
\mathbb{E}\left[K^2\left(\frac{X_1 - x_0}{h}\right)\right] &= \int K^2\left(\frac{y - x}{h}\right) f(y) dy \\
&\leq \min\{h\|f\|_\infty \int K^2(u) du, \|K\|_\infty^2 \int f(y) dy\} && \text{en posant } u = \frac{y - x}{h} \\
&= \min\{h\|f\|_\infty \|K\|^2, \|K\|_\infty^2\}
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}\left[Y_1^2 \left(K'\left(\frac{Y_1 - x_0}{h}\right)\right)^2\right] &\leq \mathbb{E}[|X_1|] \int \left(K'\left(\frac{y - x}{h}\right)\right)^2 dy && \text{par lemme} \\
&= h\mathbb{E}[|X_1|] \int (K'(u))^2 du && \text{en posant } u = \frac{y - x}{h} \\
&= h\mathbb{E}[|X_1|] \|K'\|^2
\end{aligned}$$

On obtient finalement

$$\text{Var}[\hat{f}_h(x_0)] \leq \frac{1}{nh^3} \mathbb{E}[|X_1|] \|K'\|^2 + \frac{1}{nh^2} \min\{h\|f\|_\infty \|K\|^2, \|K\|_\infty^2\}$$

Dans le **cas où** $x_0 = 0$, on a :

$$\begin{aligned}
\mathbb{E}\left[Y_1^2 \left(K'\left(\frac{Y_1 - x_0}{h}\right)\right)^2\right] &= \int (x_0 + uh) (K'(u))^2 f_Y(x_0 + uh) h du \\
&\leq h \int |x_0 + uh| (K'(u))^2 du && \text{car } |zf_Y(z)| \leq 1, \forall z \in \mathbb{R} \\
&\leq h^2 \int |u| (K'(u))^2 du && \text{pour } x_0 = 0
\end{aligned}$$

Ainsi,

$$\text{Var}[\hat{f}_h(0)] \leq \frac{\|K\|_\infty + \int |u| (K'(u))^2 du}{nh^2}$$

Par ailleurs, pour f_Y bornée et si $\int u^2(K'(u))^2 du < +\infty$ on obtient alors :

$$\mathbb{E}[Y_1^2(K'(\frac{Y_1 - x_0}{h}))^2] \leq h^3 \|f_Y\|_\infty \int u^2(K'(u))^2 du$$

Et on obtient pour $\|f\| < +\infty$

$$\text{Var}[\hat{f}_h(0)] \leq \frac{\|f\|_\infty \|K\|^2 + \|f_Y\|_\infty \int u^2(K'(u))^2 du}{nh}$$

Etudions maintenant le **risque quadratique ponctuel de l'estimateur de la fonction de survie** \hat{F}_h . Pour tout $x_0 \in \mathbb{R}$ on a :

$$\begin{aligned} \mathbb{E}[\hat{F}_h(x_0)] &= \mathbb{E} \left[\int K_h(u - x_0) \mathbb{1}_{Y_1 \geq u} du + Y_1 K_h(Y_1 - x_0) \right] \quad \text{car les } Y_i \text{ sont i.i.d} \\ &= \int \int K_h(u - x_0) \mathbb{1}_{y \geq u} f_Y(y) du dy + \int y K_h(y - x_0) f_Y(y) dy \\ &= \int K_h(u - x_0) \int f_Y(y) \mathbb{1}_{y \geq u} dy du + \int y K_h(y - x_0) f_Y(y) dy \\ &= \int K_h(u - x_0) \bar{F}_Y(u) du + \int y K_h(y - x_0) f_Y(y) dy \\ &= \bar{F} \star K_h(x_0) =: \bar{F}_h(x_0) \end{aligned}$$

Donc \hat{F}_h est un estimateur sans biais de \bar{F}_h et on a

$$\mathbb{E}[(\hat{F}_h(x_0) - \bar{F}(x_0))^2] = (\bar{F}(x_0) - \bar{F}_h(x_0))^2 + \text{Var}[\hat{F}_h(x_0)]$$

Pour majorer le **terme de biais** on utilisera donc la même méthode que pour l'estimation du biais de \hat{f}_h mais avec un développement de Taylor-Young à l'ordre $l + 1$ et on obtient ainsi :

$$(\bar{F}(x_0) - \bar{F}_h(x_0))^2 \leq h^{2(\beta+1)} \left(\frac{R}{(l+1)!} \int |u|^{\beta+1} |K(u)| du \right)^2$$

Pour le **terme de variance** on a :

$$\text{Var}[\hat{F}_h(x_0)] \leq \frac{2}{nh^2} \left(\mathbb{E} \left[\int K_h(u - x) \mathbb{1}_{Y_1 \geq u} du \right]^2 + \mathbb{E}[Y_1^2 K_h^2(Y_1 - x)] \right) \quad \text{car } \text{Var}[X] \leq \mathbb{E}[X^2] \quad (13)$$

Pour le terme de droite on a

$$\begin{aligned} \mathbb{E}[Y_1^2 K_h^2(Y_1 - x)] &\leq \mathbb{E}[X_1^2 K^2(UX - x)] \\ &\leq h \mathbb{E}[|X_1|] \|K\|^2 \end{aligned}$$

Pour le terme de gauche

$$\begin{aligned} \mathbb{E} \left[\int K_h(u-x) \mathbb{1}_{Y_1 \geq u} du \right]^2 &\leq \left[\int |K_h(u-x)| du \right]^2 \\ &\leq h^2 \|K\|^2 \end{aligned}$$

Ainsi, finalement

$$\text{Var}[\hat{F}_h(x_0)] \leq \frac{2}{n} \|K\|^2 + \frac{2}{nh} \mathbb{E}[|X_1|] \|K\|^2$$

□

Proposition 3 (Risque carré moyen intégré). *Supposons $\mathbb{E}[X_1^2] < +\infty$. Si f est de carré intégrable sur \mathbb{R} , i.e. $f \in \mathbb{L}^2(\mathbb{R})$ et sous les hypothèses (2) et (3) on a*

$$\mathbb{E}[\|\hat{f}_h - f\|^2] \leq \|f - f_h\|^2 + \frac{\|K\|^2}{nh} + \frac{\mathbb{E}[Y_1^2] \|K'\|^2}{nh^3} \quad (14)$$

Si par ailleurs X est positif, $\bar{F} \in \mathbb{L}^2(\mathbb{R}_+)$ et K est à support compact dans $[-1, 1]$ alors $\forall h \leq 1$

$$\mathbb{E} \left[\int_{\mathbb{R}_+} (\hat{F}_h(x) - \bar{F}(x))^2 dx \right] \leq \int_{\mathbb{R}_+} (\bar{F}_h(x) - \bar{F}(x))^2 dx + \frac{2\mathbb{E}[Y_1^2] \|K\|^2}{nh} + \frac{2\mathbb{E}[Y_1 + 1] \|K\|_1^2}{n} \quad (15)$$

Démonstration. Par la décomposition biais variance, nous avons pour le **risque carré moyen intégré de la densité**, pour tout $x_0 \in \mathbb{R}$:

$$\mathbb{E}[\|\hat{f}_h - f\|^2] \leq \|f - f_h\|^2 + \int_{\mathbb{R}} \text{Var}[\hat{f}_h(x_0)] dx_0$$

Ainsi en intégrant (12) on obtient :

$$\begin{aligned} \int_{\mathbb{R}} \text{Var}[\hat{f}_h(x_0)] dx_0 &\leq \frac{1}{nh^2} \left(\mathbb{E} \left[\int_{\mathbb{R}} \frac{Y_1^2}{h^2} (K'(\frac{Y_1 - x_0}{h}))^2 dx_0 \right] + \mathbb{E} \left[\int_{\mathbb{R}} K^2(\frac{X_1 - x_0}{h}) dx_0 \right] \right) \\ &\leq \frac{1}{nh^2} \left(\mathbb{E} \left[\frac{Y_1^2}{h} \int_{\mathbb{R}} (K'(u))^2 du \right] + \mathbb{E} \left[\int_{\mathbb{R}} K^2(v) dv \right] \right) \quad \text{avec } u = \frac{Y_1 - x_0}{h} \text{ et } v = \frac{X_1 - x_0}{h} \\ &\leq \frac{\mathbb{E}[Y_1^2] \|K'\|^2}{nh^3} + \frac{\|K\|^2}{nh} \end{aligned}$$

De la même façon, on utilise la décomposition biais variance pour le **risque carré moyen intégré de la fonction de survie**. En intégrant (13) on obtient :

$$\int_{\mathbb{R}} \text{Var}[\hat{F}_h(x_0)] dx_0 \leq \frac{2}{nh^2} \left(\int_{\mathbb{R}_+} \mathbb{E} \left[\int K(\frac{u-x_0}{h}) \mathbb{1}_{Y_1 \geq u} du \right]^2 dx_0 + \mathbb{E} \left[Y_1^2 \int_{\mathbb{R}_+} k^2(\frac{Y_1 - x_0}{h}) dx_0 \right] \right)$$

On majore facilement le terme de droite par un changement de variable en posant

$u = \frac{Y_1 - x_0}{h}$ et on a alors

$$\mathbb{E} \left[Y_1^2 \int_{\mathbb{R}_+} k^2 \left(\frac{Y_1 - x_0}{h} \right) dx_0 \right] = \frac{2}{nh} \|K\|^2 \mathbb{E}[Y_1^2]$$

Pour le terme de gauche, on obtient en utilisant le fait que K est à support compact sur $[-1, 1]$, $u \in [x - h, x + h] \in [-1, +\infty]$ pour $x \in \mathbb{R}_+$ et $h \leq 1$ on obtient que

$$\begin{aligned} \int_{\mathbb{R}_+} \mathbb{E} \left[\int K \left(\frac{u - x_0}{h} \right) \mathbb{1}_{Y_1 \geq u} du \right]^2 dx_0 &= \mathbb{E} \left[\int_{\mathbb{R}} \int (K_h(u - x_0) \mathbb{1}_{Y_1 \geq u} \mathbb{1}_{u \geq -1} du) dx_0 \right] \\ &= \mathbb{E}[\|K_h \star g_{Y_1}\|^2] \quad \text{avec } g_{Y_1} = \mathbb{1}_{Y_1 \geq u} \mathbb{1}_{u \geq -1} \\ &\leq \mathbb{E}[\|K\|_1^2 \|g_{Y_1}\|^2] \quad \text{Par l'inégalité de Young (6)} \end{aligned}$$

Or $\|g_{Y_1}\|^2 = \int \mathbb{1}_{u \leq Y_1}^2 \mathbb{1}_{u \geq -1}^2 du = \int_{-1}^{Y_1} du = Y_1 + 1$ et on a finalement

$$\int_{\mathbb{R}} \text{Var}[\hat{F}_h(x_0)] dx_0 \leq \frac{2}{n} \|K\|_1 \mathbb{E}[Y_1 + 1] + \frac{2}{nh} \|K\|^2 \mathbb{E}[Y_1^2]$$

Et on retrouve donc bien (15). □

Il serait possible comme cela est fait dans [6] d'évaluer le biais en utilisant les classes de régularité de Nikol'ski (voir 4). Mais comme cela est rappelé dans [2], la régularité des fonctions estimées étant inconnue, il semble préférable d'utiliser une méthode adaptative (ici la méthode de Lepski développée dans [4]) menant à une borne du risque intégré non-asymptotique. Pour cela on définit

$$\hat{f}_{h,h'}(x) = K_{h'} \star \hat{f}_h(x) \quad , \quad \text{et} \quad \hat{F}_{h,h'}(x) = K_{h'} \star \hat{F}_h(x)$$

Remarquons que comme $K_h \star K_{h'} = K_{h'} \star K_h$ on a $\hat{f}_{h,h'} = \hat{f}_{h',h}$ et $\hat{F}_{h,h'} = \hat{F}_{h',h}$. On parcourt ainsi \mathcal{H}_n , un ensemble fini de fenêtré et on choisit

$$\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}_n} (A(h) + V(h)) \quad \text{et} \quad h^* = \operatorname{argmin}_{h \in \mathcal{H}_n} (B(h) + W(h)) \quad (16)$$

avec

$$\begin{aligned} A(h) &= \sup_{h' \in \mathcal{H}_n} (\|\hat{f}_{h'} - \hat{f}_{h,h'}\|^2 - V(h'))_+ \quad \text{et} \quad B(h) = \sup_{h' \in \mathcal{H}_n} (\|\hat{F}_{h'} - \hat{F}_{h,h'}\|^2 - W(h'))_+ \\ V(h) &= \kappa_1 \|K\|_1 \left(\frac{\|K\|^2}{nh} + \frac{\mathbb{E}[Y_1^2] \|K'\|^2}{nh^3} \right) \quad \text{et} \quad W(h) = \kappa_2 \|K\|_1^2 \frac{\mathbb{E}[Y_1^2] \|K\|^2}{nh} \end{aligned} \quad (17)$$

Théorème 2. *Supposons que f appartient à $\mathbb{L}^2(\mathbb{R})$, $\mathbb{E}[X_1^8] < +\infty$ et \mathcal{H}_n est tel que*

1. $\text{Card}(\mathcal{H}_n) \leq n$,
2. $\forall a > 0, \exists \Sigma(a) > 0$ tel que $\sum_{h \in \mathcal{H}_n} h^{-2} \exp(-a/h) < \Sigma(a) < \infty$,
3. $\forall \mathcal{H}_n, \frac{1}{nh^3} \leq 1$.

Alors, en supposant que $\int |K'(u)| du < +\infty$ et que $\int |u| (K'(u))^2 du < +\infty$, il existe une constante κ_1 dans $V(h)$ défini par (17) telle que

$$\mathbb{E}[\|\hat{f}_{\hat{h}} - f\|^2] \leq c \inf_{h \in \mathcal{H}_n} (\|K\|_1^2 \|f - f_h\|^2 + V(h)) + \frac{c'}{n} \quad (18)$$

Où c est une constante et c' dépend de K et f_Y .

Si par ailleurs X est positive, $\bar{F} \in \mathbb{L}^2(\mathbb{R}_+)$, \mathcal{H}_n satisfait (1.) et (2.) et K à support compact sur $[-1, 1]$, alors il existe une constante κ_2 défini par (17) telle que

$$\mathbb{E}\left[\int_{\mathbb{R}_+} (\hat{F}_{\hat{h}^*}(x) - \bar{F}(x))^2 dx\right] \leq c_1 \inf_{h \in \mathcal{H}_n} (\|K\|_1^2 \int_{\mathbb{R}_+} (\bar{F}_h(x) - \bar{F}(x))^2 dx + W(h)) + \frac{c'_1}{n} \quad (19)$$

où c_1 est une constante et c'_1 dépend de K et f_Y .

Démonstration. Commençons par décomposer le risque à l'aide de Cauchy-Schwarz. Pour tout $h, h' \in \mathcal{H}_n$, on a :

$$\begin{aligned} \|\hat{f}_{\hat{h}} - f\| &\leq 3\|\hat{f}_{\hat{h}} - \hat{f}_{h, \hat{h}}\|^2 + 3\|\hat{f}_{h, \hat{h}} - \hat{f}_{\hat{h}}\|^2 + 3\|\hat{f}_{\hat{h}} - f\|^2 \\ &\leq 3(A(h) + V(\hat{h})) + 3(A(\hat{h}) + V(h)) + 3\|\hat{f}_{\hat{h}} - f\|^2 \end{aligned}$$

Car par définition de $A(h)$,

$$\forall h, h' \in \mathcal{H}_n, \|\hat{f}_{h, h'} - f_{h'}\| \leq A(h) + V(h')$$

$$\forall h \in \mathcal{H}_n, A(\hat{h}) + V(\hat{h}) \leq A(h) + V(h)$$

Donc

$$\|\hat{f}_{\hat{h}} - f\| \leq 6A(h) + 6V(h) + 3\|\hat{f}_{\hat{h}} - f\|^2$$

Ainsi on a

$$\mathbb{E}[\|\hat{f}_{\hat{h}} - f\|^2] \leq 3\mathbb{E}[\|\hat{f}_{\hat{h}} - f\|^2] + 6\mathbb{E}[A(h)] + 6V(h)$$

Le terme de gauche est contrôlé par (14) et le terme de droite est déterministe. Nous nous intéressons donc seulement à $\mathbb{E}[A(h)]$. Notons $f_h(x) = \mathbb{E}[\hat{f}_h(x)]$ et $f_{h, h'}(x) = \mathbb{E}[\hat{f}_{h, h'}(x)]$

$$\begin{aligned} A(h) &\leq 5 \sup_{h' \in \mathcal{H}_n} \left\{ \|\hat{f}_{h'}^{(1)} - f_{h'}^{(1)}\| - V(h')/10 \right\}_+ + 5 \sup_{h' \in \mathcal{H}_n} \left\{ \|\hat{f}_{h, h'}^{(1)} - f_{h, h'}^{(1)}\| - V(h')/10 \right\}_+ \\ &\quad + 5 \sup_{h' \in \mathcal{H}_n} \left\{ \|\hat{f}_{h'}^{(2)} - f_{h'}^{(2)}\| \right\} + 5 \sup_{h' \in \mathcal{H}_n} \left\{ \|\hat{f}_{h, h'}^{(2)} - f_{h, h'}^{(2)}\| \right\} + 5 \sup_{h' \in \mathcal{H}_n} \left\{ \|\hat{f}_{h'} - f_{h'}\| \right\} \\ &=: 5(T_1 + T_2 + T_3 + T_4 + T_5) \end{aligned}$$

Les termes T_3 et T_4 sont bornés en utilisant les hypothèses (1), (2) et (3). On a en particulier $T_3 \leq C(K)n\mathbb{E}[|Y_1|^{2+p}/c_n^p]$, $\forall p > 0$ avec $C(K)$ une constante dépendant de K . Finalement Par l'inégalité de Young (6) avec $p = 1, q = r = 2$ et avec $\|K_{h'}\| =$

$\|K\|_1$ on obtient pour T_5 :

$$\begin{aligned} T_5 &= \sup_{h' \in \mathcal{H}_n} \left\{ \|\hat{f}_{h'} - f_{h,h'}\| \right\} \\ &= \sup_{h' \in \mathcal{H}_n} \left\{ \|K_{h'} \star (f - K_h \star f)\|^2 \right\} \\ &\leq \|K\|_1^2 \|f - K_h \star f\|^2. \end{aligned}$$

Pour T_1 :

$$\begin{aligned} T_1 &= \sup_{h' \in \mathcal{H}_n} \left\{ \|\hat{f}_{h'}^{(1)} - f_{h'}^{(1)}\| - V(h')/10 \right\}_+ \\ &\leq \sum_{h \in \mathcal{H}_n} \left\{ \|\hat{f}_h^{(1)} - f_h^{(1)}\| - V(h)/10 \right\}_+. \end{aligned}$$

Notons que

$$\begin{aligned} \|\hat{f}_h^{(1)} - f_h^{(1)}\|^2 &= \sup_{t \in \mathbb{L}^2(\mathbb{R}), \|t\|=1} \langle \hat{f}_h^{(1)} - f_h^{(1)}, t \rangle^2 \\ &= \sup_{t \in \mathcal{B}(1)} \langle \hat{f}_h^{(1)} - f_h^{(1)}, t \rangle^2 \\ &=: \sup_{t \in \mathcal{B}(1)} \nu_{n,h}(\psi_t)^2 \end{aligned}$$

où $\mathcal{B}(1)$ est une famille dénombrable dense de fonctions $t \in \mathbb{L}^2(\mathbb{R})$ et telles que $\|t\| = 1$ et $\nu_{n,h}(\psi_t)$ est le processus empirique centré $\nu_{n,h}(\psi_t) = \frac{1}{n} \sum_{i=1}^n [\psi_t(Y_i) - \mathbb{E}[\psi_t(Y_i)]]$ avec :

$$\begin{aligned} \psi(t) &:= \int ((y/h^2)K'((y-x)/h) + (1/h)K((y-x)/h)) \mathbb{1}_{|y| \leq c_n} t(x) dx \\ &= [yK' \star t(y) + K_h \star t(y)] \mathbb{1}_{|y| \leq c_n} \end{aligned}$$

Ainsi on peut majorer $\mathbb{E}[T_1]$ tq :

$$\mathbb{E}[T_1] \leq \sum_{h \in \mathcal{H}_n} \mathbb{E}[\{ \sup_{t \in \mathcal{B}(1)} \nu_{n,h}^2(\psi(t)) - V(h)/10 \}_+]$$

On peut désormais appliquer le Théorème de Talagrand pour majorer cette espérance. Il nous faut pour cela calculer H, M et v . Pour M on a :

$$\begin{aligned} \sup_{t \in \mathcal{B}(1)} \|\psi_t\|_\infty &\leq \frac{\sqrt{2}}{h} \sup_{u \in \mathbb{R}} \left[\int \frac{c_n^2}{h^2} \left(K' \left(\frac{u-x}{h} \right) \right)^2 + K^2 \left(\frac{u-x}{h} \right) dx \right]^{1/2} \\ &\leq \frac{\sqrt{2}}{h} \left[\frac{c_n^2}{h} \|K\|^2 + h \|K\|^2 \right]^{1/2} \\ &\leq C(K) \frac{c_n}{h^{3/2}} := M \end{aligned}$$

On remarquera que $H^2 = V(H)/\kappa_1$ convient. On cherche enfin un v qui convient. On peut dans un premier temps majorer la variance par l'espérance au carré pour avoir :

$$\sup_{t \in \mathcal{B}(1)} \text{Var}[\psi_t(Y_1)] \leq \sup_{t \in \mathcal{B}(1)} \mathbb{E}[\psi_t^2(Y_1)]$$

avec

$$\psi_t^2(y) = [y^2(K'_h \star t)(y) + (y(K_h \star t)^2(y))'] \mathbb{1}_{|y| \leq c_n}.$$

En utilisant la Proposition (1), on obtient

$$\mathbb{E}[\psi_t^2(Y_1)] \leq \mathbb{E}[Y_1^2(k'_h \star t)^2(Y_1)] + \mathbb{E}[(H_h \star t)^2(X_1)].$$

Pour le terme de gauche, on utilise l'inégalité de Young (voir (6)) :

$$\begin{aligned} \mathbb{E}[Y_1^2(k'_h \star t)^2(Y_1)] &\leq \mathbb{E}[|X_1|] \|K'_h \star t\|^2 \\ &\leq \mathbb{E}[|X_1|] \|K'\|_1^2 \|t\|^2 = \mathbb{E}[|X_1|] \frac{\|K'\|_1^2}{h^2} \quad \text{car } \|t\| = 1 \end{aligned}$$

Pour le terme de droite on utilise l'inégalité de Young à deux reprises, avec dans un premier temps (*) $p = q = 2$ et $r = \infty$ puis dans un second temps (**) avec $p = 1$ et $q = r = 2$:

$$\begin{aligned} \mathbb{E}[(H_h \star t)^2(X_1)] &= \int (K_h \star t)^2(x) f(x) dx \\ &\leq \|K_h \star t\|_\infty \|K_h \star t\| \|f\| \quad \text{par inégalité de Young (*)} \\ &\leq \|K_h\| \|t\| \|K\|_1 \|t\| \|f\| \quad \text{par inégalité de Young (**)} \\ &= \frac{\|K\| \|K\|_1}{h^{1/2}} \|f\|. \end{aligned}$$

Ainsi on trouve $v = c(K, f)/h^2$ avec $c(K, f) = \|K'\|_1^2 \mathbb{E}[|X_1|] + \|K\|_1 \|K\| \|f\|$. Le théorème de Talagrand donne finalement, en prenant $\kappa_1 = 30$ et $\epsilon = 1/2$:

$$\mathbb{E}[\{\sup_{t \in \mathcal{B}(1)} \nu_{n,h}^2(\psi_t) - V(h)/10\}_+] \leq \frac{C_1}{n} \left(\frac{1}{h^2} \exp(-C_2/h) + \frac{c_n^2}{nh^3} \exp(-C_3 \frac{\sqrt{(n)}}{c_n}) \right)$$

Sous les hypothèse (2) et (3) on a $\frac{1}{nh^3} \leq 1$, $\sum_{h \in \mathcal{H}_n} h^{-2} \exp(-C_2/h) < \Sigma(C_2) < \infty$ et $\text{Card}(\mathcal{H}_n \leq n) \leq n$. Ainsi en choisissant $c_n = C_3 \sqrt{(n)}/(4 \log(n))$, on obtient $\mathbb{E}[T_1] \leq c/n$. On effectue la même procédure pour le second terme avec un facteur $\|K\|_1^2$ du fait d'une application de l'inégalité de Young supplémentaire. On obtient finalement un ordre de $1/n$ en choisissant $p = 6$ si $\mathbb{E}[X_1^8] < +\infty$. On retrouve donc bien (18).

La preuve de (19) suit globalement les même étapes de calcul, nous ne la détaillerons pas ici. \square

4 Discussion

Nous avons étudié dans ce rapport l'estimation de fonctions de densité et de survie dans le cadre du modèle de *multiplicative censoring* et nous avons pu observer que les méthodes employées sont dans l'ensemble similaires à celles utilisées pour l'estimation de densité par méthodes à noyau du cours [3]. Il est intéressant de voir comment une méthode découlant naturellement d'une généralisation des estimateurs par histogramme pour la fonction de répartition (proposé par Rosenblatt en 1956 dans [5]) s'applique également pour l'estimation de fonction de survie et de densité, y compris dans le cas d'un modèle bruité du type *multiplicative censoring model*.

Il semble, par ailleurs, que le livre de Tsybakov [6] fournisse un ensemble d'outils très généraux permettant l'estimation de divers fonctions décrivant le comportement de variables aléatoires de loi inconnue à l'aide de méthodes à noyaux. La méthode de Lepski [4] permet de plus de rendre ces estimateurs adaptatif en sélectionnant une fenêtre h optimale parmi un ensemble de fenêtres \mathcal{H}_n et s'appliquent parfaitement dans le cadre du *multiplicative censoring model*.

Les résultats pratiques obtenus par les auteurs semblent de plus confirmer l'intérêt de la méthode de Lepski qui semble plus efficace pour estimer la fenêtre \hat{h} optimale pour l'estimateur de la fonction de densité que la méthode par validation croisée qui montre des résultats globalement moins pertinent. Nous pouvons également noter que la méthode CPK semble plus performante que la méthode de Lepski pour l'estimation de la fonction de survie proche de 0 mais est en revanche plus coûteuse en temps de calcul.

On notera que les hypothèses permettant d'obtenir les majorations de la MSE et de la MISE sont relativement faible et qu'elle conviennent à une diversité de noyau, en particulier elles conviennent pour les noyaux gaussiens.

On remarquera par contre que la variance de l'estimateur de la fonction de densité fait apparaître un terme en $O(1/(nh^3))$ qui peut vite prendre de grandes valeurs pour h trop petit.

Enfin nous relevons le fait que les convergences L_1 et L_p pour $p > 2$ ne sont pas étudiées dans l'article [2] ce qui pourrait être l'objectif de travaux futurs.

Références

- [1] F. Comte and V. Genon-Catalot. Convolution power kernels for density estimation. *Journal of Statistical Planning and Inference*, 142(7) :1698–1715, July 2012.
- [2] F. Comte, V. Genon-Catalot, and E. BRUNEL. Nonparametric density and survival function estimation in the multiplicative censoring model. *Test*, 25(3) :570–590, Sept. 2016.
- [3] C. Dion. Estimation non-paramétrique - première partie. 2021.
- [4] A. Goldenshluger and O. Lepski. Bandwidth selection in kernel density estimation : Oracle inequalities and adaptive minimax optimality. *The Annals of Statistics*, 39(3) :1608 – 1632, 2011.
- [5] M. Rosenblatt. Remarks on Some Nonparametric Estimates of a Density Function. *The Annals of Mathematical Statistics*, 27(3) :832 – 837, 1956.
- [6] A. B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer, Dordrecht, series = Springer series in statistics, year = 2009, url = <https://cds.cern.ch/record/1315296>, doi = 10.1007/b13794.