

REPUBLIQUE ALGERIENNE DEMOCRATIQUE ET POPULAIRE
MINISTERE DE L'ENSEIGNEMENT SUPERIEUR ET DE LA RECHERCHE SCIENTIFIQUE

=====

UNIVERSITE MENTOURI – CONSTANTINE
FACULTE DES SCIENCES EXACTES

=====

DEPARTEMENT DE MATHEMATIQUES

N° d'ordre :

N° de série :

MEMOIRE PRESENTE POUR L'OBTENTION
DU
DIPLOME DE MAGISTERE
EN
MATHEMATIQUES

« Sur les processus empiriques :
Application à la statistique des données censurées »

Par
Abderrahim Kitouni

OPTION
Probabilités et Statistique

Devant le jury :

Président	M.	Z. Mohdeb	Professeur	Université Mentouri
Rapporteur	M ^{me}	F. Messaci	Professeur	Université Mentouri
Examineur	M ^{me}	S. Belaloui	M. C.	Université Mentouri
Examineur	M ^{me}	N. Nemouchi	M. C.	Université Mentouri

Soutenu le :

Remerciements

Je remercie en premier lieu Mme. Fatiha Messaci pour tout le temps qu'elle m'a consacré, ainsi que pour toute l'aide qu'elle a patiemment apportée à la construction de ce mémoire.

J'ai été très honoré que M. Zaher Mohdeb accepte la présidence du jury, et j'aimerais lui adresser de ce fait de vifs remerciements.

Enfin, je remercie sincèrement Mme. Nahima Nemouchi, et Mme. Soheir Belaloui pour avoir accepté de faire partie du jury, et pour le temps accordé à la lecture attentive du mémoire.

Table des matières

Introduction	ii
1 Introduction aux processus empiriques	1
1.1 Définitions	1
1.2 Théorème de Glivenko-Cantelli	2
1.3 Processus empirique uniforme	3
1.4 La convergence faible des processus	4
1.4.1 Définition	4
1.4.2 Les espaces \mathcal{C} et \mathcal{D}	5
1.4.3 Critère de convergence faible	6
1.5 Théorème de représentation de Skorokhod	8
1.6 Théorème de Donsker	9
2 La loi du logarithme itéré	11
2.1 LIL classique	11
2.2 LIL pour les processus empiriques	12
2.3 Lois fonctionnelles du logarithme itéré	12
3 Le processus empirique de Kaplan-Meier	14
3.1 Introduction aux données censurées	14
3.2 L'estimateur de Kaplan-Meier	16
3.3 La loi du logarithme itéré pour l'estimateur de Kaplan-Meier	17
3.4 Autres théorèmes limites pour le processus empirique de Kaplan-Meier	18
4 Fonctionnelles locales : cas de la censure à droite	20
4.1 Introduction	20
4.2 Lois fonctionnelles pour les processus basés sur des données censurées	24
4.3 Application : Estimation de la densité	39
5 Extension au cas de la censure double	43
5.1 Présentation du modèle de Patilea et Rolin (2006)	43
5.2 Quelques résultats préliminaires	45
5.3 Étude de simulation	45

Introduction

La théorie des processus empiriques est l'un des outils les plus puissants de la statistique semi-paramétrique. Ses débuts remontent aux années 1930 avec des résultats tels que le théorème de Glivenko-Cantelli (Glivenko, 1933; Cantelli, 1933) ou encore le test de Kolmogorov-Smirnov (Kolmogorov, 1933; Smirnov, 1939). Ces deux résultats ont en commun le fait qu'ils donnent la convergence de la quantité $\|F_n - F\|$ (la convergence presque sûre pour le premier et la loi asymptotique pour le second), où F est la fonction de répartition de la loi d'où l'échantillon est tiré, F_n est la fonction de répartition empirique et $\|\cdot\|$ est la norme de la convergence uniforme. Ce que l'on désigne par processus empirique n'est rien de plus que la quantité $F_n - F$ normalisée par le facteur \sqrt{n} .

L'extension du résultat de Smirnov a donné ce qu'on appelle la loi du logarithme itéré (Chung, 1949; Deheuvels, 1991), et plus tard des lois limites fonctionnelles, globales ou locales, (Finkelstein, 1971; Deheuvels, 1992; Deheuvels et Mason, 1992; Deheuvels, 2000).

Pour une description approfondie des résultats classiques concernant les processus empiriques et leurs applications, nous renvoyons aux ouvrages de Shorack et Wellner (1986); Pollard (1984); van der Vaart et Wellner (1996); Dudley (1999). Pour un texte introductif le livre de Kosorok (2008) ou la thèse de Viallon (2006) sont un bon début.

Dans les années 1980, Stute (voir Stute, 1982a,b, 1986a,b) a été l'un des premiers statisticiens (voir également Csörgő, 1981; Deheuvels, 1974), à faire un usage systématique des méthodes de la théorie des processus empiriques dans l'étude des propriétés asymptotiques d'estimateurs fonctionnels non paramétriques. Ses travaux ont concerné principalement les estimateurs à noyau.

Depuis, de nombreux auteurs, parmi lesquels nous citons Deheuvels, J.H.J. Einmahl, U. Einmahl et Mason (voir par exemple Einmahl et Mason, 2000; Deheuvels et Mason, 2004; Einmahl et Mason, 2005), ont introduit des techniques nouvelles pour aborder ces problèmes, en utilisant les lois limites fonctionnelles, ou encore des variantes locales de la théorie des processus empiriques indexés par des ensembles ou par des fonctions.

Dans ce mémoire, nous allons exposer certains résultats de base de la théorie des processus empiriques, ainsi que quelques résultats analogues dans le cas des données censurées. En particulier, nous étudions le résultat de Deheuvels et Einmahl (1996) avec plus de détails. Enfin, nous nous intéressons aux perspectives de recherche dans la généralisation

de ce résultats au cas de la censure double. Ceci est possible grâce à la loi du logarithme itéré de Messaci et Nemouchi (2011).

Dans le premier chapitre, nous exposons les définitions et les résultats de base sur les processus empiriques. Le deuxième chapitre est consacré à la loi du logarithme itéré, un résultat classique concernant la convergence des sommes de variables aléatoires i.i.d., ainsi que ses généralisations pour les processus empiriques. Le troisième chapitre décrit les données censurées, avant de donner une définition du processus empirique pour le cas de la censure aléatoire à droite en utilisant l'estimateur de Kaplan et Meier (1958), ainsi que les résultats obtenus dans ce cas, en particulier par Földes et Rejtő (1981); Einmahl et Koning (1992). Dans le quatrième chapitre, nous étudions en détails un des résultats de la théorie des processus empiriques, en l'occurrence une loi fonctionnelle du logarithme itéré pour les incréments du processus empirique de Kaplan-Meier, résultat obtenu par Deheuvels et Einmahl (1996). Nous détaillons aussi l'utilisation de ce résultat pour montrer une vitesse de convergence de l'estimateur à noyau de la densité. Le dernier chapitre est, quant à lui, consacré à l'étude du modèle de censure double de Patilea et Rolin (2006). Nous présentons la loi du logarithme itéré de Messaci et Nemouchi (2011) pour ce dernier estimateur ainsi que les perspectives de recherche concernant l'extension du résultat de Deheuvels et Einmahl (1996) à ce modèle. Nous présentons en outre une étude de simulation pour explorer les propriétés de l'estimateur à noyau de la densité obtenu pour ce même modèle.

Chapitre 1

Introduction aux processus empiriques

1.1 Définitions

Intuitivement, un processus empirique est un processus aléatoire qui dépend d'un échantillon. Par exemple, la fonction de répartition empirique.

Plus précisément, si l'on considère un espace probabilisable $(\mathcal{X}, \mathcal{A})$, et X_1, X_2, \dots, X_n un échantillon i.i.d de loi de probabilité P_X , on définit la mesure empirique P_n par :

$$P_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i},$$

où δ_x est la mesure de Dirac au point x .

Pour une famille \mathcal{S} d'ensembles mesurables, on définit alors le processus empirique $\{P_n(A), A \in \mathcal{S}\}$. Un processus empirique ainsi défini est dit *indexé par des ensembles*. Dans le cas réel, la fonction de répartition empirique peut s'écrire ainsi en prenant $\mathcal{S} = \{] - \infty, t], t \in \mathbb{R}\}$.

Pour une classe \mathcal{F} de fonctions mesurables, on peut définir un processus empirique $\{P_n f, f \in \mathcal{F}\}$ où

$$P_n f = \int f dP_n = \frac{1}{n} \sum_{i=1}^n f(X_i).$$

Un tel processus est dit *indexé par des fonctions*. Cette définition est plus générale car elle permet de décrire toutes les fonctions mesurables de l'échantillon. La définition précédente est obtenue en se limitant à des fonctions indicatrices.

De plus, pour peu que la classe \mathcal{F} soit dénombrable (resp. ait la puissance du continu), le processus empirique peut être vu comme un processus "classique" à indices entiers (resp. réels).

Souvent, ce que l'on appelle processus empirique est $\{\sqrt{n}(P_n(A) - P_X(A)), A \in \mathcal{S}\}$ ou bien $\{\sqrt{n}(P_n f - P_X f), f \in \mathcal{F}\}$. C'est-à-dire qu'au lieu de prendre la mesure empirique, on prend la différence entre cette dernière et la mesure de probabilité théorique.

Dans la suite, nous nous restreignons au cas réel, et nous appelons processus empirique un processus de la forme $\sqrt{n}(F_n(x) - F(x))$ où F_n est la fonction de répartition empirique ou un autre estimateur de la fonction de répartition. Par exemple, dans le cas des données censurées à droite, F_n est l'estimateur de Kaplan-Meier et on parle alors du *processus empirique de Kaplan-Meier*. Pour plus de détails concernant la théorie générale des processus empiriques, voir Shorack et Wellner (1986, chap. 26), Kosorok (2008) et van der Vaart et Wellner (1996).

1.2 Théorème de Glivenko-Cantelli

La loi forte des grands nombres permet de donner, en tout point de \mathbb{R} , la convergence presque sûre de la fonction de répartition empirique vers la fonction de répartition des observations. Le résultat suivant, aussi connu sous le nom de loi uniforme des grands nombres, est dû à Glivenko (1933) et Cantelli (1933) et donne une version uniforme sur \mathbb{R} de cette convergence.

Théorème 1.1. Soit X_1, X_2, \dots, X_n une suite de variables aléatoires indépendantes de même loi de probabilité, de fonction de répartition F . Et soit F_n la fonction de répartition empirique définie par

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n 1_{]-\infty, x]}(X_i).$$

Alors, $\sup_{x \in \mathbb{R}} |F_n(x) - F(x)| \rightarrow 0$ p.s.

Démonstration. voir Billingsley (1986, Théorème 20.6) □

Ce résultat inspire la notion de classe de Glivenko-Cantelli, qui est une généralisation de cette convergence pour un processus empirique quelconque :

Définition 1.1. Soit X_1, \dots, X_n un échantillon de loi P_X , et soit P_n la mesure empirique correspondante. Une classe \mathcal{F} de fonctions mesurables (resp. une famille \mathcal{S} d'ensembles mesurables) est dite de *Glivenko-Cantelli* (GC) par rapport à P_X si :

$$\sup_{f \in \mathcal{F}} |P_n f - P_X f| \rightarrow 0 \text{ p.s.}$$

(resp. $\sup_{A \in \mathcal{S}} |P_n(A) - P_X(A)| \rightarrow 0$ p.s.).

Une *classe de Glivenko-Cantelli universelle* est une classe de Glivenko-Cantelli par rapport à toutes les mesures de probabilité.

Par exemple, la classe de fonctions $\{1_{]-\infty, t]}, t \in \mathbb{R}\}$ est une classe GC universelle.

1.3 Processus empirique uniforme

Un résultat important, qui est souvent utilisé pour simplifier les preuves, est la réduction au cas uniforme. Ceci consiste à introduire un processus empirique uniforme (c'est à dire un processus empirique basé sur un échantillon de loi uniforme sur $[0, 1]$) qui est plus facile à étudier. L'extension au cas d'une loi arbitraire est souvent possible sans trop d'hypothèses, mais nous devons parfois supposer la continuité de la fonction de répartition.

Soit F une fonction de répartition, et soit F^{inv} son inverse généralisée définie pour $t \in [0, 1]$ par :

$$F^{inv}(t) = \inf\{x \in \mathbb{R} : F(x) \geq t\}.$$

On a les propriétés suivantes :

- $\forall t \in [0, 1] : F(F^{inv}(t)) \geq t$;
- $\forall x \in \mathbb{R} : F^{inv}(F(x)) \leq x$;
- $F(x) \geq t \iff F^{inv}(t) \leq x$;
- Si X suit la loi F , alors $\forall t \in [0, 1] : P(F(X) \leq t) \leq t$, et si t appartient à l'image de F alors $P(F(X) \leq t) = t$. En particulier, si F est continue alors $F(X)$ suit la loi uniforme sur $[0, 1]$.

Théorème 1.2. Soit ξ une variable aléatoire de loi uniforme sur $[0, 1]$. Et soit P_X une loi de probabilité sur \mathbb{R} , de fonction de répartition F . On définit la variable aléatoire X par $X = F^{inv}(\xi)$.

Alors X suit la loi P_X .

Démonstration. D'après les propriétés précédentes et la croissance de F , si $\xi \leq F(x)$ alors $X = F^{inv}(\xi) \leq F^{inv}(F(x)) \leq x$. Et si $X \leq x$, alors $F(x) \geq F(X) = F(F^{inv}(\xi)) \geq \xi$.

Ce qui montre que les événements $\{X \leq x\}$ et $\{\xi \leq F(x)\}$ sont équivalents, et on a alors $P(X \leq x) = P(\xi \leq F(x)) = F(x)$. □

Ce résultat élémentaire peut se généraliser au cas des processus empiriques : Considérons une suite $(U_n)_{n \geq 1}$ de v.a. i.i.d. de loi uniforme sur $[0, 1]$. Nous notons respectivement U_n et α_n la fonction de répartition empirique et le processus empirique associés

$$\alpha_n(t) = \sqrt{n}(U_n(t) - t) \quad \text{et} \quad U_n(t) = \frac{1}{n} \sum_{i=1}^n 1_{\{U_i \leq t\}}.$$

Théorème 1.3. Les suites de processus $\{F_n(t), t \in \mathbb{R}\}$ et $\{U_n(F(t)), t \in \mathbb{R}\}$ ont les mêmes lois de probabilité conjointes en n (c'est-à-dire que pour tout $k \in \mathbb{N}^*$ et pour tout $n_1, n_2, \dots, n_k \in \mathbb{N}$, $(F_{n_1}, F_{n_2}, \dots, F_{n_k})$ et $(U_{n_1}(F), U_{n_2}(F), \dots, U_{n_k}(F))$ ont la même loi). De même pour les processus $\{\sqrt{n}(F_n(t) - F(t)), t \in \mathbb{R}\}$ et $\{\alpha_n(F(t)), t \in \mathbb{R}\}$.

Démonstration. voir Shorack et Wellner (1986, Théorème 2, page 4) □

Ceci nous permet de limiter le travail probabiliste au cas des échantillons de loi uniforme, et de généraliser ces résultats en insérant la fonction de répartition F dans le résultat (On a parfois besoin de supposer que F est continue pour faire cela).

Si F est continue, et si on choisit $U_n = F(X_n)$ pour tout n dans le théorème précédent, alors non seulement on a égalité des distributions mais également l'égalité presque sûre. Ce résultat peut être généralisé sous des hypothèses assez faibles au cas où F n'est pas continue (voir par exemple Shorack et Wellner, 1986, p. 102).

1.4 La convergence faible des processus

1.4.1 Définition

Soit $(X_n)_{n \in \mathbb{N}}$ une suite de v.a. réelles. On dit que (X_n) converge en loi vers une variable X , et on note $X_n \xrightarrow{\mathcal{L}} X$, si la suite des fonctions de répartition (F_n) de (X_n) converge simplement vers la fonction de répartition F de X en tout point de continuité de cette dernière.

Ceci est équivalent à dire que la suite des lois de probabilité (P_n) des variables (X_n) converge faiblement vers la loi de probabilité de X , ce qui veut dire que pour toute fonction f continue et bornée sur \mathbb{R} , on a :

$$\int_{\mathbb{R}} f dP_n \longrightarrow \int_{\mathbb{R}} f dP.$$

Contrairement à la première, cette définition a l'avantage de pouvoir être généralisée à un espace métrique quelconque.

Dans la suite, (S, d) désigne un espace métrique muni de sa tribu borélienne \mathcal{S} (la tribu engendrée par les ouverts de S). Un tel espace est dit polonais s'il est séparable et complet.

Définition 1.2. On dit qu'une suite de mesures de probabilité (P_n) définies sur (S, \mathcal{S}) converge faiblement vers une mesure de probabilité P (définie sur le même espace (S, \mathcal{S})), et on note $P_n \Rightarrow P$, si pour toute fonction réelle f définie sur S continue et bornée, on a :

$$\int_S f dP_n \longrightarrow \int_S f dP.$$

Le théorème suivant donne d'autres définitions équivalentes de la convergence faible :

Théorème 1.4 (Portmanteau). Soient $(P_n, n \geq 1), P$ des mesures de probabilité sur (S, \mathcal{S}) . Les conditions suivantes sont équivalentes :

1. $P_n \Rightarrow P$;
2. $\int f dP_n \longrightarrow \int f dP$ pour toute fonction réelle f bornée et uniformément continue ;

3. $\limsup P_n(F) \leq P(F)$ pour tout fermé F ;
4. $\liminf P_n(O) \geq P(O)$ pour tout ouvert O ;
5. $\lim P_n(A) = P(A)$ pour tout $A \in \mathcal{S}$ tel que $P(\text{Fr } A) = 0$ ($\text{Fr } A$ désigne la frontière de l'ensemble A).

Démonstration. cf. Billingsley (1968, Théorème 2.1) □

Si l'espace S n'est pas séparable, la tribu \mathcal{S}_0 engendrée par les boules ouvertes peut être plus petite que la tribu borélienne \mathcal{S} . Dudley (1966, 1967) a introduit une théorie de la convergence faible qui utilise uniquement les ensembles de \mathcal{S}_0 et les fonctions \mathcal{S}_0 -mesurables.

1.4.2 Les espaces \mathcal{C} et \mathcal{D}

Soit M une famille de fonctions réelles définies sur un ensemble d'indices T . Pour $k \geq 1$ et $t = (t_1, \dots, t_k) \in T^k$, on note $\pi_t = \pi_{t_1, \dots, t_k}$ l'application projection de M dans \mathbb{R}^k définie par

$$\pi_t(x) = (x(t_1), \dots, x(t_k)).$$

On appelle *ensemble fini-dimensionnel* toute partie de M de la forme $\pi_t(B)$ où $t \in T^k$ et $B \in \mathcal{B}(\mathbb{R}^k)$ pour $k \geq 1$. L'ensemble de \mathcal{M}_0 de tout les ensembles fini-dimensionnels est une algèbre. On note \mathcal{M} la σ -algèbre engendrée par l'algèbre des ensembles fini-dimensionnels.

Soit $(X_t, t \in T)$ un processus stochastique à valeurs réelles. Si toutes ses trajectoires sont dans M , il peut être vu comme un élément aléatoire de M . Pour pouvoir utiliser la théorie de la convergence faible citée plus haut, il faut munir M d'une métrique dont la tribu borélienne coïncide avec la tribu engendrée par les ensembles fini-dimensionnels.

Dans la suite, nous nous restreignons au cas où $T = [0, 1]$ (car comme nous venons de le voir, on peut toujours se ramener à des variables uniformes sur $[0, 1]$), et l'ensemble M de fonctions est soit l'ensemble des fonctions continue, soit l'ensemble des fonctions "cadlag" (c'est-à-dire l'ensemble des fonctions continues à droite et ayant des limites à gauche en tout point).

On note par \mathcal{C} l'ensemble des fonctions réelles définies et continues sur $[0, 1]$, et on définit la norme uniforme d'une fonction x de \mathcal{C} par

$$\|x\| = \sup_{t \in [0, 1]} |x(t)|.$$

Cette norme induit une distance sur \mathcal{C} qui en fait un espace métrique séparable et complet. De plus, sa tribu borélienne que l'on note \mathcal{C} est engendrée par les ensembles fini-dimensionnels (Shorack et Wellner, 1986, Chapitre 2).

Seulement, les processus que l'on manipule ne sont pas toujours à trajectoires continues. C'est ce qui nous amène à travailler dans un espace plus grand. Comme une fonction de répartition est toujours continue à droite, nous allons considérer l'espace des fonctions "cadlag" définies sur $[0, 1]$, que l'on note \mathcal{D} .

Muni de la distance uniforme, \mathcal{D} est un espace complet mais n'est pas séparable. De plus, la tribu engendrée par les ensembles fini-dimensionnels est strictement incluse dans la tribu borélienne de $(\mathcal{D}, \|\cdot\|)$. Ceci fait que la norme uniforme est inappropriée pour l'étude de la convergence faible des processus à trajectoire dans \mathcal{D} .

Pour remédier à ce problème, Skorokhod (1956) a proposé la métrique définie par :

$$d(x, y) = \inf_{\lambda \in \Lambda} \max(\|x - y \circ \lambda\|, \|\lambda - I\|),$$

où Λ est l'ensemble des bijections continues et strictement croissantes de $[0, 1]$ dans lui-même et I est l'application identique.

Cette métrique induit une topologie séparable sur \mathcal{D} , qu'on appelle topologie de Skorokhod. Une suite (x_n) d'éléments de \mathcal{D} converge pour cette topologie vers un élément x de \mathcal{D} , s'il existe une suite (λ_n) d'éléments de Λ telle que : $x_n \circ \lambda_n \rightarrow x$, et $\lambda_n \rightarrow I$ uniformément sur $[0, 1]$.

Il est clair qu'une suite qui converge uniformément converge pour cette topologie, et que la réciproque est fautive en général. Ces deux notions sont cependant équivalentes si la fonction limite est continue. En particulier, la topologie trace sur \mathcal{C} de la topologie de Skorokhod coïncide avec la topologie uniforme.

Cependant, \mathcal{D} n'est pas complet pour cette métrique. En effet, la suite $(1_{[\frac{1}{2}, \frac{1}{2} + \frac{1}{n}[})_n$ est une suite de Cauchy car $d(1_{[\frac{1}{2}, \frac{1}{2} + \frac{1}{n}[}, 1_{[\frac{1}{2}, \frac{1}{2} + \frac{1}{m}[}) = |\frac{1}{n} - \frac{1}{m}|$, mais n'est pas convergente. Mais on peut définir une distance équivalente (au sens qu'elle induit la même topologie) qui fait de \mathcal{D} un espace complet. Ceci est justifié par le fait que la convergence faible des mesures de probabilité définies sur un espace S dépend uniquement de sa topologie, et non de la distance qui engendre cette topologie.

Un autre inconvénient de cette approche est qu'elle ne peut être utilisée pour les processus empiriques que si les observations sont des nombres réels. Une autre approche consiste à utiliser les notions de mesure extérieure et d'intégrale extérieure pour remédier à ce problème (voir par exemple Kosorok, 2008).

1.4.3 Critère de convergence faible

Un résultat important de la convergence faible, est que si $X_n \xrightarrow{\mathcal{L}} X$ alors pour toute fonction continue h on a $h(X_n) \xrightarrow{\mathcal{L}} h(X)$. Une conséquence de ce résultat, dans l'espace \mathcal{C} par exemple, est que la convergence faible de la loi de X_n entraîne la convergence faible des

lois fini-dimensionnelles (car les applications projections sont continues). Dans \mathcal{D} muni de la topologie de Skorokhod, les applications projections ne sont pas toutes continues, mais le résultat suivant reste valable (voir Billingsley, 1968, Théorème 5.1).

L'inverse n'est pas vrai en général, considérons par exemple la suite $(P_n)_{n \geq 0}$ de mesures de Dirac au points $x_n \in \mathcal{C}$ définis pour $n \geq 1$ par :

$$x_n(t) = \begin{cases} nt & \text{si } 0 \leq t \leq \frac{1}{n}, \\ 2 - nt & \text{si } \frac{1}{n} \leq t \leq \frac{2}{n}, \\ 0 & \text{si } \frac{2}{n} \leq t \leq 1, \end{cases}$$

et x_0 est la fonction identiquement nulle sur $[0, 1]$. Comme x_n ne converge pas uniformément (c'est-à-dire pour la topologie de \mathcal{C}) vers x_0 , la suite P_n ne converge pas faiblement vers P_0 . Il suffit pour s'en convaincre de prendre B la boule ouverte de \mathcal{C} de centre 0 (la fonction identiquement nulle) et de rayon $\frac{1}{2}$, on a alors : $P_0(\text{Fr } B) = 1$ (car P_0 est concentrée en 0) mais $P_n(B) = 0$ pour tout $n \geq 0$ (car $\|x_n - 0\| = 1$). Cependant, pour tout $k > 0$ et tout $(t_1, t_2, \dots, t_k) \in [0, 1]^k$, le vecteur $(x_n(t_1), x_n(t_2), \dots, x_n(t_k))$ converge vers le vecteur nul de \mathbb{R}^k (ce qui implique la convergence des lois fini-dimensionnelles).

Toutefois, comme les lois fini-dimensionnelles permettent de caractériser une loi (deux lois ayant les mêmes lois fini-dimensionnelles sont nécessairement identiques), la convergence des lois fini-dimensionnelles peut être utilisée pour montrer l'unicité de la limite si elle existe.

Définition 1.3. Considérons une famille Π de mesures de probabilité sur (S, \mathcal{S}) . On dit que Π est relativement compacte si toute suite d'éléments de Π admet une sous suite qui converge faiblement vers une mesure de probabilité (qui n'appartient pas nécessairement à Π)

Si (P_n) , la suite des lois de (X_n) , est relativement compacte alors de toute sous suite (P_{n_k}) on peut extraire une sous suite $(P_{n'_k})$ qui converge faiblement vers une mesure de probabilité Q . Si de plus les lois fini-dimensionnelles de P_n convergent vers les lois fini-dimensionnelles de P , alors nécessairement $P = Q$. En effet, la convergence de $P_{n'_k}$ vers Q assure la convergence des lois fini-dimensionnelles respectives, ce qui veut dire que P et Q ont les même lois fini-dimensionnelles, d'où $P = Q$. Comme toute sous suite admet une sous suite qui converge vers P alors la suite (P_n) converge elle même vers P .

Remarquons que si $P_n \Rightarrow P$ alors nécessairement (P_n) est relativement compacte. Le fait d'imposer la compacité relative n'est pas vraiment restrictif.

En conclusion, pour montrer la convergence faible d'une suite de mesures de probabilité définies sur un espace de fonction, il suffit de montrer la convergence des lois fini-dimensionnelles et la compacité relative.

La question qui se pose maintenant est de trouver un critère assez simple à vérifier pour la compacité relative, et c'est là qu'intervient la notion de tension.

Définition 1.4. Soit (S, d) un espace métrique et \mathcal{S} sa tribu borélienne.

Une mesure de probabilité P définie sur \mathcal{S} est dite tendue si :

$$\forall \varepsilon > 0, \exists K \subset S \text{ compact} : P(K) \geq 1 - \varepsilon.$$

Une famille Π de mesures de probabilité définies sur \mathcal{S} est dite tendue si :

$$\forall \varepsilon > 0, \exists K \subset S \text{ compact} : \inf_{P \in \Pi} P(K) \geq 1 - \varepsilon.$$

Remarquons que si l'espace S est séparable, alors toutes les mesures de probabilité définies sur \mathcal{S} sont tendues. On est maintenant en mesure d'énoncer le résultat suivant.

Théorème 1.5 (Prokhorov, 1956). Soit (S, d) un espace métrique et \mathcal{S} sa tribu borélienne, et soit Π une famille de mesures de probabilité définies sur \mathcal{S} .

- Si Π est tendue, alors elle est relativement compacte.
- Si S est séparable et complet, et si Π est relativement compacte, alors elle est tendue.

Démonstration. voir Billingsley (1968, théorèmes 6.1 et 6.2 p. 37) □

On peut donc énoncer le critère suivant :

Corollaire 1.1. Soit $(P_n)_{n \geq 0}$ une suite de mesures de probabilité, définies sur (C, \mathcal{C}) ou (D, \mathcal{D}) , et P une mesure de probabilité. Si les lois finie-dimensionnelles de P_n convergent vers celles de P et si la suite (P_n) est tendue, alors $P_n \Rightarrow P$.

Signalons le fait que dans la pratique, la tension est généralement prouvée non pas directement mais elle découle de conditions plus "manipulables". Pour des résultats spécifiques aux espaces \mathcal{C} et \mathcal{D} voir Billingsley (1968, Théorèmes 8.2, 8.3, 15.2 et 15.3).

1.5 Théorème de représentation de Skorokhod

Un autre résultat important de la théorie des processus empiriques est le théorème de représentation de Skorokhod.

Considérons une suite (F_n) de fonctions de répartition qui converge faiblement vers une fonction de répartition F , et ξ une v.a. de loi uniforme sur $[0, 1]$, et soit la suite (X_n^*) définie pour tout n par $X_n^* = F_n^{inv}(\xi)$ et $X^* = F^{inv}(\xi)$. Il est évident (d'après le Théorème 1.2) que X_n^* converge en loi vers X^* . De plus, la convergence faible de (F_n) entraîne que F_n^{inv} converge vers F^{inv} en tout point de continuité de cette dernière, ce qui assure que X_n^* converge presque sûrement vers X^* .

Ce résultat est un cas particulier de celui de Skorokhod (1956), mais illustre bien son utilité : Partant d'une suite de processus aléatoires on peut construire une suite de processus équivalents dont les trajectoires convergent presque sûrement. De là, on peut établir de nouveaux résultats pour cette suite et les généraliser (si possible) à la suite d'origine.

Théorème 1.6. Soit (S, d, \mathcal{S}) un espace Polonais, et soit $(X_n)_{n \geq 0}$ une suite d'éléments aléatoires de S de lois respectives $(P_n)_{n \geq 0}$ tels que $X_n \Rightarrow X_0$.

Alors il existe un espace de probabilité $(\Omega^*, \mathcal{A}^*, P^*)$ et une suite (X_n^*) d'applications mesurables de $(\Omega^*, \mathcal{A}^*, P^*)$ dans (S, \mathcal{S}) qui induisent les lois de probabilité P_n (ce qui revient à dire, dans le cas de processus, que les processus X_n et X_n^* sont équivalents pour tout n) telles que :

$$d(X_n^*, X_0^*) \longrightarrow 0 \text{ p.s.}$$

Démonstration. Billingsley (1971) a une preuve assez simple de ce résultat où il prend pour espace $(\Omega^*, \mathcal{A}^*, P^*)$ l'intervalle de Lebesgue, c'est à dire l'intervalle $[0, 1]$ muni de sa tribu borélienne et de la mesure de Lebesgue. \square

Ce résultat permet, par exemple, une preuve élégante du théorème de la fonction continue, ou une version du théorème de la convergence dominée pour la convergence en loi.

Plusieurs auteurs ont donné des extensions de ce résultat. Par exemple, Dudley (1968) a prouvé ce résultat sans l'hypothèse de complétude. Wichura (1970) et Fernandez (1974) ont montré une extension de ce résultat au cas où S n'est pas nécessairement séparable (dans ce cas, la définition de la convergence faible est donnée par rapport à la tribu engendrée par les boules ouvertes, qui coïncide avec la tribu borélienne dans le cas séparable), mais avec l'hypothèse que le support de P_0 soit séparable. Voir par exemple Dudley (1976) pour une preuve de ce résultat, aussi connu sous le nom du théorème de Skorokhod-Dudley-Wichura.

Par exemple, dans l'espace \mathcal{D} des fonctions cadlag on peut utiliser soit la distance de Skorokhod (qui a été introduite pour ce résultat même) ainsi que sa tribu borélienne, soit la distance uniforme et la tribu engendrée par les boules ouvertes, ainsi que la définition de Dudley (1966) de la convergence faible (à condition, bien sûr, que le support de la loi limite soit séparable pour la distance uniforme).

1.6 Théorème de Donsker

Le théorème de la limite centrale donne, pour tout $t \in \mathbb{R}$, une convergence en loi de $\sqrt{n}(F_n(t) - F(t))$ vers une v.a. de loi normale centrée et de variance $F(t)[1 - F(t)]$. Plus généralement, on peut montrer que pour tous $t_1, \dots, t_k \in \mathbb{R}$, le vecteur aléatoire $\sqrt{n}(F_n(t_1) - F(t_1), \dots, F_n(t_k) - F(t_k))$ converge en loi vers un vecteur Gaussien centré $(G(t_1), \dots, G(t_k))$ avec

$$\text{Cov}(G(s), G(t)) = F(\min(s, t)) - F(s)F(t) \text{ pour } s, t \in \{t_1, \dots, t_k\}.$$

Donsker (1952) a montré un résultat encore plus fort : en regardant la fonction de répartition empirique comme un élément aléatoire de l'espace \mathcal{D} des fonctions "cadlag",

$\sqrt{n}(F_n(x) - F(x))$ converge faiblement vers un processus Gaussien.

Théorème 1.7. Soit $(X_n)_{n \geq 0}$ une suite de variables aléatoires i.i.d. de fonction de répartition F , et $Y_n = \sqrt{n}(F_n - F)$ le processus empirique associé.

Alors : $Y_n \Rightarrow Y$ où Y est un processus Gaussien centré de fonction de covariance $E(Y(s)Y(t)) = F(s)(1 - F(t))$ pour $s \leq t$.

Chapitre 2

La loi du logarithme itéré

La loi du logarithme itéré (ou LIL pour “Law of the Iterated Logarithm”) est un des théorèmes limites importants de la statistique, nous commençons donc par donner la version classique de ce résultat (pour la somme d’une suite de variables aléatoires i.i.d.) avant de passer aux processus empiriques.

2.1 LIL classique

On considère une suite (X_n) de v.a. i.i.d. centrées de variance 1, et $S_n = \sum_{i=1}^n X_i$. La loi forte des grands nombres donne une convergence presque sûre de $\frac{S_n}{n}$ vers 0. Le théorème de la limite centrale, donne par contre la convergence en loi de la suite $\frac{S_n}{\sqrt{n}}$ vers la loi normale $\mathcal{N}(0, 1)$.

La loi du logarithme itéré, donne un résultat pour une suite “intermédiaire” : $\frac{S_n}{\sqrt{n \log_2 n}}$ où $\log_2 n = \log \log n$.¹ D’où le nom de “logarithme itéré”. Notons que cette suite converge en probabilité vers 0 (il suffit d’utiliser l’inégalité de Tchebychev), mais on n’a pas la convergence presque sûre :

Théorème 2.1 (Loi du logarithme itéré). *Soient (X_n) une suite de v.a. i.i.d. centrées de variance 1, et $S_n = \sum_{i=1}^n X_i$. Alors :*

$$\limsup_{n \rightarrow \infty} \frac{S_n}{\sqrt{n \log_2 n}} = \sqrt{2} \quad p.s.$$

Ce résultat, dû à Hartman et Wintner (1941) a été d’abord prouvé pour des v.a. de Bernoulli par Khinchine (1924), puis par Kolmogorov (1929) pour des v.a. de loi normale. Une preuve plus moderne de ce résultat est donnée par Acosta (1983).

1. Par souci de rigueur, on remplace parfois \log par \log_+ défini par $\log_+ x = \log(\sup(x, e))$, ce qui ne change rien vu que tous les résultats donnés sont pour la limite quand $n \rightarrow \infty$

2.2 LIL pour les processus empiriques

Dans cette section, nous allons exposer des résultats qui seront, ainsi que leur généralisation au cas des données censurées, essentiels pour la suite. Considérons d'abord le cas du processus empirique uniforme, la généralisation au cas d'une loi continue quelconque étant immédiate.

Soient U_1, \dots, U_n des v.a. i.i.d. de loi $\mathcal{U}_{[0,1]}$ et soient la fonction de répartition empirique $\mathbb{U}(t) = \frac{1}{n} \sum 1_{\{U_i \leq t\}}$ et le processus empirique $\alpha_n(t) = \sqrt{n}(\mathbb{U}(t) - t)$ associés.

Théorème 2.2 (Smirnov). *Soit $b_n = \sqrt{2 \log_2 n}$. Alors :*

$$\limsup_{n \rightarrow \infty} \frac{\|\alpha_n\|}{b_n} = \limsup_{n \rightarrow \infty} \frac{\|n(\mathbb{U}_n - \mathbb{I})\|}{\sqrt{nb_n}} = \frac{1}{2} \quad p.s.$$

où $\|x\| = \sup_{t \in [0,1]} |x(t)|$ est la norme de la convergence uniforme.

Le résultat suivant, qui a été prouvé indépendamment par Chung (1949), est plus fort au sens qu'il implique le résultat précédent (voir Shorack et Wellner, 1986, Chapitre 13) :

Théorème 2.3. *Soit (λ_n) une suite croissante de nombres positifs, alors :*

$$P \left(\limsup_{n \rightarrow \infty} \{\|\alpha_n\| \geq \lambda_n\} \right) = \begin{cases} 0 & \text{si } \sum \frac{\lambda_n}{n} \exp(-2\lambda_n^2) < \infty \\ 1 & \text{sinon} \end{cases}$$

2.3 Lois fonctionnelles du logarithme itéré

Un autre résultat très important est celui de Finkelstein (1971), il consiste à prouver une loi du logarithme itéré pour le processus empirique considéré comme un élément de \mathcal{B} , l'ensemble des fonctions réelles bornées, et pas seulement pour la norme uniforme de ce processus, comme c'est le cas pour les résultats précédents. Un tel résultat est parfois appelé loi fonctionnelle du logarithme itéré.

Avant d'aller plus loin, donnons quelques définitions

Définition 2.1. Soient $(X_n)_{n \geq 0}$ une suite d'éléments aléatoires d'un espace métrique (S, d) définies sur l'espace de probabilité $(\Omega, \mathcal{A}, \mathcal{P})$. On dit que (X_n) est presque sûrement relativement compacte dans (S, d) avec pour ensemble limite H , s'il existe $\Omega_0 \in \mathcal{A}$ avec $P(\Omega_0) = 1$ tel que pour tout $\omega \in \Omega_0$:

1. toute suite n' de nombres entiers admet une sous suite n'' telle que $X_{n''}(\omega)$ converge dans (S, d) ;
2. toute les valeurs d'adhérence de $X_n(\omega)$ appartiennent à H ;
3. pour tout $h \in H$, il existe une suite $n' = n_{h,\omega}$ telle que $X_{n'}(\omega)$ converge vers h .

Définition 2.2. Soit h une fonction définie sur un intervalle I et à valeurs réelles. On dit que h est absolument continue si $\forall \varepsilon > 0, \exists \delta > 0, \forall (]x_i, y_i])_{1 \leq i \leq N}$ intervalles disjoints de I :

$$\sum_{i=1}^N (y_i - x_i) < \delta \implies \sum_{i=1}^N |h(y_i) - h(x_i)| < \varepsilon.$$

On dit que h est absolument continue par rapport à une mesure μ (ou une fonction de répartition en sous entendant que c'est par rapport à la mesure de probabilité liée à cette fonction), si $\forall \varepsilon > 0, \exists \delta > 0, \forall (]x_i, y_i])_{1 \leq i \leq N}$ intervalles disjoints de I :

$$\sum_{i=1}^N \mu(]x_i, y_i]) < \delta \implies \sum_{i=1}^N |h(y_i) - h(x_i)| < \varepsilon.$$

Considérons l'ensemble

$$\mathcal{H} = \left\{ \begin{array}{l} h : h \text{ est absolument continue sur } [0, 1] \text{ avec} \\ h(0) = h(1) = 0 \quad \text{et} \quad \int_0^1 [h'(t)]^2 dt \leq 1 \end{array} \right\},$$

où h' est la dérivée au sens de Lebesgue de h (c'est à dire que $h(t) = \int_0^t h'(u) du$, les intégrales étant ici au sens de Lebesgue).

On donne d'abord le théorème pour le cas uniforme.

Théorème 2.4. Soit $b_n = \sqrt{2 \log_2 n}$. Alors la suite $\left\{ \frac{\alpha_n}{b_n} \right\}$ est presque sûrement relativement compacte dans $\mathcal{B}([0, 1])$ avec pour ensemble limite l'ensemble \mathcal{H} ci-dessus.

Il peut se généraliser au cas de variables aléatoires de lois continues quelconques :

Théorème 2.5. Soit $(X_n)_{n \geq 1}$ une suite de variables aléatoires i.i.d. ayant une fonction de répartition F définie et continue sur un intervalle $[a, b]$, et soit a_n le processus empirique associé. Alors la suite $\left\{ \frac{\alpha_n}{b_n} \right\}$ est presque sûrement relativement compacte dans $\mathcal{B}([0, 1])$ avec pour ensemble limite l'ensemble \mathcal{H}_F des fonctions f définies sur $[a, b]$ et vérifiant :

- $f(a) = f(b) = 0$,
- f est absolument continue par rapport à F ,
- $\int_a^b (df/dF)^2 dF \leq 1$ où df/dF est la dérivée de f par rapport à F .

Chapitre 3

Le processus empirique de Kaplan-Meier

3.1 Introduction aux données censurées

Dans de nombreuses applications statistiques, on est amené à étudier une variable de durée, c'est-à-dire un délai séparant un instant initial de l'instant où un événement d'intérêt est observé. C'est souvent le cas pour les applications médicales liées à l'étude des durées de survie. Pour cette raison, et bien que le champ d'application des méthodes que nous allons exposer ne se limite pas à la médecine ou à la biologie, nous utilisons la terminologie associées aux durées de survie. Celle-ci est la plus commune dans la littérature statistique, puisque c'est précisément dans le domaine médical que les avancées méthodologiques liées à ces variables ont été développées en premier.

Nous parlons de données censurées lorsque la durée de survie n'est connue que lorsqu'elle est dans les limites des durées d'observation. Ces limites pouvant être imposée par le type d'observation (par exemple, la durée d'hospitalisation d'un malade) ou par des événements fortuits (comme un accident ou une migration du patient au cours de son suivi). Par exemple, dans le cas dit de censure à droite, seul le minimum entre la durée de survie et une durée limite supérieure d'observation est connu, ainsi que l'indicateur exprimant si la durée de survie a été censurée ou non. Ce minimum constitue, en quelque sorte, une durée de participation qui est observée en ne donnant qu'une information partielle sur la vraie durée de survie. Il existe plusieurs autres catégories de modèles de censure, obtenus par des variations du principe de la censure à droite décrit ci-dessus. Parmi ceux-ci, mentionnons les suivants, en y incluant le modèle de censure à droite.

Censure à droite Il y a censure à droite lorsque la durée de survie est supérieure à la durée de participation. Un exemple typique est celui où l'événement considéré est le décès d'un patient malade, et la durée d'observation est une durée totale d'hospitalisation. On peut aussi observer ce genre de phénomène dans des études de fiabilité quand la panne d'un appareil ou d'un composant électronique ne permet pas de continuer l'observation pour un autre appareil ou composant. L'expérimentateur peut également fixer une date de fin d'expérience et les observations pour les individus pour lesquels on n'a pas observé l'événement d'intérêt avant cette date seront censurées à droite.

Pour ce type de censure, tout ce que l'on sait est que la vraie durée est supérieure à la durée observée.

Censure à gauche Il y a censure à gauche lorsque la durée de survie est inférieure à la durée observée.

Supposons par exemple que nous étudions la fiabilité d'un certain composant électronique qui est branché en parallèle avec un ou plusieurs autres composants. Une panne de ce composant n'entraîne pas nécessairement l'arrêt du système : le système peut continuer à fonctionner, quoique de façon aberrante, jusqu'à ce que cette panne soit détectée (par exemple lors d'un contrôle ou en cas de l'arrêt du système). La durée observée pour ce composant est alors censurée à gauche.

Un autre exemple est que l'on s'intéresse à l'âge à partir duquel une personne commence à accomplir une certaine tâche. Certaines personnes peuvent ne pas se rappeler, et donner juste une valeur supérieure (le cas extrême est que l'on prenne le début de l'étude comme observation). Cette donnée est donc censurée à gauche.

Dans ce dernier exemple, ainsi que dans beaucoup de cas, on trouve des données censurées à gauche dans un même échantillon que des données censurées à droite, ce qui conduit à la définition suivante.

Censure double (ou mixte) On dit qu'on a une censure double si on a des données censurées à droite et des données censurées à gauche dans le même échantillon. C'est l'objet d'étude du chapitre 5

Censure par intervalle Dans le cas de la censure par intervalle, on observe à la fois une borne inférieure et une borne supérieure de la durée d'intérêt. Ceci arrive dans des études de suivi médical où les patients sont contrôlés périodiquement, si un patient ne se présente pas à un ou plusieurs contrôles et se présente ensuite après que l'événement d'intérêt se soit produit. On a aussi pour ce genre d'expériences des données qui sont censurées à droite ou, plus rarement, à gauche. Un avantage de ce type est qu'il permet de représenter les données censurées à droite ou à gauche par des intervalles du type

$[a, +\infty[$ et $[0, a]$ respectivement, ce qui permet de considérer ce modèle comme étant plus générique. Turnbull (1976) présente ce genre de censure avec plus de détails.

Les trois catégories de censure décrits ci-dessus peuvent se décliner en fonction du mode ou mécanisme de censure. On obtient alors les types suivants :

Censure de type I L'expérimentateur fixe une date (non aléatoire) de fin d'expérience. La durée de participation maximale est alors fixée (non aléatoire) et vaut, pour chaque observation, la différence entre la date de fin d'expérience, et la date d'entrée du patient dans l'étude. Le nombre d'événements observés est, quant à lui, aléatoire. Ce modèle est souvent utilisé dans les études épidémiologiques.

Censure de type II L'expérimentateur fixe a priori le nombre d'événements à observer. La date de fin d'expérience devient alors aléatoire, le nombre d'événements étant, quant à lui, non aléatoire. Ce modèle est souvent utilisé dans les études de fiabilité.

Censure aléatoire C'est typiquement ce modèle qui est utilisé pour les essais thérapeutiques. Dans ce type d'expérience, la date d'inclusion du patient dans l'étude est fixée, mais la date de fin d'observation est inconnue (celle-ci correspond, par exemple, à la durée d'hospitalisation du patient). Ici, le nombre d'événements observés et la durée totale de l'expérience sont aléatoires.

Dans ce mémoire, nous nous restreignons à l'étude de la censure aléatoire, principalement pour les données censurées à droite. Dans le chapitre 5, nous présentons quelques perspectives de recherche pour les données doublement censurées.

3.2 L'estimateur de Kaplan-Meier

Soient X_1, \dots, X_n un échantillon représentant les durées d'intérêt (ces variables sont donc supposées positives), de fonction de répartition F , et C_1, \dots, C_n un échantillon représentant les temps de censure, que l'on suppose indépendants des durées d'intérêt, de fonction de répartition G . Dans le modèle de censure aléatoire à droite, on observe non pas la durée d'intérêt X_i mais plutôt la plus petite des deux valeurs $Z_i = \min(X_i, C_i)$, ainsi que l'indicateur de censure δ_i qui vaut 1 si la durée d'intérêt est observée, et 0 si elle est censurée, i.e. $\delta_i = 1_{\{X_i \leq C_i\}}$.

Dans ce genre de données, qui sont souvent des durées de survie ou des données de fiabilité, la fonction de répartition F est estimée par l'estimateur introduit par Kaplan et

Meier (1958), donné pour $z < Z_{(n)} = \max\{Z_1, \dots, Z_n\}$ par

$$F_n(z) = 1 - \prod_{i: Z_i \leq z} \left(\frac{N_n(Z_i) - 1}{N_n(Z_i)} \right),$$

où $N_n(x) = \sum_{i=1}^n 1_{\{Z_i \geq x\}}$. Pour $z \geq Z_{(n)}$, il y a plusieurs conventions pour définir $F_n(z)$: Soit on le définit par $F_n(Z_{(n)})$, ce qui fait que F_n peut ne pas être une fonction de répartition si $Z_{(n)}$ est une donnée censurée, soit on le définit par 0, soit on le laisse non défini.

Cet estimateur a des propriétés assez similaires à la fonction de répartition empirique, par exemple la convergence uniforme presque sûre (Winter *et al.*, 1978; Stute et Wang, 1993), la normalité asymptotique (Breslow et Crowley, 1974; Gill, 1983), et la loi du logarithme itéré (Földes et Rejtő, 1981). Ceci justifie que l'on s'intéresse à généraliser la théorie des processus empiriques au cas des données censurées.

Aussi, on définit le processus empirique de Kaplan-Meier par $a_n(t) = \sqrt{n}(F_n(t) - F(t))$.

3.3 La loi du logarithme itéré pour l'estimateur de Kaplan-Meier

Le résultat suivant est une loi du logarithme itéré pour l'estimateur de Kaplan-Meier de la fonction de répartition, qui est d'une certaine manière similaire au théorème de Chung-Smirnov.

On définit l'estimateur de Kaplan-Meier en utilisant la deuxième convention plus haut, c'est à dire en posant $F_n(z) = 0$ pour $z > Z_{(n)}$. Pour toute fonction de répartition L , on note par $T_L = \max\{t : L(t) < 1\}$ le point terminal du support de L .

Théorème 3.1 (Földes et Rejtő, 1981). *On suppose que F et G sont continues, et que $T_F < T_G$. Alors,*

$$P \left(\sup_{-\infty < u < +\infty} |F_n(u) - F(u)| = O \left(\sqrt{\frac{\log_2 n}{n}} \right) \right) = 1.$$

La condition $T_F < T_G$ pouvant paraître restrictive, on peut citer le théorème autrement : **Corollaire 3.1** (Földes et Rejtő, 1981). *On suppose que F et G sont continues, et on considère T tel que $G(T) > 0$. Alors,*

$$P \left(\sup_{-\infty < u \leq T^*} |F_n(u) - F(u)| = O \left(\sqrt{\frac{\log_2 n}{n}} \right) \right) = 1,$$

où $T^* = \min\{T, T_F\}$.

3.4 Autres théorèmes limites pour le processus empirique de Kaplan-Meier

Dans le modèle de censure aléatoire à droite que nous venons de voir, un rôle dominant est joué par les fonctions H et $H^{(1)}$ avec

$$H(t) = P(Z_i \leq t) = 1 - (1 - F(t))(1 - G(t)),$$

$$H^{(1)}(t) = P(Z_i \leq t, \delta_i = 1) = \int_0^t (1 - G(s)) dF(s),$$

ainsi que les versions empiriques

$$H_n(t) = \frac{1}{n} \sum_{i=1}^n 1_{[0,t]}(Z_i),$$

$$H_n^{(1)}(t) = \frac{1}{n} \sum_{i=1}^n \delta_i 1_{[0,t]}(Z_i).$$

La méthode classique pour traiter ce genre de données consistait à considérer $nH_n^{(1)}$ comme un processus de comptage (voir par exemple Aalen, 1976; Fleming et Harrington, 1991). Son compensateur est donné par

$$n \int_0^t (1 - H_n(s_-)) d\Lambda(s),$$

où $\Lambda(t) = -\log(1 - F(t))$ est la fonction de hasard de X_i . La martingale de base est définie par

$$M_n(t) = n \left(H_n^{(1)}(t) - \int_0^t (1 - H_n(s_-)) d\Lambda(s) \right).$$

C'est une martingale par rapport à la filtration $\mathcal{F}_x = \sigma\{Z_i 1_{\{Z_i \leq x\}}, \delta_i 1_{\{Z_i \leq x\}}\}$. La martingale de base joue un rôle fondamental car plusieurs processus intéressants peuvent s'écrire comme des intégrales stochastiques par rapport à cette martingale, c'est à dire sous la forme

$$Q_n(t) = \int_0^t L_n(s) dM_n(s).$$

Un processus particulièrement intéressant qui s'écrit sous cette forme est le processus empirique de Kaplan-Meier divisé par $1 - F$ qui s'écrit

$$\begin{aligned} \Pi_n(t) &= \frac{a_n(t)}{1 - F(t)} = \frac{\sqrt{n} (F_n(t) - F(t))}{1 - F(t)} \\ &= \frac{1}{\sqrt{n}} \int_0^t \frac{1 - F_n(s_-)}{1 - F(s)} \frac{1}{1 - H_n(s_-)} 1_{\{H_n(s_-) < 1\}} dM_n(s) \end{aligned}$$

Einmahl et Koning (1992) étudient une version pondérée du processus Q_n et donnent des conditions nécessaires et suffisantes pour obtenir des résultats analogues au cas des données complètes, à savoir le théorème de Chibisov (1964) et O'Reilly (1974), le théorème de Glivenko-Cantelli tel que formulé par Lai (1974) et Wellner (1977) ainsi que la loi fonctionnelle du logarithme itéré de James (1975).

Ces théorèmes, donnés dans un cadre général, s'appliquent au processus Π_n donné plus haut, et peuvent même s'appliquer (à quelque modifications près) au processus empirique de Kaplan-Meier lui-même.

Ces résultats, ainsi que d'autres résultats récents, n'utilisent cependant pas les méthodes de martingales (bien qu'ils utilisent ce processus ainsi que la martingale de base).

Bien entendu, les résultats valables pour le processus pondéré restent valables pour le processus originel (il suffit de prendre la fonction poids identiquement égale à 1).

Chapitre 4

Sur le comportement limite des fonctionnelles locales des processus empiriques dans un modèle de censure à droite

4.1 Introduction

Dans ce chapitre, nous allons étudier un résultat intéressant de la théorie des processus empiriques en la présence de données censurées à droite. Il s'agit d'une loi fonctionnelle du logarithme itéré pour les accroissements du processus empirique de Kaplan-Meier, aussi appelés processus empirique local.

Un tel résultat peut être utilisé pour donner des vitesses de convergence pour certains estimateurs à noyau. Nous nous intéressons tout particulièrement à l'estimateur à noyau de la densité, bien qu'il ne soit pas le seul estimateur auquel ce résultat s'applique.

L'estimateur à noyau de la densité est donné par :

$$f_n(z) = \int_{-\infty}^{+\infty} \frac{1}{h_n} K\left(\frac{t-z}{h_n}\right) dF_n(t),$$

où h_n (appelée fenêtre) est une suite de nombres strictement positifs et K (appelé noyau) est une fonction définie de \mathbb{R} dans \mathbb{R} , sur lesquels des conditions sont imposées plus loin. Cet estimateur est une généralisation de l'estimateur de Parzen-Rosenblatt (Parzen, 1962; Rosenblatt, 1956) au cas des données censurées à droite.

Les résultats de ce chapitre sont dûs à Deheuvels et Einmahl (1996). Remarquons qu'ils impliquent la convergence ponctuelle d'estimateurs de la densité, résultat qui a été

d'abord étendu au cas uniforme sur un compact dans Deheuvels (2000) ; ensuite dans Viallon (2006) ce dernier a été aussi généralisé à une convergence uniforme par rapport à une fenêtre appartenant à un compact bien spécifié.

Notations Nous utilisons les mêmes notations que précédemment : Soient X_1, \dots, X_n les durées d'intérêt, indépendantes et de fonction de répartition F , et indépendamment d'elles, soient C_1, \dots, C_n les durées de censure, indépendantes et de fonction de répartition G . Et on observe les couples $(Z_1, \delta_1), \dots, (Z_n, \delta_n)$ avec $Z_i = \min\{X_i, C_i\}$ et $\delta_i = 1_{\{X_i \leq C_i\}}$.

Pour toute fonction R , on note :

$$R_+(x) = \lim_{\varepsilon \downarrow 0} R(x + \varepsilon) \quad \text{et} \quad R_-(x) = \lim_{\varepsilon \downarrow 0} R(x - \varepsilon),$$

si ces limites existent, et pour toute fonction de répartition L , on note :

$$T_L = \sup\{t : L(t) < 1\} \quad \text{et} \quad L_-(\infty) = \lim_{x \rightarrow \infty} L(x).$$

On suppose que $F_-(\infty) = 1$, mais on accepte pour la distribution de Y que $G_-(\infty) = 1 - P(Y = \infty) \leq 1$. En particulier, si $P(Y = \infty) = 1$, on revient au cas non censuré.

Posons $\Theta = \min\{T_F, T_G\}$. On suppose que $\Theta > 0$ (si $\Theta = 0$, alors F_n est presque sûrement dégénérée).

On note f_+ et g_+ (resp. f_- et g_-) les dérivées à droite (resp. à gauche) de F et G (resp. de F_- et G_-) si elles existent.

Remarquons que l'existence de f_+ et f_- (resp. g_+ et g_-) en un point donné ne garantit pas la continuité de F (resp. de G) en ce point. Si F (resp. G) est continue, alors f_+ et f_- (resp. g_+ et g_-) ne sont que les dérivées à droite et à gauche de F (resp. G).

Si \mathcal{J} est un ensemble, on note $\mathbb{B}(\mathcal{J})$ l'ensemble des fonctions réelles bornées définies sur \mathcal{J} , muni de la norme de la convergence uniforme.

Hypothèses et résultats La fenêtre h_n doit vérifier les conditions suivantes :

H1 $h_n \downarrow 0$; $nh_n \uparrow \infty$.

H2 $\frac{nh_n}{\log_2 n} \rightarrow \infty$ où $\log_2 u = \log \log u$.

Ces conditions sont utilisées pour les théorèmes 4.1 et 4.2 ci dessous. Un exemple de fenêtre vérifiant ces hypothèses est de la forme $n^{-\alpha}$ avec $0 < \alpha < 1$. La condition (H1) peut être affaiblie en :

H1' Il existe une suite h_n^* vérifiant (H1) et telle que :

$$0 < \liminf_{n \rightarrow \infty} \left(\frac{h_n}{h_n^*} \right) \leq \limsup_{n \rightarrow \infty} \left(\frac{h_n}{h_n^*} \right) < \infty.$$

Les hypothèses sur le noyau K sont :

K1 K est à variation bornée

K2 K est à support compact

K3 $\int_{-\infty}^{+\infty} K(u) du = 1$.

Remarquons que la seule condition qui puisse être considérée comme restrictive est (K2). Parmi les noyaux les plus utilisés, seul le noyau Gaussien ne vérifie pas cette condition. De plus, ces hypothèses sont uniquement utilisées pour le Théorème 4.2.

L'article de Deheuvels et Einmahl (1996) est structuré comme suit : Dans un premier temps, ils donnent un résultat analogue (mais plus simple, vu qu'il ne considère que le cas réel) au résultat de Deheuvels et Mason (1994) dans le cas des données censurées à droite. Puis, ils appliquent ce théorème pour obtenir la vitesse de convergence de l'estimateur à noyau de la densité. L'avantage de cette approche réside dans le fait que les résultats ainsi obtenus peuvent être utilisés pour donner les vitesses de convergence d'autres types d'estimateurs.

Pour $z \in]0, \Theta[$ fixé, on considère la suite de fonctions aléatoires définie par :

$$\xi_n(u) = \frac{1}{b_n} (a_n(z + h_n u) - a_n(z)) \quad \text{pour } -M \leq u \leq M.$$

où M est une constante donnée, $b_n = \sqrt{2h_n \log_2 n}$, et a_n est le processus empirique de Kaplan-Meier défini par :

$$a_n(x) = \sqrt{n} (F_n(x) - F(x)).$$

Nous allons maintenant énoncer le résultat principal, objet de notre étude. Il donne le comportement asymptotique de la suite (ξ_n) :

Théorème 4.1 (Deheuvels et Einmahl, 1996, Théorème 1.2). *Soit $z \in]0, \Theta[$ fixé. On suppose que F est continue dans un voisinage de z et que $f_+(z)$ et $f_-(z)$ existent. Alors, sous les hypothèses (H1) et (H2), la suite $(\xi_n)_{n \geq 1}$ est presque sûrement relativement compacte dans $\mathbb{B}([-M, M])$ avec pour ensemble limite l'ensemble des fonctions $h \in \mathbb{B}([-M, M])$ de la forme*

$$h(u) = \int_0^u \psi(s) ds; \quad -M \leq u \leq M,$$

avec

$$\frac{1 - G_-(z)}{f_-(z)} \int_{-M}^0 \psi^2(s) ds + \frac{1 - G(z)}{f_+(z)} \int_0^M \psi^2(s) ds \leq 1.$$

Remarque 4.1. Si $f_-(z) = 0$, on remplace la condition précédente par :

$$\int_{-M}^0 \psi^2(s) ds = 0 \quad ; \quad \int_0^M \psi^2(s) ds \leq \frac{f_+(z)}{1 - G(z)};$$

et de même pour $f_+(z) = 0$.

Si F est dérivable en z et G continue en z , la condition précédente devient

$$h(u) = \int_0^u \psi(s) ds; \quad -M \leq u \leq M,$$

avec

$$\int_{-M}^M \psi^2(s) ds \leq \frac{f(z)}{1 - G(z)}.$$

La preuve de ce théorème se déduit d'une succession de lemmes à la section 4.2.

Nous allons maintenant énoncer le résultat pour l'estimateur à noyau de la densité. Posons

$$\mathbb{E}f_n(z) = \int_{-\infty}^{+\infty} \frac{1}{h_n} K\left(\frac{t-z}{h_n}\right) dF(t).$$

En général $\mathbb{E}f_n(z) \neq E(f_n(z))$, sauf dans le cas non censuré.

Théorème 4.2 (Deheuvels et Einmahl, 1996, Théorème 1.1). *Soit $z \in]0, \Theta[$ fixé. On suppose que F est continue dans un voisinage de z et admet des dérivées à droite et à gauche en z . Alors sous les hypothèses (H1), (H2), (K1), (K2) et (K3) on a :*

$$\limsup_{n \rightarrow \infty} \pm \sqrt{\frac{nh_n}{2 \log_2 n}} (f_n(z) - \mathbb{E}f_n(z)) = \left(\frac{f_-(z)}{1 - G_-(z)} \int_{-\infty}^0 K^2(t) dt + \frac{f_+(z)}{1 - G(z)} \int_0^{+\infty} K^2(t) dt \right)^{1/2} \quad p.s.$$

La preuve de ce théorème est reportée à la section 4.3.

Remarque 4.2. Si on suppose que F est dérivable en z , et G continue en z , alors la relation précédente devient :

$$\limsup_{n \rightarrow \infty} \pm \sqrt{\frac{nh_n}{2 \log_2 n}} (f_n(z) - \mathbb{E}f_n(z)) = \left(\frac{f(z)}{1 - G(z)} \int_{-\infty}^{+\infty} K^2(t) dt \right)^{1/2} \quad p.s$$

On suppose F éventuellement discontinue en z , mais que $f_+(z)$ et $f_-(z)$ existent, alors, sous (H1), (K1) et (K2) on a :

$$\lim_{n \rightarrow \infty} \mathbb{E}f_n(z) = \frac{1}{2}(f_+(z) + f_-(z)) + \frac{1}{2}(f_+(z) - f_-(z)) \int_0^{+\infty} (K(t) - K(-t)) dt,$$

qui se réduit à $f(z)$ si $f_-(z) = f_+(z) = f(z)$.

4.2 Lois fonctionnelles pour les processus basés sur des données censurées

En plus des notations précédentes, on utilisera pour la suite les notations suivantes : Soit H la fonction de répartition de Z

$$H(x) = P(Z \leq x) = 1 - (1 - F(x))(1 - G(x)),$$

elle est décomposée en somme de $H^{(1)}$ et $H^{(0)}$ avec :

$$H^{(1)}(x) = P(Z \leq x, \delta = 1) = \int_0^x (1 - G_-(t)) dF(t), \quad (4.1)$$

$$H^{(0)}(x) = P(Z \leq x, \delta = 0) = \int_0^x (1 - F(t)) dG(t),$$

Posons :

$$p = P(\delta = 1) = \int_0^{+\infty} (1 - G_-(t)) dF(t) = H_-^{(1)}(\infty).$$

Les hypothèses excluent le cas $p = 0$, par contre le cas $p = 1$ est possible et correspond au cas non censuré. Ce dernier a été traité de manière similaire par Deheuvels et Mason (1994), on peut l'exclure sans perte de généralité. On suppose donc que $0 < p < 1$.

Soient les fonctions $Q^{(1)}$ et $Q^{(0)}$ définies par :

$$\begin{aligned} Q^{(1)}(s) &= \inf\{x : H^{(1)}(x) \geq s\} \quad \text{pour } 0 < s < p, \\ Q^{(0)}(s) &= \inf\{x : H^{(0)}(x) \geq s\} \quad \text{pour } 0 < s < 1 - p. \end{aligned}$$

Ce sont les inverses généralisées (ou fonctions quantile) de $H^{(1)}$ et $H^{(0)}$ respectivement. Ces définitions impliquent :

$$\begin{aligned} Q^{(1)}(s) \leq x &\iff s \leq H^{(1)}(x) \quad \text{pour } 0 < s < p \\ Q^{(0)}(s) \leq x &\iff s \leq H^{(0)}(x) \quad \text{pour } 0 < s < 1 - p \end{aligned}$$

Les fonctions de répartition empiriques correspondant à H , $H^{(1)}$, $H^{(0)}$ sont données par

$$\begin{aligned} H_n(x) &= \frac{1}{n} \sum_{i=1}^n 1_{\{Z_i \leq x\}}, \\ H_n^{(1)}(x) &= \frac{1}{n} \sum_{i=1}^n \delta_i 1_{\{Z_i \leq x\}}, \\ H_n^{(0)}(x) &= \frac{1}{n} \sum_{i=1}^n (1 - \delta_i) 1_{\{Z_i \leq x\}}. \end{aligned}$$

Réduction au cas uniforme Comme on l'a fait dans le cas des données complètes (section 1.3), on peut construire une suite de variables aléatoires de loi uniforme sur $[0, 1]$, et telle que la suite des observations (Z_i, δ_i) puisse s'écrire en fonction de cette suite. Notons $Z = Z_1$ et $\delta = \delta_1$.

Lemme 4.1. *Sur un espace de probabilité suffisamment riche¹, on peut définir une variable aléatoire U de loi uniforme sur $[0, 1]$, telle que presque sûrement :*

$$\delta = 1_{\{0 < U < p\}} = 1 - 1_{\{p < U < 1\}};$$

$$Z = \begin{cases} Q^{(1)}(U) & \text{si } 0 < U < p, \\ Q^{(0)}(U - p) & \text{si } p < U < 1. \end{cases}$$

Démonstration. Si les fonctions $H^{(1)}$ et $H^{(0)}$ sont continues, on peut poser

$$U = \delta H^{(1)}(Z) + (1 - \delta)(p + H^{(0)}(Z)).$$

Alors la variable U suit la loi uniforme sur $[0, 1]$ et vérifie la relation du théorème.

Si $H^{(1)}$ et $H^{(0)}$ sont quelconques, on peut remarquer que la fonction de répartition de la variable Z conditionnellement à $\{\delta = 1\}$ est $\frac{1}{p}H^{(1)}$, et que conditionnellement à $\{\delta = 0\}$, la fonction de répartition est $\frac{1}{1-p}H^{(0)}$. Leurs inverses généralisées sont donc $Q^{(1)}(ps)$ et $Q^{(0)}((1-p)s)$ respectivement.

Ainsi, le théorème 1.2 implique que pour une variable aléatoire V de loi uniforme sur $[0, 1]$, la variable définie par $Q^{(1)}(pV)$ (resp. $Q^{(0)}((1-p)V)$) suit la même loi que Z conditionnellement à $\{\delta = 1\}$ (resp. $\{\delta = 0\}$).

L'égalité en distribution plus haut peut être renforcée en égalité presque sûre et on aura alors

$$Z = \delta Q^{(1)}(pV_1) + (1 - \delta)Q^{(0)}((1-p)V_2),$$

pour V_1 et V_2 de loi uniforme sur $[0, 1]$. En posant

$$U = \delta pV_1 + (1 - \delta)(p + (1-p)V_2),$$

on peut facilement voir que U vérifie les relations du théorème. □

D'après ce lemme, on peut supposer que la suite originelle est définie sur un espace de probabilité tel qu'il existe une suite $\{U_n, n \geq 1\}$ de variables aléatoires indépendantes, de même loi uniforme sur $[0, 1]$ et vérifiant pour tout i :

$$\delta_i = 1_{\{0 < U_i < p\}} = 1 - 1_{\{p < U_i < 1\}};$$

1. Ceci peut s'exprimer par l'existence d'une suite de variables aléatoires indépendantes de loi uniforme sur $[0, 1]$ (ou de n'importe quelle autre loi continue).

$$Z = \begin{cases} Q^{(1)}(U_i) & \text{si } 0 < U_i < p, \\ Q^{(0)}(U_i - p) & \text{si } p < U_i < 1. \end{cases}$$

Considérons maintenant la fonction de répartition empirique

$$\mathbb{U}_n(s) = \frac{1}{n} \sum_{i=1}^n 1_{\{U_i \leq s\}},$$

et le processus empirique correspondant

$$\alpha_n(s) = \sqrt{n} (\mathbb{U}_n(s) - s).$$

D'après ce qui précède, on peut réécrire $H_n^{(1)}$ et $H_n^{(0)}$ comme suit :

$$\begin{aligned} H_n^{(1)}(x) &= \mathbb{U}_n(H^{(1)}(x)) && \text{si } 0 < H^{(1)}(x) < p, \\ H_n^{(0)}(x) &= \mathbb{U}_n(H^{(0)}(x) + p) - \mathbb{U}_n(p) && \text{si } 0 < H^{(0)}(x) < 1 - p. \end{aligned}$$

En effet, si $0 < H^{(1)}(x) < p$ on a pour tout x :

$$\begin{aligned} H_n^{(1)}(x) &= \frac{1}{n} \sum_{i=1}^n \delta_i 1_{\{Z_i \leq x\}} \\ &= \frac{1}{n} \sum_{i=1}^n 1_{\{0 < U_i < p, Z_i \leq x\}} \\ &= \frac{1}{n} \sum_{i=1}^n 1_{\{0 < U_i < p, Q^{(1)}(U_i) \leq x\}} \\ &= \frac{1}{n} \sum_{i=1}^n 1_{\{0 < U_i < p, U_i \leq H^{(1)}(x)\}} \\ &= \frac{1}{n} \sum_{i=1}^n 1_{\{0 < U_i \leq H^{(1)}(x)\}} \\ &= \mathbb{U}_n(H^{(1)}(x)). \end{aligned}$$

Et si $0 < H^{(0)}(x) < 1 - p$:

$$\begin{aligned}
 H_n^{(0)}(x) &= \frac{1}{n} \sum_{i=1}^n (1 - \delta_i) 1_{\{Z_i \leq x\}} \\
 &= \frac{1}{n} \sum_{i=1}^n 1_{\{U_i > p, Z_i \leq x\}} \\
 &= \frac{1}{n} \sum_{i=1}^n 1_{\{U_i > p, Q^{(0)}(U_i - p) \leq x\}} \\
 &= \frac{1}{n} \sum_{i=1}^n 1_{\{U_i > p, U_i - p \leq H^{(0)}(x)\}} \\
 &= \frac{1}{n} \sum_{i=1}^n 1_{\{U_i \leq p + H^{(0)}(x)\}} - \frac{1}{n} \sum_{i=1}^n 1_{\{U_i \leq p\}} \\
 &= U_n(p + H^{(0)}(x)) - U_n(p).
 \end{aligned}$$

Soient $s_1 \in]0, p[$, $s_0 \in]p, 1[$ et $M_1 > 0$ fixés, et définissons les suites de fonctions aléatoires $(f_n^{(1)})_{n \geq 1}$ et $(f_n^{(2)})_{n \geq 1}$ par :

$$f_n^{(j)}(u) = \frac{1}{b_n} (\alpha_n(s_j + h_n u) - \alpha_n(s_j)) \quad \text{pour } -M_1 \leq u \leq M_1,$$

où $b_n = \sqrt{2h_n \log_2 n}$, et les suites $(w_n^{(1)})$ et $(w_n^{(0)})$ par :

$$\begin{aligned}
 w_n^{(1)} &= \frac{1}{\sqrt{2 \log_2 n}} \alpha_n(s_1), \\
 w_n^{(0)} &= \frac{1}{\sqrt{2 \log_2 n}} (\alpha_n(s_0) - \alpha_n(p)).
 \end{aligned}$$

Le lemme suivant, qui est un cas particulier du Théorème 1.2 de Deheuvels et Mason (1994), est le point de départ de ce travail.

Lemme 4.2. *Sous (H1) et (H2), la suite $\left\{ (w_n^{(1)}, w_n^{(0)}, f_n^{(1)}, f_n^{(0)}), n \geq 1 \right\}$ est presque sûrement relativement compacte dans $\mathbb{R}^2 \times \mathcal{B}([-M, M])$ avec pour ensemble limite l'ensembles de tous les quadruplets $(w^{(1)}, w^{(0)}, f^{(1)}, f^{(0)})$ vérifiant :*

$$w^{(1)} = \int_0^{s_1} \Psi(s) ds \quad ; \quad w^{(0)} = \int_p^{s_0} \Psi(s) ds$$

$$f^{(j)}(u) = \int_0^u \Psi_j(s) ds \quad ; \quad j = 0, 1 \quad ; \quad -M_1 \leq u \leq M_1;$$

avec

$$\int_0^1 \Psi(s) ds = 0 \quad ; \quad \int_0^1 \Psi^2(s) ds + \int_{-M_1}^{M_1} (\Psi_1^2(s) + \Psi_0^2(s)) ds \leq 1.$$

Soit $M > 0$ fixé, $(\gamma_n^{(1)})$ et $(\gamma_n^{(0)})$ deux suites de fonctions définies sur $[-M, M]$ et telles que :

$$\lim_{n \rightarrow \infty} \sup_{u \in [-M, M]} |\gamma_n^{(j)}(u) - \gamma_j u| = 0 \quad ; \quad j = 0, 1;$$

où γ_1 et γ_0 sont deux constantes positives. Posons pour $j = 0, 1$ et $u \in [-M, M]$

$$g_n^{(j)} = f_n^{(j)}(\gamma_n^{(j)}(u)) = \frac{1}{b_n} (\alpha_n(s_j + h_n \gamma_n^{(j)}(u)) - \alpha_n(s_j)).$$

Nous allons faire usage d'une version modifiée du lemme précédent :

Lemme 4.3. *Sous (H1) et (H2), la suite $(w_n^{(1)}, w_n^{(0)}, g_n^{(1)}, g_n^{(0)})$ est presque sûrement relativement compacte dans $\mathbb{R}^2 \times \mathbb{B}^2([-M, M])$ avec pour ensemble limite l'ensemble de tous les quadruplets $(w^{(1)}, w^{(0)}, g^{(1)}, g^{(0)})$ tels que :*

$$w^{(1)} = \int_0^{s_1} \phi(s) ds \quad ; \quad w^{(0)} = \int_p^{s_0} \phi(s) ds;$$

$$g^{(j)}(u) = \int_0^u \phi_j(s) ds \quad ; \quad j = 0, 1 \quad ; \quad -M \leq u \leq M;$$

avec

$$\int_0^1 \phi(s) ds = 0 \quad ; \quad \int_0^1 \phi^2(s) ds + \int_{-M}^M \frac{1}{\gamma_1} \phi_1^2(s) + \frac{1}{\gamma_0} \phi_1^2(s) ds \leq 1.$$

Démonstration. Soit $M_1 > 0$ vérifiant $|\gamma_n^{(j)}| \leq M_1, \forall u, j$ pour n suffisamment grand.

D'après le Lemme 4.2, de toute sous suite $(w_{n_k}^{(1)}, w_{n_k}^{(0)}, f_{n_k}^{(1)}, f_{n_k}^{(0)})$, on peut extraire une sous suite convergente, c'est-à-dire qu'il existe une sous suite (n'_k) de (n_k) et $(w^{(1)}, w^{(0)}, f^{(1)}, f^{(0)})$ vérifiant :

$$\sum_{j=0,1} \left(\left| w_{n'_k}^{(j)} - w^{(j)} \right| + \sup_{v \in [-M_1, M_1]} \left| f_{n'_k}^{(j)}(v) - f^{(j)}(v) \right| \right) \rightarrow 0,$$

d'où :

$$\sum_{j=0,1} \left(\left| w_{n'_k}^{(j)} - w^{(j)} \right| + \sup_{v \in [-M_1, M_1]} \left| f_{n'_k}^{(j)}(\gamma_{n'_k}^{(j)}(u)) - f^{(j)}(\gamma_{n'_k}^{(j)}(u)) \right| \right) \rightarrow 0.$$

De plus, les fonctions $f^{(j)}, j = 0, 1$ sont uniformément équicontinues. En effet, d'après l'inégalité de Schwarz :

$$\begin{aligned} |f^{(j)}(v') - f^{(j)}(v'')| &= \left| \int_{v'}^{v''} \psi_j(s) ds \right| \\ &\leq |v' - v''|^{1/2} \left| \int_{v'}^{v''} \psi_j^2(s) ds \right|^{1/2} \\ &\leq |v' - v''|^{1/2}, \quad \text{car} \quad \left| \int_{v'}^{v''} \psi_j^2(s) ds \right| \leq 1; \end{aligned}$$

il en résulte que : $\forall v', v'' : -M \leq v', v'' \leq M$

$$\sum_{j=0,1} \left(\left| w_{n'_k}^{(j)} - w^{(j)} \right| + \sup_{u \in [-M, M]} \left| g_{n'_k}^{(j)}(u) - f^{(j)}(\gamma_j u) \right| \right) \xrightarrow[k \rightarrow \infty]{} 0,$$

car :

$$\begin{aligned} |g_n^{(j)}(u) - f^{(j)}(\gamma_j u)| &= |f_n^{(j)}(\gamma_n^{(j)}(u)) - f^{(j)}(\gamma_n^{(j)}(u)) + f^{(j)}(\gamma_n^{(j)}(u)) - f^{(j)}(\gamma_j u)| \\ &\leq |f_n^{(j)}(\gamma_n^{(j)}(u)) - f^{(j)}(\gamma_n^{(j)}(u))| + |f^{(j)}(\gamma_n^{(j)}(u)) - f^{(j)}(\gamma_j u)| \\ &\leq |f_n^{(j)}(\gamma_n^{(j)}(u)) - f^{(j)}(\gamma_n^{(j)}(u))| + |\gamma_n^{(j)}(u) - \gamma_j(u)|^{\frac{1}{2}}. \end{aligned}$$

En posant :

$$\begin{aligned} \phi(u) &= \psi(u) \\ g^{(j)} &= f^{(j)}(\gamma_j u) \\ \phi_j(u) &= \gamma_j \psi_j(\gamma_j u) \quad j = 0, 1; -M < u < M \end{aligned}$$

on obtient :

$$\begin{aligned} \int_0^u \phi_j(s) ds &= \int_0^u \gamma_j \psi_j(\gamma_j s) ds = \int_0^{\gamma_j u} \psi_j(s) ds \\ &= f^{(j)}(\gamma_j u) = g^{(j)}(u); \\ \int_{-M}^M \left(\frac{\phi_1^2(s)}{\gamma_1} + \frac{\phi_0^2(s)}{\gamma_0} \right) ds &= \int_{-M}^M (\gamma_1 \psi_1^2(s) + \gamma_0 \psi_0^2(s)) ds; \\ \int_{-M}^M \left(\frac{\phi_1^2(s)}{\gamma_1} + \frac{\phi_0^2(s)}{\gamma_0} \right) ds &= \int_{-\gamma_1 M}^{\gamma_1 M} \psi_1^2(s) ds + \int_{-\gamma_0 M}^{\gamma_0 M} \psi_0^2(s) ds \\ &\leq \int_{-M_1}^{M_1} (\psi_1^2(s) + \psi_0^2(s)) ds; \end{aligned}$$

$|\gamma_j M| \leq M_1$ car $\lim_{n \rightarrow \infty} \gamma_n^{(j)}(M) = \gamma_j M$ et $|\gamma_n^{(j)}(M)| \leq M_1, \forall j = 0, 1$.

Ainsi, la suite $\left\{ \left(w_n^{(1)}, w_n^{(0)}, g_n^{(1)}, g_n^{(0)} \right), n \geq 1 \right\}$ est presque sûrement relativement compacte et l'ensemble de ses valeurs d'adhérence est inclus dans l'ensemble défini dans le lemme. L'autre inclusion se démontre de la même manière. \square

Considérons $z_1, z_0 \in]0, \Theta[$ fixés et posons :

$$s_1 = H^{(1)}(z_1) \in]0, p[\quad \text{et} \quad s_0 = p + H^{(0)}(z_0) \in]p, 1[.$$

Pour $M > 0$ fixé, on considère les suites $\left(k_n^{(1)} \right)$ et $\left(k_n^{(0)} \right)$ de fonctions aléatoires définies sur $[-M, M]$ par :

$$k_n^{(j)}(u) = \frac{\sqrt{n}}{b_n} \left(H_n^{(j)}(z_j + h_n u) - H^{(j)}(z_j + h_n u) - H_n^{(j)}(z_j) + H^{(j)}(z_j) \right),$$

et les suites de variables aléatoires $(v_n^{(0)})$, $(v_n^{(1)})$ définies par :

$$v_n^{(j)} = \frac{\sqrt{n}}{\sqrt{2 \log_2 n}} (H_n^{(j)}(z_j) - H^{(j)}(z_j)) \text{ pour } j = 0, 1.$$

D'après (H1), il existe $n_0 \geq 1$ tel que pour tout $n \geq n_0$ et tout $j = 0, 1$, $-M \leq u \leq M$, on ait $z_j + h_n u \in]0, \Theta[$. Ceci assure que $(v_n^{(1)}, v_n^{(0)}, k_n^{(1)}, k_n^{(0)}) \in \mathbb{R}^2 \times \mathbb{B}^2([-M, M])$ est bien définie. Le théorème suivant décrit le comportement asymptotique de cette suite.

Théorème 4.3 (Deheuvels et Einmahl, 1996, Théorème 2.1). *On suppose que F est continue et G dérivable en z_0 avec $g(z_0) = G'(z_0) > 0$. On suppose aussi que G est continue et F dérivable en z_1 avec $f(z_1) = F'(z_1) > 0$.*

Alors sous les hypothèses (H1) et (H2), la suite $\left\{ (v_n^{(1)}, v_n^{(0)}, k_n^{(1)}, k_n^{(0)}) : n \geq 1 \right\}$ est presque sûrement relativement compacte dans $\mathbb{R}^2 \times \mathbb{B}^2([-M, M])$ avec pour ensemble limite l'ensemble des quadruplets $(v^{(1)}, v^{(0)}, k^{(1)}, k^{(0)})$ vérifiant :

$$v^{(1)} = \int_0^{H^{(1)}(z_1)} \phi(s) ds \quad ; \quad v^{(0)} = \int_p^{p+H^{(0)}(z_0)} \phi(s) ds;$$

$$k^{(j)}(u) = \int_0^u \phi_j(s) ds;$$

pour $j = 0, 1$ et $-M \leq u \leq M$, avec :

$$\int_0^1 \phi(s) ds = 0;$$

$$\int_0^1 \phi^2(s) ds + \int_{-M}^M \left(\frac{\phi_1^2(s)}{f(z_1)(1-G(z_1))} + \frac{\phi_0^2(s)}{g(z_0)(1-F(z_0))} \right) ds \leq 1.$$

Démonstration. On a :

$$\begin{aligned} k_n^{(1)}(u) &= \frac{1}{b_n} (\sqrt{n} (\mathbb{U}_n (H^{(1)}(z_1 + h_n u)) - H^{(1)}(z_1 + h_n u)) - \sqrt{n} (\mathbb{U}_n (H^{(1)}(z_1)) - H^{(1)}(z_1))) \\ &= \frac{1}{b_n} (\alpha_n (H^{(1)}(z_1 + h_n u)) - \alpha_n (H^{(1)}(z_1))) \\ &= \frac{1}{b_n} \left(\alpha_n \left(H^{(1)}(z_1) + h_n \frac{1}{h_n} (H^{(1)}(z_1 + h_n u) - H^{(1)}(z_1)) \right) - \alpha_n (H^{(1)}(z_1)) \right) \\ &= f_n^{(1)} \left(\frac{1}{h_n} (H^{(1)}(z_1 + h_n u) - H^{(1)}(z_1)) \right); \end{aligned}$$

et de la même manière : $k_n^{(0)}(u) = f_n^{(0)} \left(\frac{1}{h_n} (\mathbb{H}^{(0)}(z_0 + h_n u) - \mathbb{H}^{(0)}(z_0)) \right)$. On a aussi pour $v_n^{(1)}$ et $v_n^{(0)}$:

$$\begin{aligned} v_n^{(1)} &= \frac{1}{\sqrt{2 \log_2 n}} \sqrt{n} (\mathbb{U}_n (\mathbb{H}^{(1)}(z_1)) - \mathbb{H}^{(1)}(z_1)) \\ &= \frac{1}{\sqrt{2 \log_2 n}} \alpha_n (\mathbb{H}^{(1)}(z_1)) = w_n^{(1)}; \end{aligned}$$

et $v_n^{(0)} = w_n^{(0)}$. De plus, si on considère les fonctions :

$$\gamma_n^{(j)}(u) = \frac{1}{h_n} (\mathbb{H}^{(j)}(z_j + h_n u) - \mathbb{H}^{(j)}(z_j)),$$

pour $j = 0, 1$ et $u \in [-M, M]$, alors en vertu de (4.1) :

$$\begin{aligned} &\sup_{-M \leq u \leq M} \left| \mathbb{H}^{(1)}(z_1 + hu) - \mathbb{H}^{(1)}(z_1) - f(z_1)(1 - G(z_1))hu \right| \\ &= \sup_{-M \leq u \leq M} \left| \int_{z_1}^{z_1+hu} (1 - G_-(t)) dF(t) - f(z_1)(1 - G(z_1))hu \right| \\ &\leq \sup_{-M \leq u \leq M} \left| \int_{z_1}^{z_1+hu} (1 - G_-(t)) dF(t) - \int_{z_1}^{z_1+hu} (1 - G(z_1)) dF(t) \right| \\ &\quad + \sup_{-M \leq u \leq M} \left| (1 - G(z_1))(F(z_1 + hu) - F(z_1)) - (1 - G(z_1))f(z_1)hu \right| \\ &\leq \sup_{-M \leq u \leq M} \left| \int_{z_1}^{z_1+hu} (G(z_1) - G_-(t)) dF(t) \right| \\ &\quad + (1 - G(z_1)) \sup_{-M \leq u \leq M} |F(z_1 + hu) - F(z_1) - f(z_1)hu| \\ &\leq \sup_{z_1 - hM \leq t \leq z_1 + hM} |G(z_1) - G_-(t)| \sup_{-M \leq u \leq M} |F(z_1 + hu) - F(z_1)| \\ &\quad + (1 - G(z_1)) \sup_{-M \leq u \leq M} |F(z_1 + hu) - F(z_1) - f(z_1)hu|. \end{aligned}$$

Or, puisque $\sup_{z_1 - hM \leq t \leq z_1 + hM} |G(z_1) - G_-(t)|$ tend vers 0 quand h tend vers 0, $\frac{1}{h} \sup_{-M \leq u \leq M} |F(z_1 + hu) - F(z_1)|$ est bornée, et $F(z_1 + hu) = F(z_1) + f(z_1)hu + o(h)$, on obtient :

$$\frac{1}{h} \sup_{-M \leq u \leq M} \left| \mathbb{H}^{(1)}(z_1 + hu) - \mathbb{H}^{(1)}(z_1) - f(z_1)(1 - G(z_1))hu \right| \longrightarrow 0,$$

quand $h \rightarrow 0$, ce qui montre que

$$\sup_{-M \leq u \leq M} \left| \gamma_n^{(1)}(u) - \gamma_1 u \right| \longrightarrow 0,$$

avec $\gamma_1 = f(z_1)(1 - G(z_1))$. De la même manière, on peut montrer que

$$\sup_{-M \leq u \leq M} |\gamma_n^{(0)}(u) - \gamma_0 u| \rightarrow 0,$$

avec $\gamma_0 = g(z_0)(1 - F(z_0))$, d'où le résultat en appliquant le Lemme 4.3. □

Remarque 4.3. Comme il apparaît dans la preuve ci-dessus, la continuité de F en z_0 et l'existence et la positivité de la dérivée de G en z_0 ne sont nécessaires que pour la compacité relative de la suite $\{k_n^{(0)}\}$. Si on veut uniquement la compacité relative de la suite $\{k_n^{(1)}\}$ dans $\mathbb{B}([-M, M])$ (ce qui est le cas dans la suite de ce travail) avec pour ensemble limite l'ensemble des fonctions $k^{(1)} \in \mathbb{B}([-M, M])$ de la forme

$$k^{(1)}(u) = \int_0^u \phi_1(s) ds \text{ avec } \int_{-M}^M \phi_1^2(s) ds \leq f(z_1)(1 - G(z_1)),$$

alors il suffit que G soit continue en z_1 , et que F soit dérivable en z_1 .

La deuxième partie de ce chapitre consiste à trouver une relation entre ξ_n et $k_n^{(1)}$, ce qui nous permettra de généraliser les résultats du Théorème 4.3 et de la remarque qui le suit au processus ξ_n , objet de notre étude.

Comme nous n'avons plus besoin de $k_n^{(0)}$ pour la suite, nous notons $z = z_1$. Nous supposons aussi, pour simplifier, que F est continue dans un voisinage de z , et dérivable en z , et que G est continue en z .

Rappelons que la martingale de base est définie par :

$$M_n(x) = n \left(H_n^{(1)}(x) - \int_0^x (1 - H_n(s_-)) d\Lambda(s) \right), \quad (4.2)$$

où Λ est la mesure de hasard de X (voir par exemple Gu et Lai, 1990).

$H^{(1)}$ peut être réécrite comme suit :

$$\begin{aligned} H^{(1)}(x) &= \int_0^x (1 - G_-(s))(1 - F(s)) \frac{dF(s)}{1 - F(s)} \\ &= \int_0^x (1 - H_-(s)) d\Lambda(s); \end{aligned}$$

d'où

$$\frac{M_n(x)}{\sqrt{n}} = \sqrt{n} (H_n^{(1)}(x) - H^{(1)}(x)) + \sqrt{n} \int_0^x (H_{n-}(s) - H_-(s)) d\Lambda(s). \quad (4.3)$$

Pour $z \in]0, \Theta[$ fixé, considérons l'accroissement :

$$\mu_n(u) = \frac{1}{b_n \sqrt{n}} (M_n(z + h_n u) - M_n(z)).$$

Lemme 4.4. *On suppose que F est continue dans un voisinage de z , et dérivable en z , et que G est continue en z . Alors, sous les hypothèses (H1) et (H2), on a :*

$$\sup_{u \in [-M, M]} |\mu_n(u) - k_n^{(1)}(u)| = O(\sqrt{h_n}) \rightarrow 0 \text{ p.s.}$$

Démonstration. On a d'après (4.3)

$$\begin{aligned} \mu_n(u) &= \frac{1}{b_n \sqrt{n}} (M_n(z + h_n u) - M_n(z)) \\ &= \frac{1}{b_n} \left(\sqrt{n} (H_n^{(1)}(z + h_n u) - H^{(1)}(z + h_n u)) + \sqrt{n} \int_0^z (H_{n-}^{(1)}(s) - H_-^{(1)}(s)) d\Lambda(s) \right) \\ &\quad - \frac{1}{b_n} \left(\sqrt{n} (H_n^{(1)}(z) - H^{(1)}(z)) + \sqrt{n} \int_0^z (H_{n-}^{(1)}(s) - H_-^{(1)}(s)) d\Lambda(s) \right) \\ &= \frac{\sqrt{n}}{b_n} (H_n^{(1)}(z + h_n u) - H_n^{(1)}(z) - H^{(1)}(z + h_n u) + H^{(1)}(z)) \\ &\quad + \frac{\sqrt{n}}{b_n} \int_z^{z+h_n u} (H_{n-}^{(1)}(s) - H_-^{(1)}(s)) d\Lambda(s) \\ &= k_n^{(1)}(u) + \frac{\sqrt{n}}{b_n} \int_z^{z+h_n u} (H_{n-}^{(1)}(s) - H_-^{(1)}(s)) d\Lambda(s). \end{aligned}$$

En appliquant la loi du logarithme itéré de Chung-Smirnov au processus empirique continu à gauche² $\sqrt{n} (H_{n-}(t) - H_-(t))$ on obtient

$$\lim_{n \rightarrow \infty} \sup_{x \in [-M, M]} \frac{|\sqrt{n}(H_{n-}(t) - H_-(t))|}{\sqrt{2 \log_2 n}} = \frac{1}{2}.$$

Donc, pour n assez grand,

$$\sup_{x \in [-M, M]} |\sqrt{n} (H_{n-}(t) - H_-(t))| \leq \sqrt{\log_2 n},$$

d'où

$$\sup_{u \in [-M, M]} |\mu(u) - k_n^{(1)}(u)| \leq \frac{1}{b_n} \sqrt{\log_2(n)} (\Lambda(z + h_n M) - \Lambda(z - h_n M)),$$

ce qui donne que

$$\sup_{u \in [-M, M]} |\mu_n(u) - k_n^{(1)}(u)| = O(\sqrt{h_n}).$$

□

2. Ceci est équivalent à la définition de la fonction de répartition par $P(X < t)$ et ne change rien aux théorèmes que nous avons vus précédemment.

Lemme 4.5. *Sous les hypothèses du Lemme 4.4, la suite (μ_n) est presque sûrement relativement compacte dans $\mathbb{B}([-M, M])$, avec pour ensemble limite l'ensemble de tous les $\mu \in \mathbb{B}([-M, M])$ satisfaisant :*

$$\mu(u) = \int_0^u \Phi(s) ds \quad -M \leq u \leq M,$$

avec

$$\int_{-M}^M \Phi^2(s) ds \leq f(z)(1 - G(z)).$$

Démonstration. D'après la remarque 4.3, les hypothèses du lemme permettent de conclure que $\{k_n^{(1)}\}$ est presque sûrement relativement compacte dans $\mathcal{B}([-M, M])$ avec pour ensemble limite l'ensemble des éléments $\mu \in \mathcal{B}([-M, M])$ vérifiant :

$$\mu(u) = \int_0^u \phi(s) ds \text{ pour } u \in [-M, M],$$

avec

$$\int_{-M}^M \phi^2(s) ds \leq f(z_1)(1 - G(z_1)).$$

D'après le lemme précédent, il en est de même pour $\{\mu_n\}_{n \geq 1}$. □

Nous allons maintenant utiliser la représentation intégrale suivante du processus empirique en utilisant la martingale de base (4.2) (voir aussi la section 3.4). On a, pour tout $x \leq Z_{(n)} = \max\{Z_1, Z_2, \dots, Z_n\}$

$$\begin{aligned} \Pi_n(x) &= \frac{\sqrt{n}(F_n(x) - F(x))}{1 - F(x)} \\ &= \frac{1}{\sqrt{n}} \int_0^x \frac{dM_n(t)}{1 - H_-(t)} \\ &\quad + \frac{1}{\sqrt{n}} \int_0^x \left(\frac{1 - F_{n-}(t)}{1 - F(t)} \frac{1}{1 - H_{n-}(t)} - \frac{1}{1 - H_-(t)} \right) dM_n(x), \end{aligned}$$

d'après l'équation de Duhamel (voir Gill et Johansen, 1990). On peut alors décomposer $\Pi_n(x)$ en somme $\Pi_{n,1}(x) + \Pi_{n,2}(x)$ avec

$$\Pi_{n,1} = \frac{1}{\sqrt{n}} \int_0^x \frac{dM_n(t)}{1 - H_-(t)},$$

et

$$\Pi_{n,2} = \frac{1}{\sqrt{n}} \int_0^x \left(\frac{1 - F_{n-}(t)}{1 - F(t)} \frac{1}{1 - H_{n-}(t)} - \frac{1}{1 - H_-(t)} \right) dM_n(x).$$

Considérons maintenant l'accroissement :

$$\pi_n(u) = \frac{1}{b_n} (\Pi_n(z + h_n u) - \Pi_n(z)).$$

Lemme 4.6. *Sous les hypothèses du Lemme 4.4, on a :*

$$\sup_{u \in [-M, M]} \left| \pi_n(u) - \frac{\mu_n(u)}{1 - H(z)} \right| \rightarrow 0 \quad p.s$$

Démonstration. Comme $Z_n = \max\{Z_1, \dots, Z_n\} \rightarrow \Theta$ p.s. alors pour tout $\theta < \Theta$, il existe presque sûrement un rang n_θ à partir duquel on a la représentation intégrale ci-dessus pour tout $x \leq \theta$ et tout $n \geq n_\theta$. On peut donc supposer sans perte de généralité que $n \geq n_\theta$ avec $\theta > z + h_n M$. On peut alors poser

$$\begin{aligned} \pi_{n,1}(u) &= \frac{1}{b_n} (\Pi_{n,1}(z + h_n u) - \Pi_{n,1}(z)) \\ &= \frac{1}{b_n \sqrt{n}} \int_z^{z+h_n u} \frac{dM_n(t)}{1 - H_-(t)}; \end{aligned}$$

et

$$\begin{aligned} \pi_{n,2}(u) &= \frac{1}{b_n} (\Pi_{n,2}(z + h_n u) - \Pi_{n,2}(z)) \\ &= \frac{1}{b_n \sqrt{n}} \int_z^{z+h_n u} \left(\left(\frac{1 - F_n(t_-)}{1 - F(t_-)} \right) \frac{1}{1 - H_{n-}(t)} - \frac{1}{1 - H_-(t)} \right) dM_n(t). \end{aligned}$$

En intégrant par parties dans la définition de $\pi_{n,1}$, on obtient (en utilisant le fait que $dM_n(t) = d(M_n(t) - M_n(z))$)

$$\begin{aligned} \pi_{n,1}(u) &= \frac{1}{b_n \sqrt{n}} \frac{M_n(z + h_n u) - M_n(z)}{1 - H(z + h_n u)} \\ &\quad - \frac{1}{b_n \sqrt{n}} \int_z^{z+h_n u} (M_n(t) - M_n(z)) d \left(\frac{1}{1 - H(t)} \right) \\ &= \frac{1}{b_n \sqrt{n}} \frac{M_n(z + h_n u) - M_n(z)}{1 - H(z + h_n u)} \\ &\quad - \frac{1}{b_n \sqrt{n}} \int_0^u (M_n(z + h_n v) - M_n(z)) d \left(\frac{1}{1 - H(z + h_n v)} \right) \\ &= \frac{\mu_n(u)}{1 - H(z + h_n u)} - \int_0^u \mu_n(v) d \left(\frac{1}{1 - H(z + h_n v)} \right); \end{aligned}$$

d'où

$$\begin{aligned} \sup_{u \in [-M, M]} \left| \pi_{n,1}(u) - \frac{\mu_n(u)}{1 - H(z)} \right| &\leq \sup_{u \in [-M, M]} \left| \mu_n(u) \left(\frac{1}{1 - H(z + h_n u)} - \frac{1}{1 - H(z)} \right) \right| \\ &\quad + \sup_{u \in [-M, M]} \left| \int_0^u \mu_n(v) d \left(\frac{1}{1 - H(z + h_n v)} \right) \right|. \end{aligned}$$

Or d'après le lemme précédent, la suite $\{\mu_n\}$ est relativement compacte dans $\mathcal{B}([-M, M])$, donc elle est bornée

$$\exists C > 0, \forall n \geq 1 : \sup_{u \in [-M, M]} |\mu_n(u)| \leq C;$$

ce qui donne

$$\begin{aligned} \sup_{u \in [-M, M]} \left| \pi_{n,1}(u) - \frac{\mu_n(u)}{1 - H(z)} \right| &\leq C \sup_{u \in [-M, M]} \left| \frac{1}{1 - H(z + h_n u)} - \frac{1}{1 - H(z)} \right| \\ &\quad + C \sup_{u \in [-M, M]} \left| \frac{1}{1 - H(z + h_n u)} - \frac{1}{1 - H(z)} \right| \\ &\leq 2C \left(\frac{1}{1 - H(z + h_n M)} - \frac{1}{1 - H(z - h_n M)} \right) \\ &\xrightarrow{n \rightarrow \infty} 0 \quad \text{p.s.} \end{aligned}$$

car H est continue en z .

Il reste à montrer que

$$\sup_{u \in [-M, M]} |\pi_{n,2}(u)| \xrightarrow{n \rightarrow \infty} 0 \quad \text{p.s.}$$

D'autre part, d'après la définition de M_u :

$$\begin{aligned} \pi_{n,2}(u) &= \frac{\sqrt{n}}{b_n} \int_z^{z+h_n u} \left(\left(\frac{1 - F_{n-}(t)}{1 - F(t)} \right) \frac{1}{1 - H_n(t)} - \frac{1}{1 - H(t)} \right) dH_n^{(1)}(t) \\ &\quad + \frac{\sqrt{n}}{b_n} \int_z^{z+h_n u} \left(\frac{1 - F_{n-}(t)}{1 - F(t)} - \frac{1 - H_{n-}(t)}{1 - H_{n-}(t)} \right) d\Lambda(t). \end{aligned}$$

Comme $h_n \rightarrow 0$, et en utilisant la continuité locale de F , on a pour n suffisamment grand $F_{n-}(t) = F(t)$ pour tout $t \in [z - h_n M, z + h_n M]$. D'où, en utilisant la loi du logarithme itéré de Földes et Rejtő (1981)

$$\sup_{t \in [z - h_n M, z + h_n M]} \left| \frac{1 - F_{n-}(t)}{1 - F(t)} - 1 \right| = O \left(\sqrt{\frac{\log_2 n}{n}} \right) \quad \text{p.s.}$$

De la même manière, la loi du logarithme itéré de Chung-Smirnov donne

$$\sup_{t \in [z - h_n M, z + h_n M]} \left| \frac{1 - H_{n-}(t)}{1 - H_{n-}(t)} - 1 \right| = O \left(\sqrt{\frac{\log_2 n}{n}} \right) \quad \text{p.s.}$$

On obtient alors pour un certain $A_n = O(\sqrt{\log_2 n})$ que :

$$\begin{aligned} \sup_{u \in [-M, M]} |\pi_{n,2}(u)| &= \frac{A_n}{b_n} (H_n^{(1)}(z + h_n M) - H_n^{(1)}(z - h_n M)) \\ &\quad + \frac{A_n}{b_n} (\Lambda(z + h_n M) - \Lambda(z - h_n M)). \end{aligned}$$

De plus, sachant que

$$k_n^{(1)}(u) = \frac{\sqrt{n}}{b_n} \left(H_n^{(1)}(z + h_n u) - H^{(1)}(z + h_n u) - H_n^{(1)}(z) + H^{(1)}(z) \right),$$

on trouve

$$\begin{aligned} & \frac{1}{b_n} (H_n^{(1)}(z + h_n M) - H_n^{(1)}(z - h_n M)) \\ &= \frac{1}{\sqrt{n}} (k_n^{(1)}(M) - k_n^{(1)}(-M)) \\ & \quad + \frac{1}{b_n} (H^{(1)}(z + h_n M) - H^{(1)}(z - h_n M)) \\ &= O\left(\frac{1}{\sqrt{n}}\right) + O\left(\sqrt{\frac{h_n}{\log_2 n}}\right), \end{aligned}$$

car la suite $\{k_n^{(1)}(M) - k_n^{(1)}(-M)\}$ est bornée d'après le théorème 4.3. Sachant (d'après l'hypothèse (H2)) que $\frac{nh_n}{\log_2 n} \rightarrow \infty$, on a $\frac{1}{\sqrt{n}} = O\left(\sqrt{\frac{h_n}{\log_2 n}}\right)$ d'où

$$\frac{1}{b_n} (H^{(1)}(z + h_n M) - H^{(1)}(z - h_n M)) = O\left(\sqrt{\frac{h_n}{\log_2 n}}\right).$$

D'autre part, sachant que $\Lambda(s) = -\log(1 - F(s))$ et que $b_n = \sqrt{2h_n \log_2 n}$, la dérivabilité de F assure que

$$\frac{1}{b_n} (\Lambda(z + h_n M) - \Lambda(z - h_n M)) = O\left(\sqrt{\frac{h_n}{\log_2 n}}\right).$$

En substituant ces deux derniers résultats, on obtient

$$\sup_{u \in [-M, M]} |\pi_{n,2}(u)| = O\left(\sqrt{h_n}\right) \rightarrow 0 \quad p.s$$

ce qui termine la preuve du lemme. □

Lemme 4.7. *Sous les hypothèses du Lemme 4.4, on a :*

$$\sup_{u \in [-M, M]} \left| \xi_n(u) - \frac{k_n^{(1)}(u)}{1 - G(z)} \right| \rightarrow 0 \quad p.s$$

Démonstration. On a :

$$\begin{aligned}
 \xi_n(u) &= \frac{1}{b_n} (a_n(z + h_n) - a_n(z)) \\
 &= \frac{\sqrt{n}}{b_n} (F_n(z + h_n u) - F(z + h_n u) - F_n(z) + F(z)) \\
 &= \frac{1}{b_n} (1 - F(z + h_n u)) \frac{\sqrt{n} (F_n(z + h_n u) - F(z + h_n u))}{1 - F(z + h_n u)} \\
 &\quad - \frac{1}{b_n} (1 - F(z)) \frac{\sqrt{n} (F_n(z) - F(z))}{1 - F(z)} \\
 &= \frac{1}{b_n} ((1 - F(z + h_n u)) \Pi_n(z + h_n u) - (1 - F(z)) \Pi_n(z)) \\
 &= \frac{1}{b_n} (1 - F(z)) (\Pi_n(z + h_n u) - \Pi_n(z)) \\
 &\quad - (1 - F(z)) \Pi_n(z + h_n u) + (1 - F(z + h_n u)) \Pi_n(z + h_n u) \\
 &= (1 - F(z)) \pi_n(u) + \frac{1}{b_n} (F(z) - F(z + h_n u)) \Pi_n(z + h_n u).
 \end{aligned}$$

De plus

$$\sup_{u \in [-M, M]} |\Pi(z + h_n u)| = \sup_{u \in [-M, M]} \left| \frac{1}{\sqrt{n}} \int_0^{z+h_n u} \frac{1 - F_{n-}(t)}{1 - F(t)} \frac{1}{1 - H_{n-}(t)} dM_n(t) \right|;$$

et comme on a d'après la loi du logarithme itéré de Földes et Rejtő (1981)

$$\sup_{u \in [z-h_n M, z+h_n M]} \left| \frac{1 - F_{n-}(t)}{1 - F(t)} - 1 \right| = O\left(\sqrt{\frac{\log_2 n}{n}}\right);$$

on en déduit que

$$\sup_{u \in [-M, M]} |\Pi(z + h_n u)| = O\left(\sqrt{\log_2 n}\right).$$

D'autre part, la dérivabilité de F en z et le fait que $h_n \rightarrow 0$ donnent

$$\sup_{u \in [-M, M]} |F(z) - F(z + h_n u)| = O(h_n).$$

Ce qui fait qu'au final on ait

$$\sup_{u \in [-M, M]} |\xi_n(u) - (1 - F(z)) \pi_n(u)| = O(\sqrt{h_n}) \rightarrow 0 \quad \text{p.s}$$

Le résultat du lemme découle enfin des lemmes 4.4 et 4.6 sachant que $\frac{1-F(z)}{1-H(z)} = \frac{1}{1-G(z)}$. \square

Preuve du Théorème 4.1. Si on suppose que F est dérivable en z et G continue en z , alors, sous les hypothèses (H1) et (H2), et par le Théorème 4.3, la suite $\left\{ \frac{k_n^{(1)}}{1-G(z)} : n \geq 1 \right\}$ est presque sûrement relativement compacte avec pour ensemble limite l'ensemble des fonctions $h \in \mathbb{B}([-M, M])$ de la forme :

$$h(u) = \int_0^u \frac{\phi(s)}{1-G(z)} ds \quad -M \leq u \leq M,$$

avec

$$\int_{-M}^M \phi^2(s) ds \leq f(z)(1-G(z)).$$

Le résultat en découle en posant $\psi(s) = \frac{\phi(s)}{1-G(z)}$ et en appliquant le Lemme 4.7. \square

4.3 Application : Estimation de la densité

Le but du Théorème 4.1 est de donner des résultats pour quelques estimateurs qui dépendent localement du processus empirique. L'exemple que nous allons traiter est l'estimateur à noyau de la densité, mais la méthode que nous exposerons peut être utilisée pour les autres estimateurs du même type (par exemple l'estimateur du taux de hasard).

Considérons une fonctionnelle Γ définie et continue sur un ensemble fermé \mathcal{S} de $\mathbb{B}([-M, M])$ et satisfaisant la condition $\xi_n \in \mathcal{S}, \forall n \geq 1$. On définit alors la statistique $T_n = \Gamma(\xi_n)$.

Théorème 4.4 (Deheuvels et Einmahl, 1996, Théorème 3.1). *Soit $z \in]0, \Theta[$ et $M > 0$ fixés. On suppose que F est continue dans un voisinage de z et dérivable en z , et que G est continue en z . Alors, sous (H1) et (H2), la suite $(T_n)_{n \geq 1}$ est presque sûrement relativement compacte dans \mathbb{R} , avec pour ensemble limite l'intervalle*

$$\left[\inf_{h \in \mathbb{L}_M} \Gamma(h), \sup_{h \in \mathbb{L}_M} \Gamma(h) \right],$$

où \mathbb{L}_M est l'ensemble limite de (ξ_n) (voir le théorème 4.1).

Démonstration. Le fait que T_n soit presque sûrement relativement compacte dans \mathbb{R} avec pour ensemble limite $\Gamma(\mathbb{L}_M)$ découle directement du Théorème 4.1.

L'ensemble \mathbb{L}_M est connexe dans $\mathbb{B}([-M, M])$. En effet, l'ensemble des fonctions ψ vérifiant $\int_{-M}^M \psi^2(s) ds \leq \frac{f(z)}{1-G(z)}$ est une boule fermée de $L^2([-M, M])$, elle est donc connexe. L'opérateur qui à toute fonction $\psi \in L^2([-M, M])$ associe la fonction Ψ définie pour $u \in [-M, M]$ par $\Psi(u) = \int_0^u \psi(s) ds$ étant continu, \mathbb{L}_M est aussi connexe.

De plus, si $\Psi \in \mathbb{L}_M$, alors il existe $\psi \in L^2([-M, M])$ tel que $\Psi(u) = \int_0^u \psi(s) ds$, et on a alors, pour tout $x_1, x_2 \in [-M, M]$:

$$\begin{aligned} |\Psi(x_2) - \Psi(x_1)| &= \left| \int_{x_1}^{x_2} \psi(s) ds \right| \\ &\leq |x_2 - x_1|^{1/2} \left(\int_{-M}^M \psi^2(s) ds \right)^{1/2} \\ &\leq \left(\frac{f(z)}{1 - G(z)} \right)^{1/2} |x_2 - x_1|^{1/2}; \end{aligned}$$

ce qui montre que \mathbb{L}_M est uniformément équicontinu. Et comme \mathbb{L}_M est borné (c'est l'image par un opérateur borné d'une boule), il est compact d'après le théorème d'Arzelà-Ascoli.

Son image par l'application continue Γ est donc un intervalle fermé. \square

En choisissant la fonctionnelle Γ , on peut montrer des taux de convergences pour quelques estimateurs qui dépendent localement du processus empirique, la preuve ci dessous est un exemple pour l'estimateur à noyau de la densité.

Preuve du Théorème 4.2. En vertu de l'hypothèse (K2), il existe un réel positif M tel que $\forall u, |u| \geq \frac{M}{2} \implies K(u) = 0$. Donc,

$$\begin{aligned} f_n(z) - \mathbb{E}f_n(z) &= \int_{-\infty}^{+\infty} \frac{1}{h_n} K\left(\frac{t-z}{h_n}\right) dF_n(t) - \int_{-\infty}^{+\infty} \frac{1}{h_n} K\left(\frac{t-z}{h_n}\right) dF(t) \\ &= \frac{1}{h_n} \int_{-\infty}^{+\infty} K(u) d(F_n(z + h_n u) - F(z + h_n u)) \\ &= \frac{1}{h_n} \int_{-M}^M K(u) d(F_n(z + h_n u) - F_n(z) - F(z + h_n u) + F(z)) \quad (4.4) \end{aligned}$$

$$\begin{aligned} &= \frac{-1}{h_n} \int_{-M}^M (F_n(z + h_n u) - F_n(z) - F(z + h_n u) + F(z)) dK(u) \quad (4.5) \\ &= \frac{-b_n}{h_n \sqrt{n}} \int_{-M}^M \xi_n(u) dK(u), \end{aligned}$$

où (4.4) est due au fait que $F_n(z)$ et $F(z)$ sont constants, et (4.5) est obtenue en utilisant l'intégration par parties (ce qui est possible grâce à l'hypothèse (K1)).

Définissons la fonctionnelle Γ en posant pour toute fonction $h \in \mathbb{B}([-M, M])$ à variation bornée :

$$\Gamma(h) = - \int_{-M}^M h(u) dK(u),$$

et, du fait que $b_n = \sqrt{2h_n \log_2 n}$ on a alors :

$$T_n = \Gamma(\xi_n) = \frac{h_n \sqrt{n}}{b_n} (f_n(z) - \mathbb{E}f_n(z)) = \sqrt{\frac{nh_n}{2 \log_2 n}} (f_n(z) - \mathbb{E}f_n(z)).$$

Le Théorème 4.4 s'applique et l'ensemble limite de la suite T_n est l'intervalle

$$\left[\inf_{h \in \mathbb{L}_M} \Gamma(h), \sup_{h \in \mathbb{L}_M} \Gamma(h) \right].$$

Il reste à calculer $\sup_{h \in \mathbb{L}_M} \pm \Gamma(h)$ (car $\inf_x f(x) = -\sup_x -f(x)$). Or par les hypothèses (K1) et (K2), une intégration par parties permet d'écrire

$$\begin{aligned} \sup_{h \in \mathbb{L}_M} \pm \Gamma(h) &= \sup_{h \in \mathbb{L}_M} \pm \int_{-M}^M K(u) dh(u) \\ &= \sup_{h \in \mathbb{L}_M} \left\{ \pm \int_{-M}^M K(u) \Psi(u) du : h(u) = \int_0^u \Psi(s) ds \right\}. \end{aligned}$$

Pour simplifier le calcul de cette expression, on pose :

$$c_- = \frac{1 - G_-(z)}{f_-(z)} ; \quad c_+ = \frac{1 - G(z)}{f_+(z)};$$

$$K^*(t) = \begin{cases} \frac{1}{\sqrt{c_-}} K(t) & \text{si } t < 0, \\ \frac{1}{\sqrt{c_+}} K(t) & \text{sinon;} \end{cases}$$

$$\Psi^*(t) = \begin{cases} \sqrt{c_-} \Psi(t) & \text{si } t < 0, \\ \sqrt{c_+} \Psi(t) & \text{sinon;} \end{cases}$$

ce qui donne que $K(t)\Psi(t) = K^*(t)\Psi^*(t)$ et $h \in \mathbb{L}_M \iff \int_{-M}^M (\Psi^*(u))^2 du \leq 1$. Enfin :

$$\begin{aligned} \sup_{h \in \mathbb{L}_M} \pm \Gamma(h) &= \sup \left\{ \pm \int_{-M}^M K^*(u) \Psi^*(u) du : \int_{-M}^M (\Psi^*(u))^2 du \leq 1 \right\} \\ &\leq \left(\int_{-M}^M (K^*(u))^2 du \right)^{1/2}. \end{aligned}$$

d'après l'inégalité de Schwarz. D'autre part, le choix particulier de

$$\Psi^*(u) = \frac{K^*(u)}{\left(\int_{-M}^M (K^*(t))^2 \right)^{1/2}}$$

montre l'égalité $\sup_{h \in \mathbb{L}_M} \pm \Gamma(h) = \left(\int_{-M}^M (K^*(u))^2 du \right)^{1/2}$. D'où

$$\sup_{h \in \mathbb{L}_M} \pm \Gamma(h) = \left(\frac{f_-(z)}{1 - G_-(z)} \int_{-M}^0 K^2(u) du + \frac{f_+(z)}{1 - G(z)} \int_0^M K^2(u) du \right)^{1/2}.$$

Le résultat visé en découle en rappelant que pour toute suite (x_n) , $\limsup x_n$ n'est autre que la borne supérieure de l'ensemble des valeurs d'adhérence de (x_n) . \square

Chapitre 5

Extension au cas de la censure double

Plusieurs modèles non paramétriques ont été proposés pour l'étude de la censure double. Par exemple, le modèle de Turnbull (1974) est le plus utilisé, et plusieurs travaux sont basés sur ce modèle. Cependant, bien que la définition de ce modèle soit plus intuitive, l'estimateur qui est proposé n'est pas pour autant facile à utiliser, car il est donné par une équation intégrale dont la solution n'est pas connue explicitement.

D'autres modèles (qui englobent parfois la censure par intervalles) ont été proposés par Peto (1973), Samuelsen (1989) ou encore Huang (1999). Dans des rapports techniques, Patilea et Rolin (2001, 2004) discutent les avantages et les inconvénients de ces modèles, et en proposent d'autres.

Dans la suite nous nous intéressons au modèle de Patilea et Rolin (2006) (qui reprend d'ailleurs le modèle de Patilea et Rolin (2001)), et nous discutons de la possibilité d'étendre les résultats du chapitre 4 à ce modèle. Nous terminons par une étude de simulation pour voir la performance de l'estimateur de la densité obtenu pour ce dernier modèle.

5.1 Présentation du modèle de Patilea et Rolin (2006)

Considérons trois variables aléatoires positives indépendantes X , L et R de fonctions de répartition respectives F , F_L et F_R , et de fonctions de survie¹ respectives S , S_L et S_R , où X représente la durée d'intérêt et L et R sont les durées de censure à gauche et à droite respectivement. Dans le modèle I de Patilea et Rolin (2006), au lieu d'observer un échantillon de X on observe un échantillon du couple (Z, A) où $Z = \max(\min(X, R), L)$

1. Si F est la fonction de répartition d'une variable aléatoire X , alors sa fonction de survie est $S = 1 - F$.

et

$$A = \begin{cases} 0 & \text{si } L < X \leq R, \\ 1 & \text{si } L < R < X, \\ 2 & \text{si } \min(X, R) \leq L. \end{cases}$$

Ce modèle considère la censure à droite et la censure à gauche comme deux phénomènes qui agissent indépendamment l'un de l'autre mais que l'un peut censurer l'autre. Un exemple de ce modèle est donné par un système formé par trois composants, dont deux sont placés en série (le composant dont le temps de fonctionnement nous intéresse et un autre). Un troisième est placé en parallèle avec ce système en série.

Le modèle II proposé par les mêmes auteurs est similaire, mais le rôle de la censure à droite et à gauche est inversé. On observe un échantillon du couple (Z, A) où $Z = \min(\max(X, L), R)$ et

$$A = \begin{cases} 0 & \text{si } L < X < R, \\ 1 & \text{si } R < \max(X, L), \\ 2 & \text{si } X \leq L \leq R. \end{cases}$$

Le traitement des deux modèles étant très similaire, nous nous contentons de parler du premier. Considérons H la fonction de répartition de Z , elle se décompose en $\sum_{k=0}^2 H^{(k)}(t)$ où

$$H^{(k)}(t) = P(Z \leq t, A = k), \quad \text{pour } k = 0, 1, 2.$$

Ces fonctions peuvent s'écrire :

$$\begin{aligned} H^{(0)}(t) &= \int_0^t F_{L-}(t) S_{R-}(t) dF(t), \\ H^{(1)}(t) &= \int_0^t F_{L-}(t) S(t) dF_R(t), \\ H^{(2)}(t) &= \int_0^t \{1 - S(t) S_R(t)\} dF_L(t), \end{aligned}$$

et c'est à partir de ces équations que l'estimateur est obtenu.

L'idée est de considérer dans un premier temps $S = \min(X, R)$ et L dans un modèle de censure à gauche (c'est-à-dire que l'on considère une donnée complète si $A = 0$ ou $A = 1$ et censurée à gauche si $A = 2$), et d'estimer la fonction de répartition F_S , puis utiliser la fonction de répartition ainsi estimée au lieu de la fonction de répartition empirique de S pour estimer la fonction de répartition de la variable d'intérêt X en considérant un modèle de censure à droite.

L'estimateur de la fonction de survie S ainsi obtenu, en remplaçant à la fin les fonctions $H^{(0)}$, $H^{(1)}$ et $H^{(2)}$ par leurs estimateurs empiriques, obtenus à partir d'un échantillon

$(Z_i, A_i)_{1 \leq i \leq n}$, est donné par :

$$1 - F_n(Z'_j) = \prod_{1 \leq l \leq j} \left\{ 1 - \frac{D_{0l}}{U_{l-1} - N_{l-1}} \right\},$$

où $(Z'_j)_{1 \leq j \leq M}$ sont les valeurs distinctes des Z_i prises dans l'ordre croissant, et

$$D_{kj} = \sum_{1 \leq i \leq n} 1_{\{Z_i = Z'_j, A_i = k\}},$$

$$N_j = \sum_{1 \leq i \leq n} 1_{\{Z_i \leq Z'_j\}},$$

$$U_{j-1} = n \prod_{j \leq l \leq M} \left\{ 1 - \frac{D_{2l}}{N_l} \right\},$$

pour $0 \leq l \leq 2$ et $1 \leq j \leq M$.

Soulignons le fait que si $L \equiv 0$ (pas de censure à gauche), $1 - F_n$ se réduit à l'estimateur de Kaplan-Meier qui lui-même se réduit au complément à 1 de la fonction de répartition empirique si $R \equiv \infty$.

L'estimateur à noyau de la densité peut être défini en fonction de cet estimateur de la même manière que pour les données complètes ou censurées à droite, autrement dit il s'écrit

$$f_n(z) = \int_{-\infty}^{+\infty} \frac{1}{h_n} K\left(\frac{t-z}{h_n}\right) dF_n(t).$$

5.2 Quelques résultats préliminaires

Une étude attentive du chapitre précédant met en évidence le fait que les principaux outils utilisés par Deheuvels et Einmahl (1996) pour arriver à leurs résultats sont une loi fonctionnelle du logarithme itéré due à Deheuvels et Mason (1994) (Lemme 4.2) ainsi qu'une loi du logarithme itéré pour l'estimateur de Kaplan-Meier (Corollaire 3.1). Le premier est donné dans un cadre bien plus général et semble tout à fait s'appliquer au cas du modèle de censure mixte de Patilea et Rolin (2006). Le second a été étendu par Messaci et Nemouchi (2011) à l'estimateur F_n et s'énonce comme suit. Rappelons que pour toute v.a. U de fonction de répartition F , $I_U = \inf\{t : F(t) > 0\}$ et $T_U = \sup\{t : F(t) < 1\}$.

Théorème 5.1. *On suppose que F , F_L et F_R sont continues et que $I_L < I_X \leq I_R$ et $T_X < T_R$. Alors,*

$$P\left(\sup_{u \in \mathbb{R}} |F_n(u) - F(u)| = O\left(\sqrt{\frac{\log_2 n}{n}}\right)\right) = 1.$$

Ceci montre qu'il est tout à fait envisageable d'étendre les résultats du chapitre précédent à l'estimateur de Patilea et Rolin (2006), travail qui fera l'objet d'une recherche ultérieure. Nous nous contentons dans le paragraphe suivant de mener une étude de simulation afin d'évaluer les performances de l'estimateur à noyau de la densité, obtenu dans le cadre de la censure double.

5.3 Étude de simulation

On va maintenant mener une étude de simulation pour évaluer les performances de l'estimateur à noyau de la densité dans ce cadre.

Pour l'estimation de la densité, il est raisonnable de supposer qu'il n'y pas d'ex-æquo. Ainsi, l'estimateur de Patilea et Rolin (2006) devient :

$$1 - F_n(Z_j) = \prod_{1 \leq k \leq j} \left\{ 1 - \frac{1_{\{A_k=0\}}}{U_{k-1} - k + 1} \right\},$$

avec

$$U_{j-1} = n \prod_{j \leq k \leq n} \left\{ 1 - \frac{1_{\{A_k=2\}}}{k} \right\}.$$

En substituant dans la définition de f_n donnée ci-dessus nous obtenons :

$$f_n(z) = \sum_{1 \leq j \leq n} \frac{1_{\{A_j=0\}} F_n(Z_{j-1})}{U_{j-1} - j + 1} \frac{1}{h_n} K\left(\frac{Z_j - z}{h_n}\right).$$

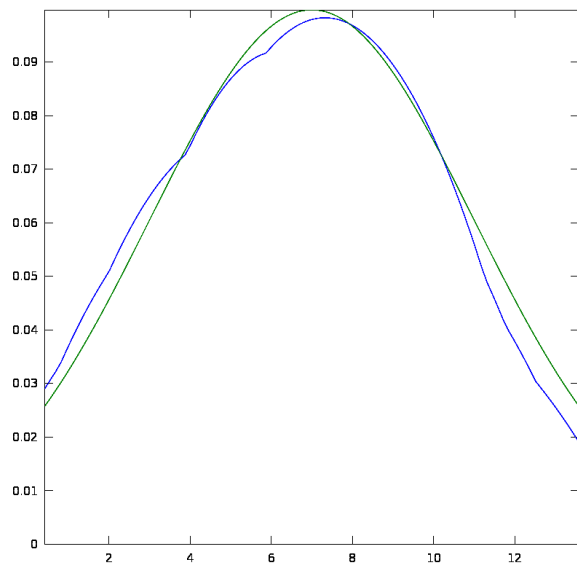
C'est cet estimateur que nous allons étudier pour deux modèles, d'abord un modèle normal puis un modèle Weibull.

Dans le premier modèle, la variable X suit une loi normale de moyenne 7 et d'écart type 4, R suit une loi normale de moyenne 11 et d'écart type 4, et L suit une loi normale de moyenne 3 et d'écart type 4. Le taux de censure n'est pas choisi à l'avance mais est calculé à partir de l'échantillon. Bien que la loi normale ne soit pas positive, ce qui est une condition pour le modèle de Patilea et Rolin (2006), on obtient de bons résultats, ce qui suggère que la condition de positivité n'est pas nécessaire.

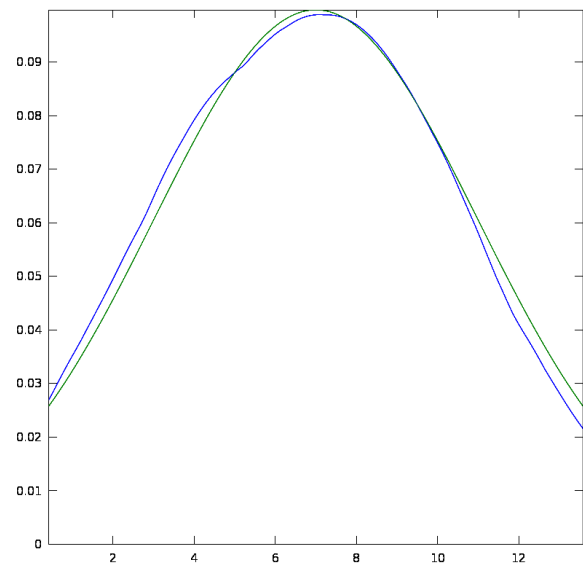
Les graphes suivants ont été réalisés en utilisant le noyau parabolique donné par

$$K(x) = \frac{3}{4}(1 - x^2)1_{\{|x| < 1\}}.$$

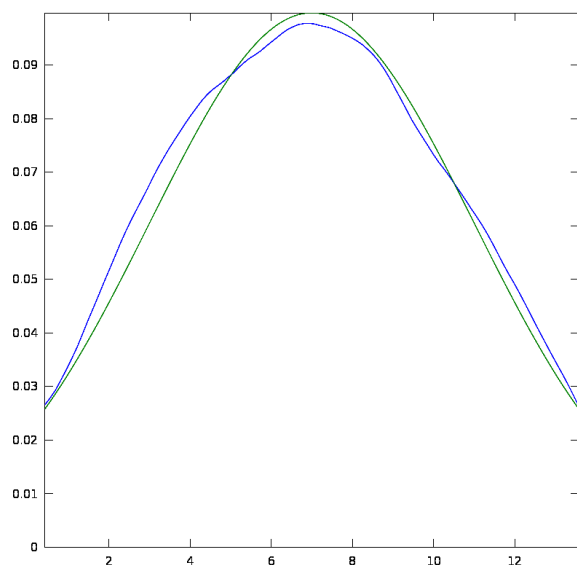
La densité théorique est en vert, et la densité estimée en bleu.



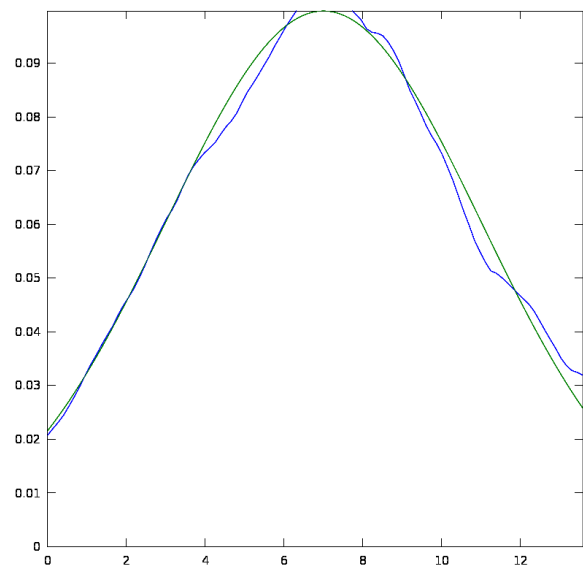
$n = 10$
 20% de censure à droite
 20% de censure à gauche



$n = 50$
 14% de censure à droite
 18% de censure à gauche



$n = 100$
 15% de censure à droite
 23% de censure à gauche

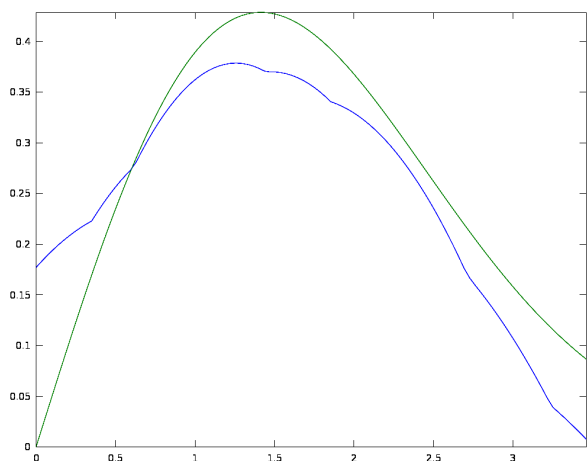


$n = 300$
 21% de censure à droite
 27% de censure à gauche

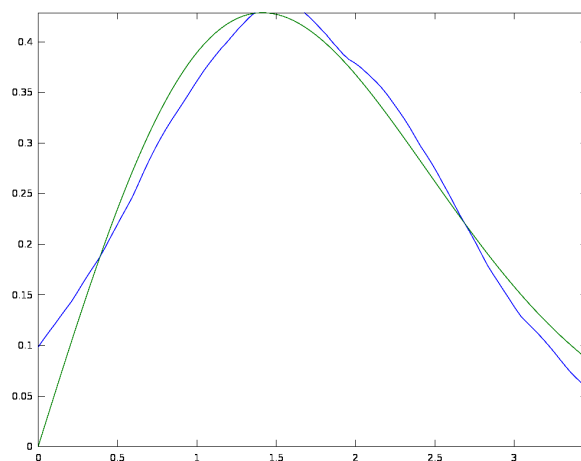
Dans le second modèle, X suit une loi de Weibull de paramètres $\alpha = 2, \beta = 2$, R suit une loi de Weibull de paramètres $\alpha = 3, \beta = 2$, et L suit une loi de Weibull de paramètres $\alpha = 1, \beta = 2$.

La loi de Weibull est définie par sa densité, donnée pour $x > 0$ par :

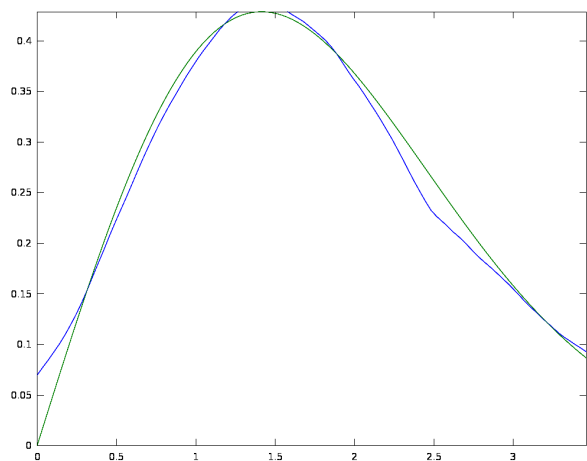
$$\frac{\beta}{\alpha^\beta} x^{\beta-1} \exp\left(-\left(\frac{x}{\alpha}\right)^\beta\right)$$



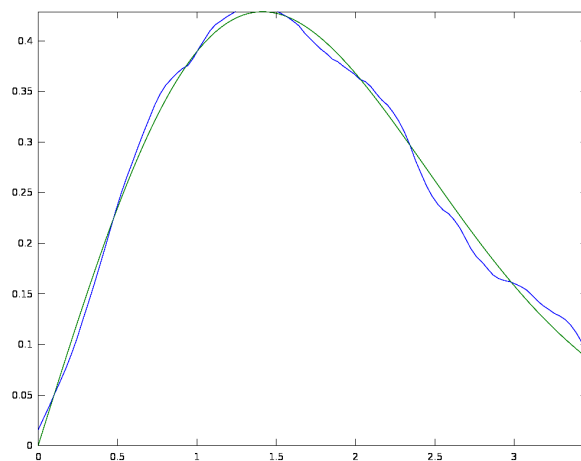
$n = 10$
30% de censure à droite
10% de censure à gauche



$n = 50$
16% de censure à droite
20% de censure à gauche



$n = 100$
21% de censure à droite
20% de censure à gauche



$n = 300$
22% de censure à droite
28% de censure à gauche

Conclusion Les graphes obtenus pour le modèle gaussien montrent une bonne performance de l'estimateur de la densité et ce même pour les petites tailles de l'échantillon. De plus, la qualité de l'estimation s'améliore quand la taille de l'échantillon augmente, ce qui est tout à fait prévisible.

Quant au modèle de Weibull, bien que la qualité de l'estimation ne soit pas bonne pour $n = 10$, les graphes montrent une amélioration tout à fait satisfaisante quand la taille de l'échantillon augmente.

Notons aussi que le taux de censure se situe autour de 40%, ce qui nous semble appréciable.

Bibliographie

Odd AALEN : *Statistical theory for a family of counting processes*. Institute of Mathematical Statistics, Univ. of Copenhagen, 1976.

Patrick BILLINGSLEY : *Convergence of Probability Measures*. Wiley, New York, 1968.

Patrick BILLINGSLEY : *Weak Convergence of Measures : Applications in Probability*. CBMS-NSF Regional Conference Series in Applied Mathematics. SIAM, Philadelphia, Pennsylvania, 1971.

Patrick BILLINGSLEY : *Probability and Measure*. John Wiley & Sons, 2e édition, 1986.

N. BRESLOW et J. CROWLEY : A large sample study of the life table and product limit estimates under random censorship. *The Annals of Statistics*, 2(3):437–453, 1974.

F. P. CANTELLI : Sulla determinazione empirica delle leggi di probabilità. *Giornale dell'Istituto Italiano degli Attuari*, 4:421–424, 1933.

D. M. CHIBISOV : Some theorems on the limiting behavior of empirical distribution functions. *Selected Transl. Math. Statist. Prob.*, 6:147–156, 1964.

Kai-Lai CHUNG : An estimate concerning the Kolmogoroff limit distribution. *Transactions of the American Mathematical Society*, 67(1):36–50, September 1949.

P. CSÖRGŐ, M. et Révész : *Strong Approximation in Probability and Statistics*. Academic Press, New York, 1981.

Alejandro de ACOSTA : A new proof of the Hartman-Wintner law of the iterated logarithm. *The Annals of Probability*, 11(2):270–276, May 1983.

Paul DEHEUVELS : Conditions nécessaires et suffisantes de convergence ponctuelle presque sûre des estimateurs de la densité. *C. R. Acad. Sci., Paris, Sér. A*, 278:1217–1220, 1974.

Paul DEHEUVELS : Laws of the iterated logarithm for density estimators. In G. ROUSSAS, éditeur : *Nonparametric Functional estimators and Related Topics*, NATO adv.Sci.Ins. C, pages 19–29. Kluwer Academic, 1991.

- Paul DEHEUVELS : Functional laws of the iterated logarithm for large increments of empirical and quantile processes. *Stochastic Processes Appl.*, 43(1):133–163, 1992.
- Paul DEHEUVELS : Chung type functional laws of the iterated logarithm for tail empirical processes. *Ann. I. H. P. B.*, 36:583–616, 2000.
- Paul DEHEUVELS et John H. J. EINMAHL : On the strong limiting behavior of local functionals of empirical processes based upon censored data. *The Annals of Probability*, 24(1):504–525, 1996.
- Paul DEHEUVELS et David MASON : Functional laws of the iterated logarithm for the increments of empirical and quantile processes. *Ann. Prob.*, 20:1248–1287, 1992.
- Paul DEHEUVELS et David MASON : General confidence bounds for nonparametric functional estimators. *Stat. Inf. for Stoch. Proc.*, 7:225–277, 2004.
- Paul DEHEUVELS et David M. MASON : Functional laws of the iterated logarithm for local empirical processes indexed by sets. *The Annals of Probability*, 22(3):1619–1661, 1994.
- M. D. DONSKER : Justification and extension of Doob’s heuristic approach to the Kolmogorov-Smirnov theorems. *The Annals of Mathematical Statistics*, 23:277–281, 1952.
- R. M. DUDLEY : Weak convergence of probabilities on nonseparable metric spaces and empirical measures on Euclidean spaces. *Illinois Journal of Mathematics*, 10:109–126, 1966.
- R. M. DUDLEY : measures on non-separable metric spaces. *Illinois Journal of Mathematics*, 11(3):449–453, 1967.
- R. M. DUDLEY : Distances of probability measures and random variables. *The Annals of Mathematical Statistics*, 39:1563–1572, 1968.
- R. M. DUDLEY : *Probabilities and metrics : Convergence of laws on metric spaces, with a view to statistical testing*, volume 45 de *Lecture Notes Series*. Matematisk Institut, Aarhus Universitet, 1976.
- R. M. DUDLEY : *Uniform central limit theorems*. Cambridge University Press, 1999.
- John H. J. EINMAHL et A. J. KONING : Limit theorems for a general weighted process under random censoring. *Canadian Journal of Statistics*, 20, 1992.
- Uwe EINMAHL et David MASON : An empirical process approach to the uniform consistency of kernel type estimators. *Journ. Theoretic. Probab.*, 13:1–13, 2000.
- Uwe EINMAHL et David M. MASON : Uniform in bandwidth consistency of kernel-type function estimators. *The Annals of Statistics*, 33(3):1380–1403, 2005.

- Pedro J. FERNANDEZ : Almost surely convergent versions of sequences which converge weakly. *Bulletin of the Brazilian Mathematical Society*, 5:51–61, 1974.
- Helen FINKELSTEIN : The law of the iterated logarithm for empirical distribution. *The Annals of Mathematical Statistics*, 42(2):607–615, 1971.
- Thomas R. FLEMING et David P. HARRINGTON : *Counting Processes and Survival Analysis*. Wiley, New York, 1991.
- A. FÖLDES et L. REJTÓ : A LIL type result for the product limit estimator. *Probability Theory and Related Fields*, 56(1):75–86, 1981.
- Richard GILL : Large sample behaviour of the product-limit estimator on the whole line. *The Annals of Statistics*, 11(1):49–58, 1983.
- Richard D. GILL et Soren JOHANSEN : A survey of product-integration with a view toward application in survival analysis. *The Annals of Statistics*, 18(4):1501–1555, 1990.
- V. GLIVENKO : Sulla determinazione empirica delle leggi di probabilità. *Giornale dell'Istituto Italiano degli Attuari*, 4:92–99, 1933.
- Ming Gao GU et Tze Leung LAI : Functional laws of the iterated logarithm for the product limit-estimator of a distribution function under random censorship or truncation. *The Annals of Probability*, 18(1):160–189, 1990.
- Philip HARTMAN et Aurel WINTNER : On the law of the iterated logarithm. *American Journal of Mathematics*, 63(1):169–176, January 1941.
- J. HUANG : Asymptotic properties of nonparametric estimation based on partly interval-censored data. *Statistica Sinica*, 9:501–519, 1999.
- B. R. JAMES : A functional law of the iterated logarithm for weighted empirical distributions. *The Annals of Probability*, 3:762–772, 1975.
- E. L. KAPLAN et P. MEIER : Nonparametric estimation from incomplete observations. *Journal of the American Statistical Association*, 53:457–481, 1958.
- A. KHINCHINE : Über einen Satz der Wahrscheinlichkeitsrechnung. *Fundamenta Mathematica*, 6:9–20, 1924.
- A. KOLMOGOROV : Über das Gesetz des iterierten Logarithmus. *Mathematische Annalen*, 101:126–135, 1929.
- A. N. KOLMOGOROV : Sulla determinazione empirica di una legge di distribuzione. *Giornale dell'Istituto Italiano degli Attuari*, 4:83–91, 1933.
- Michael R. KOSOROK : *Introduction to Empirical Processes and Semiparametric Inference*. Springer Verlag, New York, 2008.

- T. L. LAI : Convergence rates in the strong law of large numbers for random variables taking values in Banach spaces. *Bull. Inst. Math. Acad. Sinica*, 2:67–85, 1974.
- Fatiha MESSACI et Nahima NEMOUCHI : A law of the iterated logarithm for the product limit estimator with doubly censored data. *Statistics & Probability Letters*, 81(8):1241–1244, 2011.
- N. E. O'REILLY : On the weak convergence of empirical processes in sup-norm metrics. *The Annals of Probability*, 2:642–651, 1974.
- Emanuel PARZEN : On estimation of a probability density function and mode. *The Annals of Mathematical Statistics*, 33:1065–1076, 1962.
- Valentin PATILEA et Jean-Marie ROLIN : Product limit estimators of the survival function for doubly censored data. Rapport technique, Institut de Statistique, Université Catholique de Louvain, July 2001.
- Valentin PATILEA et Jean-Marie ROLIN : Product-limit estimators of the survival function with left or right censored data. Rapport technique, Institut de Statistique, Université Catholique de Louvain, 2004.
- Valentin PATILEA et Jean-Marie ROLIN : Product limit estimators of the survival function with twice censored data. *The Annals of Statistics*, 34(2):925–938, 2006.
- Richard PETO : Experimental survival curves for interval censored data. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 22(1):86–91, 1973.
- D. POLLARD : *Convergence of stochastic processes*. Springer-Verlag, New York, 1984.
- Yu. V. PROKHOROV : Convergence of random processes and limit theorems in probability theory. *Theory of Probability and its Applications*, 1(2):157–214, 1956.
- Murray ROSENBLATT : Remarks on some nonparametric estimates of density function. *The Annals of Mathematical Statistics*, 27:832–837, 1956.
- S. O. SAMUELSEN : Asymptotic theory for non-parametric estimators from doubly censored data. *Scandinavian Journal of Statistics*, 16:1–21, 1989.
- Galen R. SHORACK et Jon W. WELLNER : *Empirical processes with applications to statistics*. John Wiley & Sons, 1986.
- A. V. SKOROKHOD : Limit theorems for stochastic processes. *Theory of Probability and its Applications*, 1(3):261–290, 1956.
- N. SMIRNOV : Sur les écarts de la courbe de distribution empirique. *Recueil Mathématique [Matematicheskii Sbornik]*, 6(48)(1):3–26, 1939.

- W. STUTE : The law of the iterated logarithm for kernel density estimators. *Ann. Prob.*, 10:414–422, 1982a.
- W. STUTE : The oscillation behavior of empirical processes. *Ann. Prob.*, 10:86–107, 1982b.
- W. STUTE : Conditional empirical processes. *Ann. Prob.*, 14(2):638–647, 1986a.
- W. STUTE : On almost sure convergence of conditional empirical distribution functions. *Ann. Prob.*, 14(3):891–901, 1986b.
- W. STUTE et J.-L. WANG : The strong law under random censorship. *The Annals of Statistics*, 21(3):1591–1607, 1993.
- Bruce W. TURNBULL : Nonparametric estimation of a survivorship function with doubly censored data. *Journal of the American Statistical Association*, 69(345):169–173, March 1974.
- Bruce W. TURNBULL : The empirical distribution function with arbitrarily grouped, censored and truncated data. *Journal of the Royal Statistical Society. Series B (Methodological)*, 38(3):290–295, 1976.
- A. W. van der VAART et Jon A. WELLNER : *Weak Convergence and Empirical Processes : With Applications in Statistics*. Springer Verlag, New York, 1996.
- Vivian VIALON : *Processus empiriques, estimation non paramétrique et données censurées*. Thèse de doctorat, Université Paris 6, 2006.
- Jon A. WELLNER : A Glivenko-Cantelli theorem and strong law of large numbers for functions of order statistics. *The Annals of Statistics*, pages 473–780, 1977.
- M. J. WICHURA : On the construction of almost uniformly convergent random variables with given weakly convergent image laws. *The Annals of Mathematical Statistics*, 41:284–291, 1970.
- B. B. WINTER, A. FÖLDES et L. REJTŐ : Glivenko-Cantelli theorems for the product limit estimate. *Problems of Control and Information Theory*, 7:213–225, 1978.

ملخص

في هذه الرسالة، نقدم بعض النتائج الأساسية لنظرية العمليات التجريبية، كما نقدم بعض النتائج المماثلة في حالة المعطيات المحجوبة¹. نهتمّ بشكل خاص بنتيجة (Deheuvels and Einmahl 1996)، وهو قانون دالي للوغارتم المكرر بالنسبة لتزايدات العملية التجريبية في حالة المعطيات المحجوبة من اليمين. و في الأخير، نناقش آفاق البحث في تعميم هذه النتيجة إلى قضية الحجب المزدوج.

إذا أخذنا عينة $(X_i)_{1 \leq i \leq n}$ ذات دالة توزيع F ، وإذا كانت F_n دالة التوزيع التجريبية فإن العملية التجريبية تُعرف بـ $a_n(t) = \sqrt{n}(F_n(t) - F(t))$. نقدم أولاً بعض النتائج النظرية للعمليات التجريبية، وقانون اللوغارتم المكرر لهذه العملية. ثم نركّز على عملية تجريبية في حالة المعطيات المحجوبة. في هذا النموذج، لا يمكن أن نشاهد بدقة المعطيات التي نريد، ولكن نلاحظ عوضاً عنها كلاً من $Z_i = \min(X_i, C_i)$ و $\delta_i = 1_{\{X_i \leq C_i\}}$ حيث $(C_i)_{1 \leq i \leq n}$ هي الأزمنة الحاجبة. ويكفي في هذه الحال أن نعوض دالة التوزيع التجريبية بمقدّر Kaplan-Meier. نقدم بعض النتائج الأولية لهذا النموذج، ونهتمّ بشكل خاص بقانون دالي للوغارتم المكرر لتزايدات عملية Kaplan-Meier، أي بالعملية $\xi_n(u) = \frac{1}{b_n}(a_n(z + hu) - a_n(z))$ حيث $b_n = \sqrt{2h \log \log n}$ ، ونقتصر في ذلك على الحالة أين يكون $h = h_n$ متتالية لا تتعلق إلا بـ n . علينا أيضاً أن ننظر في تطبيق هذه النتيجة لإعطاء سرعة تقارب مقدر نواة من الكثافة. أخيراً، يتم إجراء دراسة محاكاة لاستكشاف خصائص امتدادها لهذا مقدر لحالة بيانات رقابة مزدوجة، على أساس نموذج

Patilea and Rolin (2006)، لتوجيه الحدس حول إمكانية تمديد نتيجة Deheuvels-Einmahl لهذا النموذج.

العبارات الرئيسية عملية تجريبية؛ معطيات محجوبة؛ قانون دالي للوغارتم المكرر؛ حجب مزدوج

¹آثرت استعمال هذا المصطلح بدلاً من "مراقبة" الذي يستعمله البعض لأن هذا الأخير لا يؤدي المعنى الإحصائي لكلمة "censoring".

Abstract

In this work, we present some basic results of the theory of empirical processes, and some analogous results in the case of censored data. We are particularly interested in the result of Deheuvels and Einmahl (1996), which is a functional law of the iterated logarithm for the increments of the empirical process in the setting of right-censored data. Finally, we discuss the possibility of extending this result to the setting of twice censoring.

Given a sample $(X_i)_{1 \leq i \leq n}$ with distribution function F , the empirical process is defined by $a_n(t) = \sqrt{n}(F_n(t) - F(t))$ where F_n is the empirical distribution function. First, we present some results of the theory of empirical processes, and some laws of iterated logarithm for this process. Then we look at the empirical process in the case of right censored data. In this model, we cannot observe the data of interest, but rather $Z_i = \min(X_i, C_i)$ and $\delta_i = 1_{\{X_i \leq C_i\}}$ where $(C_i)_{1 \leq i \leq n}$ are censoring time. It is enough to replace the empirical distribution function by the Kaplan-Meier estimator. For this model, we present some preliminary results, and we are particularly interested in a functional law of iterated logarithm for the increments of the Kaplan-Meier empirical process, that is, for $\xi_n(u) = \frac{1}{b_n}(a_n(z + hu) - a_n(z))$ with $b_n = \sqrt{2h \log \log n}$. We consider only the case where $h = h_n$ is a deterministic sequence which only depends on n , we also consider the application of this result to give a rate of convergence for the kernel density estimator. Finally, a simulation study is conducted to explore the properties of this estimator in the setting of twice censored data, based on the model of Patilea and Rolin (2006), with the goal of investigating the possibility of extending the result of Deheuvels and Einmahl (1996) to this model.

Keywords Empirical process; Censored data; Functional law of the iterated logarithm; Twice censoring

Résumé

Dans ce mémoire, nous exposons certains résultats de base de la théorie des processus empiriques, ainsi que quelques résultats analogues dans le cas des données censurées. Nous nous intéressons particulièrement au résultat de Deheuvels et Einmahl (1996), qui est une loi fonctionnelle du logarithme itéré pour les accroissements du processus empirique dans le cadre des données censurées à droite. Enfin, nous abordons les perspectives de recherche dans la généralisation de ce résultat au cas de la censure double.

Étant donné un échantillon $(X_i)_{1 \leq i \leq n}$ de fonction de répartition F , le processus empirique est défini par $a_n(t) = \sqrt{n}(F_n(t) - F(t))$ où F_n est la fonction de répartition empirique. Nous exposons dans un premier temps les résultats de la théorie des processus empiriques, ainsi que des lois du logarithme itéré pour ce processus. Ensuite, nous nous intéressons au processus empirique dans le cas de données censurées à droite. Dans ce modèle, on n'observe pas les données d'intérêt, mais plutôt $Z_i = \min(X_i, C_i)$ et $\delta_i = 1_{\{X_i \leq C_i\}}$ où $(C_i)_{1 \leq i \leq n}$ sont des temps de censure, et il suffit alors de remplacer la fonction de répartition empirique par l'estimateur de Kaplan-Meier. Pour ce modèle, nous exposons quelques résultats préliminaires, et nous nous intéressons particulièrement à une loi fonctionnelle du logarithme itéré pour les accroissements du processus de Kaplan-Meier, c'est-à-dire pour $\xi_n(u) = \frac{1}{b_n}(a_n(z + hu) - a_n(z))$ avec $b_n = \sqrt{2h \log \log n}$, nous considérons uniquement le cas où $h = h_n$ est une suite déterministe qui ne dépend que de n . Nous considérons aussi l'application de ce résultat pour donner une vitesse de convergence de l'estimateur à noyau de la densité. Enfin, une étude de simulation est menée pour explorer les propriétés d'une extension de cet estimateur au cas des données doublement censurées, d'après le modèle de Patilea et Rolin (2006), afin de guider l'intuition sur la possibilité d'étendre le résultat de Deheuvels et Einmahl (1996) à ce modèle.

Mots-clés Processus empirique ; Données censurées ; Loi fonctionnelle du logarithme itéré ; Censure double