

République Algérienne Démocratique et Populaire
Ministère de l'Enseignement Supérieur et de la Recherche Scientifique
UNIVERSITÉ KASDI MERBAH OUARGLA

FACULTÉ DE MATHÉMATIQUES ET SCIENCES DE LA MATIÈRE

DÉPARTEMENT DE MATHÉMATIQUES



Mémoire présenté En Vue De L'obtention Du

DIPLÔME DE MASTER

EN MATHÉMATIQUES

Option : Probabilité et Statistique

Par

Abir HALASSA

Intitulé

L'estimation des indices extrêmes conditionnelles dans le cas des données censurées

Membres du jury

Mohamed AGTI	M. A. A	UKMO	Président
Abdelmalek BOUSSAAD	M. A. A	UKMO	Examineur
Fatima MEDDI	M. A. A	UKMO	Rapporteur

septembre 2020

Dédicace

Je dédie ce modeste travail

*A ma très cher Mère et mon très cher Père
A mes chers sœurs et frères
A toutes la familles HLASSA*

*Sans oublier de dédier ce mémoire à mes très chères amies intimes
A tous mes collègues de ma promotion de Proba-Stat 2019*

REMERCIEMENTS

Je remercie Dieu le tout-puissant de m'avoir donné la volonté, la force, et le courage patience et la santé pour bien mener et finir mon travail de mémoire.

Je voudrais d'abord et avant tout remercier ma encadreur vertueuse MEDDI fatima pour efforts en vue d'établir ce mémoire, elle a eu le rôle fondamental et essentiel et la grand mérite dans tout ce qui a été réalisé, me guider et corriger mes erreurs et me donner de précieux conseils et les conseils appropriés et les alertes considérés, je la répète mes remerciements pour tout ce qu'elle a fait pour moi à travers la mise en plan à toutes les étapes de la préparation de ce mémoire depuis le début jusqu'à la fin, étape par étape, était que le premier et dernier facteur pour le succès de ce travail, et je lui dis encore une fois merci beaucoup pour votre appréciation profonde.

Avec un grand honneur, j'aimerais présenter mes remerciements et ma gratitude aux membres du jury, Monsieur

Abdelmalek BOUSAAD, Mohamed AGTI , et Monsieur Said Zibar tout d'abord d'avoir accepté d'examiner mon mémoire, qui sans eux ce mémoire ne pourra jamais voir le jour pour l'intérêt et apport qu'ils ont apporté à mon travail.

J'exprime ma gratitude à ma famille qui m'a toujours soutenue et encouragée dans la voie que je m'étais fixée. Je remercie particulièrement mes parents qui m'ont stimulée et encouragé pendant mes études. qui étaient toujours prêts à fournir tous les moyens physique et morales pour la réussite de ce travail.

Table des matières

0.1	Abréviations et Notations	9
0.2	Introduction	11
1	Rappels sur la théorie des valeurs extrêmes et sur la censure	14
1.1	Comportement du maximum d'un échantillon	14
1.1.1	Densités et loi des statistiques d'ordre	14
1.1.2	Comportement asymptotique des extrêmes	16
1.2	Présentation de la théorie des valeurs extrêmes	17
1.2.1	Loi des valeurs extrêmes	17
1.2.2	Caractérisation des Domaines d'attraction	18
1.2.3	Estimation de l'indice des valeurs extrêmes γ	23
1.3	Quelques rappels sur la modélisation des valeurs extrêmes conditionnelles .	28
1.3.1	Fonction de répartition conditionnelle	28
1.3.2	Estimation de la fonction répartition conditionnelle	28
1.3.3	Estimation des quantiles conditionnels	30
1.4	Quelques généralités sur la censure	30
1.4.1	données de survie	30
1.4.2	données censurées	31
1.4.3	Estimation de la fonction de survie	33
2	Estimation de l'indice des valeurs extrêmes conditionnelles en présence de données Complètes et de données censurées	35
2.1	Estimation de l'indice des valeurs extrêmes conditionnelles en design fixe en présence de données Complètes	35
2.1.1	Définitions des estimateurs	36
2.1.2	Propriétés asymptotiques de l'estimateur de l'indice	36
2.2	Estimation de l'indice des valeurs extrêmes conditionnelles en design aléatoire en présence de données Complètes	37
2.2.1	Définitions des estimateurs	38
2.2.2	Propriétés asymptotiques de l'estimateur de l'indice	38
2.3	Estimation de l'indice des valeurs extrêmes conditionnelles en présence de données censurées	39
2.4	Construction d'un nouvel estimateur de l'indice des valeurs extrêmes conditionnelles	40
2.4.1	Sélection de l'échantillon conditionnel	41
2.4.2	Estimateur pour des données complètes	42
2.4.3	Estimateur pour des données censuré	43

3	Etude de simulation	45
3.1	introduction du logiciel R	45
3.1.1	Qu'est-ce que le logiciel R?	45
3.1.2	les avantages de logiciel R	45
3.2	Princip du Bootstrap	45
3.3	Simulation pour des données complètes	46
3.3.1	Échantillon initial et paramètres de simulation	46
3.3.2	Sélectionnement de l'échantillon conditionnel	46
3.3.3	Simulation de nouvel estimateur γ^h en fonction k	50
3.3.4	En plait de bootsrap pour de $\hat{\gamma}$	57
3.4	Simulation pour des données censurées	60
3.4.1	Échantillon initial et paramètres de simulation	60
3.4.2	Sélectionnement de l'échantillon conditionnel censure	60
3.4.3	Simulation de nouvel estimateur de l'indice en fonction de k	65
3.4.4	emploi du bootsrap	72

Table des figures

1.1	Densité de $\Phi_\gamma(x), \Psi_\gamma(x), \Lambda(x)$	18
1.2	Représentation de la densité et fonction de répartition : Gumbel ($\gamma = 0$), Fréchet ($\gamma = 1$) et Weibull ($\gamma = -1$)	21
1.3	Comportement asymptotique de l'estimateur de Pickands avec des intervalles de confiance au niveau de confiance 0.95 sur les données d'assurance incendie Danoise pour la période allant de 01/01/1980 au 31/12/1990	24
1.4	Comportement asymptotique de l'estimateur de Hill avec des intervalles de confiance au niveau de confiance 0.95 sur les données d'assurance incendie Danoise pour la période allant de 01/01/1980 au 31/12/1990.	26
1.5	Comportement asymptotique de l'estimateur des moments avec des intervalles de confiance au niveau de confiance 0.95 sur les données d'assurance incendie Danoise pour la période allant de 01/01/1980 au 31/12/1990	27
1.6	Comportement asymptotique de l'estimateur UH avec des intervalles de confiance au niveau de confiance 0.95 sur les données d'assurance incendie Danoise pour la période allant de 01/01/1980 au 31/12/1990.	28
2.1	Les différentes étapes de sélection des données. En ordonnée on a la variable Y et en abscisse la covariable X.	41
3.1	Echantillon conditionnel Z pour ($\gamma_1 = 0.2$) et ($\gamma_2 = 0.8$)	48
3.2	Échantillon conditionnel Z pour ($\gamma_1 = 0.5$) et ($\gamma_2 = 0.8$)	48
3.3	Échantillon conditionnel Z pour ($\gamma_1 = 1$) et ($\gamma_2 = 0.8$)	49
3.4	Échantillon conditionnel Z pour ($\gamma_1 = 1.5$) et ($\gamma_2 = 0.8$)	49
3.5	Comportement graphique de $\hat{\gamma}$ vs k ,pour n=200 et beta =2	52
3.6	Comportement graphique de $\hat{\gamma}$ vs k ,pour n=200 et beta =5	52
3.7	Comportement graphique de $\hat{\gamma}$ vs k ,pour n=200 et beta 10	53
3.8	Comportement graphique de $\hat{\gamma}$ vs k ,pour n=1000 et beta =2	53
3.9	Comportement graphique de $\hat{\gamma}$ vs k ,pour n=1000 et beta =5	54
3.10	Comportement graphique de $\hat{\gamma}$ vs k ,pour n=1000 et beta =10	54
3.11	Comportement graphique de $\hat{\gamma}$ vs k ,pour n=5000 et beta =2	55
3.12	Comportement graphique de $\hat{\gamma}$ vs k ,pour n=5000 et beta =5	55
3.13	Comportement graphique de $\hat{\gamma}$ vs k ,pour n=5000 et beta =10	56
3.14	QQ-norm de la distribution limite bootstrap de 2000 répétition de Paréto ($\gamma_1 = 0.2$), $n = 1000$	59
3.15	Échantillon conditionnel Y, ($\gamma_1 = 0.35$) et ($\gamma_2 = 2.5$) (10% censurées)	63
3.16	Échantillon conditionnel Y, ($\gamma_1 = 0.35$) et ($\gamma_2 = 1$) (25% censurées)	63
3.17	Échantillon conditionnel Y, ($\gamma_1 = 0.35$) et ($\gamma_2 = 0.5$) (40% censurées)	64
3.18	Comportement graphie de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 2.5$), (10% de censure), pour n=200	68
3.19	Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 2.5$), (10% de censure), pour n=1000	68

3.20	Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 2.5$), (10% de censure), pour n=5000	69
3.21	Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 1$), (25% de censure), pour n=200	69
3.22	Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 1$), (25% de censure), pour n=1000	70
3.23	Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 1$), (25% de censure), pour n=5000	70
3.24	Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 0.5$), (40% de censure), pour n=200	71
3.25	Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 0.5$), (40% de censure), pour n=1000	71
3.26	Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 0.5$), (40% de censure), pour n=5000	72
3.27	QQ-norm de la distribution limite bootstrap de 2000 répétition de Paréto ($\gamma_1 = 0.35$),censurée par($\gamma_2 = 1$), (25% de censure),n=1000	77
3.28	QQ-norm de la distribution limite bootstrap de 2000 répétition de Paréto ($\gamma_1 = 0.35$),censurée par($\gamma_2 = 0.5$), (40% de censure),n=1000	77

Liste des tableaux

1.1	Quelques distributions appartenant au domaine d'attraction de Fréchet . . .	22
1.2	Quelques distributions appartenant au domaine d'attraction de Gumble . . .	22
1.3	Quelques distributions appartenant au domaine d'attraction de Weibull . . .	22
3.1	Résultats de simulation pour n=1000	59
3.2	Résultats de simulation pour n=1000	76

0.1 Abréviations et Notations

EVD	Distribution des valeurs extrêmes
EVI, γ	Indice des valeurs extrêmes
F	Fonction de répartition
F_n	Fonction de répartition empirique
F^{\leftarrow}	Inverse généralisée de F
GEV	Distribution des valeurs extrêmes généralisée
GPD	Distribution de pareto généralisée
G_γ	Famille des la loi de valeurs extrêmes généralisée
$i.i.d$	Indépendantes et identiquement distribuées.
$\mathbb{I}_{\{A\}}$	Fonction indicatrice de l'ensemble A
Λ	Loi de Gumbel
$\ell(x)$	Fonction à variation lente
DA	Domaine d'attraction de maxcimum
$M_n = X_{n,n}$	Maxcimum de X_1, \dots, X_n
VR	Vriation régulier
$v.a$	Variable aléatiore
$p.s$	Prèsque sûre
Φ	Loi de frêchet
Ψ	Loi de weibull
$resp$	Respectivement
$S = \bar{F}$	$1 - F$ fonction de survie
TEV	Théorème des valeurs extrêmes
$X_{1:n}, \dots, X_{n:n}$	Statistique d'ordre associées à X_1, \dots, X_n
$X \wedge Y$	$\min(X, Y)$
x_F	Point terminal
\mathcal{L}	Égalité en loi
$=$	Égalité en définition
$:=$	Égalité en définition
\xrightarrow{D}	Converge en distribution
\xrightarrow{l}	Converge en loi
\xrightarrow{p}	Converge en probabilité
$\xrightarrow{p.s}$	Converge presque sûre
$\xrightarrow{o_P}(\cdot)$	Converge vers 0 en probabilité

$O_P(\cdot)$	Être borné en probabilité
VR_α	Variation régulière d'indice α
$\hat{\gamma}^{(c,H)}$	Estimateur de Hill avec les données censurées
MSE	L'erreur quadratique moyenne
$Z^* = (Z_1^*, \dots, Z_n^*)$	Échantillon Bootstrap
se	Erreur standard
$al.$	Autres
$(\Omega, \mathcal{A}, \mathbb{P})$	Espace probabilisé
$\sup A$	Supremum de l'ensemble A
$s.o$	Statistique d'ordre

0.2 Introduction

Dans la nuit du 31 janvier au 1^{er} février 1953, de très fortes tempêtes traversant la mer du Nord balayèrent les côtes flamandes et néerlandaises d'ouest en est. Après que plusieurs digues eurent cédé, les provinces néerlandaises de la Hollande et de la Zélande furent particulièrement touchées. Les conséquences de ce raz-de-marée furent désastreuses. On dénombra plus de 2500 morts, 47000 habitations inondées et 10000 détruites, 200000 hectares de terres inondées, 30 000 têtes de bétail noyées, environ 9% des fermes des Pays-Bas inondées et plus de 400 brèches dans les digues.

A la suite de cette tempête, le gouvernement néerlandais décida la mise en place du "plan Delta" pour se prémunir contre une nouvelle inondation. Il fallut construire un nouveau réseau de digues renforcées de plus de 500 km le long de la côte de la mer du Nord.

Un comité composé de nombreux scientifiques, parmi lesquels des statisticiens, se réunit afin d'étudier le phénomène et proposer ainsi des recommandations sur les hauteurs des digues. Il fallut prendre en compte des facteurs économiques (coût de construction, coût des inondations,...), des facteurs physiques (rôle du vent sur la marée,...), mais aussi les hauteurs de marées enregistrées lors des précédentes inondations. En 1953, les inondations avaient entraîné une montée des eaux à 3.85 mètres au-dessus du niveau de la mer, soit largement en-deçà des 4 mètres atteints le 1^{er} novembre 1570, soit 382 ans auparavant.

Le but était de construire des digues assez grandes, de telle sorte qu'aucune vague ne les dépasse dans un horizon de 10000 ans. Autrement dit, il convenait de déterminer quelle serait la hauteur de la plus grande vague dans un horizon de temps de 10000 ans.

La difficulté résidait dans le fait que, pour calculer la hauteur des digues et donc la hauteur maximale d'une vague qui n'a jamais eu lieu, le comité d'experts devait se baser sur les informations des années précédentes; or il n'y avait que très peu de données disponibles en particulier pour des événements de cette ampleur.

Le 4 octobre 1986 marqua l'achèvement du plan Delta aux Pays-Bas. Il s'agit du plus grand chantier de génie civil de tous les temps. Il permit de relier toutes les îles côtières de la province de Zélande par des digues.

Après calcul, le comité d'expert estima que ces digues devraient mesurer au moins 5 mètres, estimation fondée sur l'utilisation de techniques statistiques empruntées à la théorie des valeurs extrêmes qui constitue le sujet de la présente thèse.

Il est d'un grand intérêt de se prémunir contre les risques extrêmes, qu'ils résultent d'une crise financière, d'un accident nucléaire ou d'une catastrophe naturelle, compte tenu des répercussions humaines, économiques et financières que ces derniers peuvent avoir.

Les inondations de 1570 et de 1953 peuvent nous conduire à nous interroger sur une possible récurrence de ces événements. Le but serait d'être en mesure de prédire l'apparition de tels phénomènes, leurs impacts et le retour de ces derniers. Pour cela, il convient de définir précisément ce qu'est un événement extrême.

Un événement extrême est un événement qui a une faible probabilité de se produire mais qui, lorsqu'il se produit, prend de très petites ou de très grandes valeurs et a un grand impact. On notera la différence avec un événement rare qui, par définition, est un événement dont la probabilité d'occurrence est faible. Le fait qu'un événement soit rare n'implique pas qu'il soit extrême; il est dépourvu de la notion de quantifiabilité (petites ou grandes valeurs). A l'inverse, tout événement extrême est rare au sens où il a une faible probabilité de se produire.

La théorie des valeurs extrêmes trouve à s'appliquer dans de nombreux domaines tels qu'en fiabilité, en métallurgie [9] et en astrophysique [17]. Elle intéresse également les sciences de l'environnement, avec la modélisation de grands feux de forêts[2] ainsi que la climatologie[46] et la météorologie [[14], [15], [49]].

Dans un premier article, Einmahl et Magnus [21] proposèrent une application aux temps limites de records en athlétisme ; dans un second, ils s'intéressèrent plus particulièrement à l'estimation du temps minimal possible sur 100m [38].

Aarssen et de Haan proposèrent des résultats afin de calculer l'âge limite possible de l'être humain. Pour d'autres exemples d'applications, se référer au livre de Reiss et Thomas[43] . Au sujet de la mise en garde d'une mauvaise utilisation de la théorie des valeurs extrêmes, on citera Bouleau[11] .

Le domaine d'application historique reste l'hydrologie , notamment suite aux travaux de Jules Emile Gumbel en 1954[35] et son ouvrage [36]en 1958. Plus récemment l'utilisation de la théorie des valeurs extrêmes a été vivement recommandée par le rapport Flood Study Report NERC [41] afin de modéliser les lois de probabilités des maxima annuels de précipitations et de crues. Ce rapport, établi en 1975, utilise et teste de nombreux outils statistiques issus de la théorie des valeurs extrêmes sur un très grand nombre de jeux de données de crues de rivières et de pluviométrie en Grande-Bretagne.

La modélisation des valeurs extrêmes censurées voit le jour en première fois en 1997 dans la littérature des extrêmes avec la sortie du livre Reiss et Thomas. Il a fallu qu'en 2007 Beirlant et al. abordent réellement la statistique non paramétrique des valeurs extrêmes avec des données censurées. Leur estimateur est basé sur un estimateur standard de l'indice de queue divisé par l'estimateur de la proportion de données non censurées dépassant un certain seuil donné. Ils ont appliqué cette théorie sur des données du SIDA. Puis, Einmahl et al. 2008 ont utilisé le même concept pour proposer un estimateur de l'indice de queue sur les k-plus grandes valeurs ensuite déterminer ses propriétés asymptotiques et enfin illustrer son comportement sur ces mêmes données du SIDA. Puis, la recherche sur la théorie des valeurs extrêmes censurées est devenue une actualité.

Ce mémoire est réparti en trois chapitre.

Nous présentons dans 1^{eme} *chapitre* quelques rappels essentielles sur la théorie des valeurs extrêmes et la modélisation des valeurs extrêmes conditionnelles et de la notion de censure qui permettra de faciliter la lecture du mémoire. Ainsi, il s'agira de présenter brièvement les résultats essentiels rencontrés dans la littérature. Nous définissons rapidement les notions de domaine d'attraction, fonctions à variations régulières puis nous présentons ensuite quelques estimateurs de l'indice de queue, ainsi que ces propriétés asymptotiques. Puis nous définissons définir la fonction de répartition conditionnelle. Quand à la censure nous présenterons quelques définitions liées à la statistique des durées de survie.

Dans le 2^{eme} *chapitre*, nous considérons le problème de l'estimation de l'indice des valeurs extrêmes conditionnels pour covariable fixé et aléatoire pour des données complètes et censurées, nous présentons ensuite un aperçu historique des quelques estimateurs de l'indice de queue conditionnels, ainsi que leur propriétés asymptotiques dans le cas des données complètes et censurées. Puis nous présenterons la définition de l'estimateur harmonique sur laquelle nous travaillerons dans le cas d'un échantillon conditionnel pour des données complètes et censuré, avant cela nous avons expliqué la méthode de fenêtre mobile pour sélectionner notre échantillon conditionnel par l'adaptation de l'estimateur de BEARAN .J et al 2013. Pour la 3^{eme} *chapitre*, Nous présentons une définition simplifiée du logiciel R avec certains de ses avantages mentionnés et nous présentons le principe de bootstrap. Puis nous simulons les échantillons conditionnels et notre nouvel estimateur harmonique et nous appliquons la méthode du bootstrap sur ce dernier dans le cas des données complètes et censurées pour une loi à variation régulière d'indice positif. Nous avons présenté à chaque algorithme sont programme correspondant aboutissant des graphiques significatifs.

Notamment une distribution limite empirique a été donnée par la méthode du bootstrap dans le but de calculer des paramètres de dispersion empirique comme le biais est l'erreur quadratique moyenne.

Chapitre 1

Rappels sur la théorie des valeurs extrêmes et sur la censure

La théorie des valeurs extrêmes communément appelée « Extrême Value Theory » (EVT) en anglais, est une vaste théorie dont le but est d'étudier les événements rares c'est-à-dire les événements dont la probabilité d'apparition est faible. Autrement dit elle essaie d'amener des éléments de réponses aux intempéries, aux inondations, aux catastrophes naturelles, aux problèmes financiers, etc. en prédisant leurs occurrences dans les années à venir. En d'autres termes on veut estimer des petites probabilités ou des quantités dont la probabilité d'observation est très faible c'est-à-dire proche de zéro.

1.1 Comportement du maximum d'un échantillon

1.1.1 Densités et loi des statistiques d'ordre

Définition 1.1.1 (*Statistiques d'ordre*). soit (X_1, \dots, X_n) une suite de variable aléatoire indépendantes et identiquement distribuées, classée par ordre croissant. On écrit cette suite d'observation sous la notation $X_{i,n}$ tel que :

$$X_{1,n} \leq X_{2,n} \leq \dots \leq X_{n,n},$$

où

- $X_{i,n}$: la $i^{\text{ème}}$ statistique d'ordre (statistique d'ordre i) dans un échantillon de taille n .

- $X_{1,n}$: la plus petite valeur observée (ou statistique de minimum) avec

$$X_{1,n} = \min(X_1, \dots, X_n)$$

- $X_{n,n}$: la plus grande statistique d'ordre (ou statistique de maximum) avec

$$X_{n,n} = \max(X_1, \dots, X_n)$$

Dans un échantillon de taille n , deux statistiques d'ordre sont particulièrement intéressantes pour l'étude des événements extrêmes, le minimum et le maximum : et

David [1970] et Balakrishnan et Clifford Cohen [1991] montrent que l'expression de la distribution de $X_{i,n}$ est

$$F_{i,n} = \mathbb{P}\{X_{i,n} \leq x\} = \sum_{r=i}^n \binom{n}{r} (F(x))^r (1 - F(x))^{n-r}.$$

Nous en déduisons que la fonction de densité est :

$$f_{i,n}(x) = \frac{n!}{(i-1)!(n-i)!} [F(x)]^{i-1} [1-F(x)]^{n-i} f(x),$$

où $f(x)$ est la densité de probabilité de X_i et F sa fonction de répartition associée. En utilisant la propriété d'indépendance des variables aléatoires X_1, \dots, X_n , on obtient :

Loi de $X_{1,n}$.

$$F_{1,n}(x) = \mathbb{P}\{X_{1,n} \leq x\} = 1 - (1 - F(x))^n,$$

d'où

$$f_{1,n}(x) = n f(x) (1 - F(x))^{n-1}.$$

Loi de $X_{n,n}$.

$$F_{n,n}(x) = \mathbb{P}\{X_{n,n} \leq x\} = (F(x))^n,$$

d'où

$$f_{n,n}(x) = n f(x) (F(x))^{n-1}.$$

Remarque 1.1.1

$$\mathbb{P}\{X_{n,n} \leq x\} = (F(x))^n \rightarrow 0 \text{ ou } 1 \text{ quand } n \rightarrow \infty.$$

Les expressions de $F_{1:n}$ et $F_{n:n}$ peuvent s'obtenir très facilement en considérant les relations[37]

$$\begin{aligned} \{X_{1:n} \geq x\} &\Leftrightarrow \{\min(X_1, \dots, X_n) \geq x\} \\ &\Leftrightarrow \bigcap_{i=1}^n \{X_i \geq x\} \end{aligned}$$

et

$$\begin{aligned} \{X_{n:n} \leq x\} &\Leftrightarrow \{\max(X_1, \dots, X_n) \leq x\} \\ &\Leftrightarrow \bigcap_{i=1}^n \{X_i \leq x\} \end{aligned}$$

En utilisant la propriété d'indépendance des variables aléatoires X_1, \dots, X_n nous en déduisons que

$$\begin{aligned} F_{1:n}(x) &= \mathbb{P}\{X_{1:n} \leq x\} \\ &= 1 - \mathbb{P}\{X_{1:n} \geq x\} \\ &= 1 - \mathbb{P}\left\{\bigcap_{i=1}^n \{X_i \geq x\}\right\} \\ &= 1 - \prod_{i=1}^n \mathbb{P}\{X_i \geq x\} \\ &= 1 - \prod_{i=1}^n [1 - \mathbb{P}\{X_i \leq x\}] \\ &= 1 - [1 - F(x)]^n \end{aligned}$$

et

$$\begin{aligned}
 F_{n:n}(x) &= \mathbb{P}\{X_{n:n} \leq x\} \\
 &= \mathbb{P}\left\{\bigcap_{i=1}^n \{X_i \leq x\}\right\} \\
 &= \prod_{i=1}^n \mathbb{P}\{X_i \leq x\} \\
 &= [F(x)]^n.
 \end{aligned}$$

Définition 1.1.2 (La fonction de répartition empirique). La fonction de répartition empirique de l'échantillon (X_1, \dots, X_n) notée F_n est donnée par :

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{]-\infty, x[}(X_i), \quad x \in \mathbb{R}$$

Il existe une autre version de la définition de F_n en utilisant les (s.o) comme suit :

$$F_n(x) = \begin{cases} 0 & si, \quad x \leq X_{1,n} \\ \frac{i-1}{n} & si, \quad X_{i-1,n} < x \leq X_{i,n}, \quad 2 \leq i < n \\ 1 & si, \quad x > X_{n,n} \end{cases}$$

Définition 1.1.3 (Fonction de quantile) Pour tout $0 < s < 1$, la fonction de quantile associée à F est définie par

$$Q(s) = \inf\{t : F(t) \geq s\}, = F^{-1}(s)$$

où F^{-1} représente la fonction inverse généralisée de F . On l'exprime en termes de la fonction de survie par :

$$Q(s) = \inf\{t : \bar{F}(t) < 1 - s\}, \quad 0 < s < 1.$$

Définition 1.1.4 (Quantile empirique) la fonction de quantile empirique de l'échantillon X_1, \dots, X_n est défini par :

$$Q_n(s) = F_n^{-1}(s) = \inf\{t : F_n(t) \geq s\}, \quad 0 < s < 1$$

où

$$Q_n(s) = \bar{F}_n^{-1}(1 - s) = \inf\{t : \bar{F}_n(t) < 1 - s\}, \quad 0 < s < 1$$

1.1.2 Comportement asymptotique des extrêmes

On définit la variable aléatoire M_n , qui traduit le maximum d'un n-échantillon d'une variable aléatoire X , (les variables aléatoires X_i sont indépendantes et suivent la même loi que X) par :

$$M_n = \max(X_i)_{1 \leq i \leq n}$$

On pourrait aussi s'intéresser au minimum en utilisant la relation :

$$\min(X_i)_{1 \leq i \leq n} = -\max(-X_i)_{1 \leq i \leq n}$$

Le résultat central de la théorie des valeurs extrêmes concerne la distribution asymptotique du maximum (ou le minimum) en fonction de celle de la variable aléatoire X . Notons la fonction de répartition F_X de la variable aléatoire de loi de probabilité \mathbb{P} , à savoir $F_X(x) = \mathbb{P}(X < x)$. La fonction de répartition de M_n est alors définie par :

$$\begin{aligned} F_{M_n}(x) &= \mathbb{P}(M_n \leq x) \\ &= \mathbb{P}(X_1 \leq x, \dots, X_n \leq x) \\ &= \mathbb{P}(X_1 \leq x) \times \dots \times \mathbb{P}(X_n \leq x) \\ &= [F_X(x)]^n \end{aligned}$$

De ces résultats, nous tirons la conclusion que le maximum M_n est une variable aléatoire dont la fonction de répartition est égale à $(F_X)^n$. La fonction de répartition de X étant souvent inconnue et généralement pas possible d'être déterminée. Notons $x_F = \sup\{x \in \mathbb{R} : F_X(x) < 1\}$ le point terminal à droite (right-end point) de la fonction de répartition F_X . Ce point terminal peut être infini ou fini (Embrechts et al. 1997). On s'intéresse ici à la distribution asymptotique du maximum, en faisant tendre n vers l'infini,

$$\lim_{n \rightarrow \infty} F_{M_n}(x) = \lim_{n \rightarrow \infty} [F_X(x)]^n = \begin{cases} 0 & \text{si } F(x) < 1 \\ 1 & \text{si } F(x) = 1 \end{cases}$$

On constate que la distribution asymptotique du maximum, donne une loi dégénérée, une masse de Dirac en x_F , puisque pour certaines valeurs de x , la probabilité peut être égale à 1 dans le cas où x_F est fini. et donc M_n tend vers x_F presque sûrement, Ce fait ne fournit pas assez d'informations, d'où l'idée d'utiliser une transformation afin d'obtenir des résultats plus exploitables pour les loi limites des maxima M_n . On s'intéresse par conséquent à une loi non dégénérée pour le maximum, la théorie des valeurs extrêmes permet de donner une réponse à cette problématique. Les premiers résultats sur la caractérisation du comportement asymptotique des maxima M_n convenablement normalisés et donnés par la suite.

1.2 Présentation de la théorie des valeurs extrêmes

1.2.1 Loi des valeurs extrêmes

Comme la fonction de répartition obtenue précédemment conduit à une loi dégénérée lorsque n tend vers l'infini, on recherche une loi non dégénérée pour le maximum de X . Cette loi limite non dégénérée est fournie par le "théorème des distributions extrêmes" qui donne une condition nécessaire et suffisante pour l'existence d'une loi limite non dégénérée pour le maximum. Ce théorème est proposé par Gnedenko (1943) qui donne la forme des lois limites et Jenkinson (1955) qui en donne l'expression générale.

Théorème 1.2.1 (Fisher et Tippett, 1928, Gnedenko, 1943). Soit X_1, \dots, X_n une suite de n variables aléatoires réelles indépendantes et identiquement distribuées de loi continue P et $M_n = \max(X_i)_{1 \leq i \leq n}$. S'il existe deux suites réelles $(a_n)_{n \geq 1}$ et $(b_n)_{n \geq 1}$ avec $b_n > 0$, et une fonction de répartition non-dégénérée G_γ telle que,

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[\frac{M_n - b_n}{a_n} \leq x \right] = \lim_{n \rightarrow \infty} F^n(a_n x + b_n) = G_\gamma(x) \quad \forall x \in \mathbb{R} \quad (1.1)$$

Alors G est du même type qu'une des trois lois suivantes :

$$\text{loi de Gumbel : } \Lambda_\gamma(x) = \exp(-\exp(-x)) \quad -\infty < x < +\infty \quad (1.2)$$

$$\text{loi de Fréchet : } \Phi_\gamma(x) = \begin{cases} 0 & x < 0 \\ \exp(-x^{-1/\gamma}) & x \geq 0, \gamma > 0 \end{cases} \quad (1.3)$$

$$\text{loi de Weibull : } \Psi_\gamma(x) = \begin{cases} \exp(-(-x)^{-1/\gamma}) & x < 0, \gamma < 0 \\ 1 & x \geq 0, \end{cases} \quad (1.4)$$

avec G_γ est la loi des valeurs extrêmes et γ est l'indice des valeurs extrêmes. a_n et b_n sont des paramètres de normalisation. Ce théorème est proposé par Gnedenko (1943) qui donne la forme des lois limites.

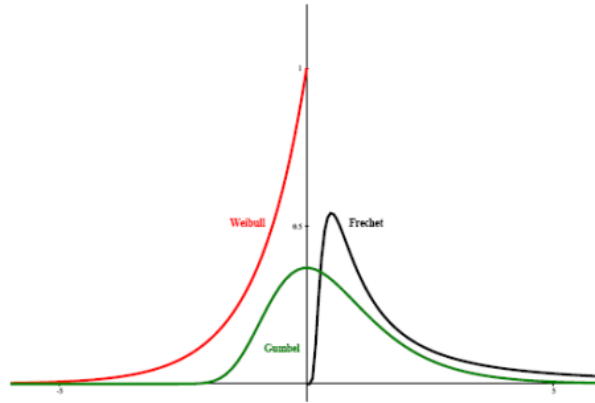


FIG. 1.1 – Densité de $\Phi_\gamma(x)$, $\Psi_\gamma(x)$, $\Lambda(x)$

Jenkinson (1955) donne l'expression générale notée GEV (Generalized Extreme Value Distribution) des trois distribution par,

$$G_\gamma(x) = \begin{cases} \exp(-(1 + \gamma x)^{-1/\gamma}), \forall x \in \mathbb{R}, 1 + \gamma x > 0 \text{ si } \gamma \neq 0 \\ \exp(-\exp(-x)), \quad \forall x \in \mathbb{R}, & \text{si } \gamma = 0 \end{cases} \quad (1.5)$$

1.2.2 Caractérisation des Domaines d'attraction

On définit les notions de fonctions à variations régulières et de fonctions à variations lentes qui nous seront utiles par la suite. Pour plus de détails, se référer à Bingham et al. [10] où de nombreux résultats sur les fonctions à variations régulières sont donnés.

Fonctions à variation régulière

Définition 1.2.1 . Une fonction mesurable $g : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ est à variation régulière à l'infini si et seulement si, il existe un réel α tel que, pour tout $x > 0$

$$\lim_{t \rightarrow \infty} \frac{g(tx)}{g(t)} = x^\alpha.$$

Et on note $g \in VR_\alpha$, α est appelé indice (ou exposant) de la fonction à variation régulière.

Dans le cas particulier où $\alpha = 0$, on dit que g est à variation lente à l'infini, c'est à dire,

$$\lim_{t \rightarrow \infty} \frac{g(tx)}{g(t)} = 1 \quad \forall x > 0$$

Les fonctions à variation lente sont génériquement notées $\ell(x)$. Pour toute fonction à variation lente ℓ à l'infini, on a :

$$\lim_{x \rightarrow \infty} \frac{\log(\ell(x))}{\log(x)} = 0$$

Proposition 1.2.1 Soient $\alpha \in \mathbb{R}$ et $g \in VR_\alpha$. Alors il existe une fonction à variation lente ℓ à l'infini telle que :

$$\forall x > 0, \quad g(x) = x^\alpha \ell(x)$$

Théorème 1.2.2 Une fonction $\ell : (0, \infty) \rightarrow (0, \infty)$ est à variation lente si et seulement si elle peut être écrite comme :

$$\ell(x) = c(x) \exp \left\{ \int_{x_0}^x \frac{\varepsilon(u)}{u} du \right\}, \quad x \geq x_0,$$

pour certain $x_0 > 0$, où $\lim_{x \rightarrow \infty} c(x) = c \in (0, \infty)$ et $\lim_{x \rightarrow \infty} \varepsilon(x) = 0$.

Des exemples typiques de fonctions à variation lente sont des constantes positives ou des fonctions convergeant vers une constante positive, logarithmes, puissances des logarithmes et logarithmes itérés.

Définition 1.2.2 . Une fonction de répartition F sur \mathbb{R} appartient à une classe à variation régulière VR s'il existe $\alpha \geq 0$ tel que $1 - F \in VR_{-\alpha}$ sur \mathbb{R} , ou d'une manière équivalente :

$$1 - F(x) \sim x^{-\alpha} \ell(x), \quad \text{quand } x \rightarrow \infty$$

pour certaines $\ell \in RV_0$.

Domaine d'attraction de Fréchet

Théorème 1.2.3 F appartient au domaine d'attraction de Fréchet avec un indice de valeur extrêmes $\gamma > 0$ si et seulement si $x_F = +\infty$ et $1 - F$ est une fonction à variation régulière d'indice $-1/\gamma$ c'est-à-dire,

$$1 - F = x^{-1/\gamma} \ell(x)$$

où ℓ est une fonction à variation lente. Dans ce cas, un choix possible pour les suites a_n et b_n est :

$$a_n = F^{-1} \left(1 - \frac{1}{n} \right) \quad \text{et} \quad b_n = 0.$$

Ce théorème permet de caractériser très simplement les distributions appartenant au domaine d'attraction de Fréchet. En effet, elles doivent vérifier

$$\lim_{t \rightarrow \infty} \frac{1 - F(tx)}{1 - F(t)} = x^{-\alpha}$$

Prenons par exemple le cas de la distribution Pareto. Nous avons

$$F(x) = 1 - x^{-1/\gamma}$$

Nous en déduisons que

$$\begin{aligned} \lim_{t \rightarrow \infty} \frac{1 - F(tx)}{1 - F(t)} &= \lim_{t \rightarrow \infty} \frac{(tx)^{-1/\gamma}}{t^{-1/\gamma}} \\ &= x^{-1/\gamma} \end{aligned}$$

Donc $1 - F \in RV_{-1/\gamma}$

Domaine d'attraction de Weibull

Théorème 1.2.4 F appartient au domaine d'attraction de Weibull avec un indice de valeur extrême $\gamma < 0$ si et seulement si $x_F < +\infty$ et $1 - F^*$ est une fonction à variation régulière d'indice $1/\gamma$ c'est-à-dire,

$$1 - F = (x_F - x)^{-1/\gamma} \ell \left[(x_F - x)^{-1} \right].$$

avec,

$$F^*(x) = \begin{cases} 0 & \text{si } x \leq 0, \\ F(x_F - x^{-1}) & \text{si } x > 0. \end{cases}$$

Dans ce domaine d'attraction les suites de normalisation a_n et b_n sont déterminées comme suit :

$$a_n = x_F - F^{-1} \left(1 - \frac{1}{n} \right) \quad \text{et} \quad b_n = x_F.$$

Domaine d'attraction de Gumbel

Définition 1.2.3 Soit F une fonction de répartition de point terminal x_F fini ou infini. S'il existe $z < x$ tel que

$$1 - F(x) = c \exp \left\{ - \int_z^x \frac{1}{a(t)} dt \right\}, \quad z < x < x_F,$$

où $c > 0$ et a une fonction positive absolument continue de densité a' vérifiant $\lim_{x \uparrow x_F} a'(x) = 0$.

Alors F est une fonction de Von-Mises et a est sa fonction auxiliaire.

Dans ce cas, un choix possible pour les suites (a_n) et (b_n) pour tout $n > 0$ est :

$$a_n = F^{-1} \left(1 - \frac{1}{n} \right) \quad \text{et} \quad b_n = \frac{1}{\overline{F}(a_n)} \int_{a_n}^{y_F} \overline{F}(z) dz.$$

Théorème 1.2.5 F appartient au domaine d'attraction de Gumbel si et seulement si il existe une fonction de Von-Mises F^* telle que pour $z < x < x_F$ on ait :

$$1 - F(x) = c(x) [1 - F^*(x)] = c(x) \exp \left\{ - \int_z^x \frac{1}{a(t)} dt \right\},$$

où $c(x) \rightarrow c > 0$ lorsque $x \rightarrow x_F$.

Donc selon le signe de γ on distingue les trois cas de domaines d'attraction (**D.A.**) :

1. Si $\gamma > 0$, F appartient au **D.A. de Fréchet**, et l'on note $F \in D.A.(Fréchet)$. Il contient toutes les lois dont la fonction de survie décroît comme une fonction puissance. Ce sont les lois à «**queue lourde**». Les distributions du domaine de Fréchet sont beaucoup utilisées en fiabilité mécanique, dans les phénomènes climatiques tels que la météorologie, l'hydrologie, la vitesse du vent enregistrée en continu dans les aéroports et en finance dans les études de risque.
2. Si $\gamma = 0$, F appartient au **D.A. de Gumbel**, et l'on note $F \in D.A.(Gumbel)$. Ce sont les lois dont la fonction de survie décroît vers zéro à une vitesse exponentielle. Ce sont les lois à «**queue légères**». Ces distributions sont souvent utilisées pour faire des prévisions dans les événements environnementaux tels que le séisme (le tremblement de terre), l'hydrologie (les inondations, la destruction des barrages), etc.
3. Si $\gamma < 0$, F appartient au **D.A. de Weibull**, et l'on note $F \in D.A.(Weibull)$. Ce domaine regroupe toutes les lois dont le point terminal, $x_F = \inf \{x, F(x) \geq 1\}$ est fini. Ce sont les lois à «**queue finie**». Les distributions de type de Weibull sont souvent utilisées pour décrire la résistance mécanique d'un matériau ou encore le temps de fonctionnement d'un appareil électronique ou mécanique.

Les Figures 1.1 et 1.2[13] ci-dessous illustre le comportement de différentes distributions GEV correspondant à différentes valeurs de γ . Les Tableaux (Tableau 1.1, Tableau 1.2, Tableau 1.3, ([50]) donnent quelques lois et leur domaine d'attraction.

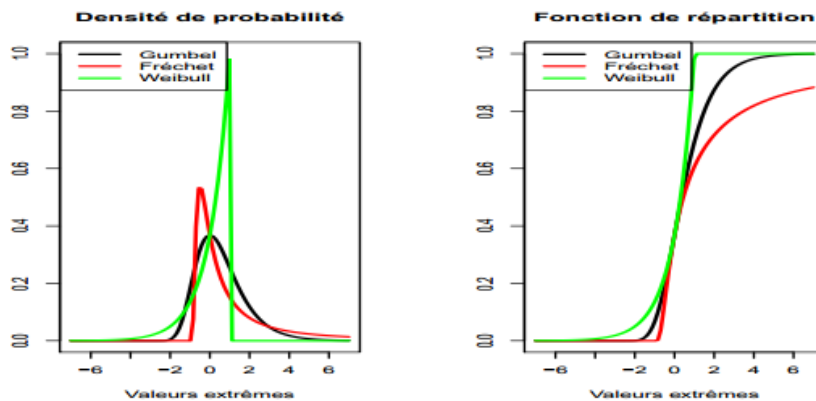


FIG. 1.2 – Représentation de la densité et fonction de répartition : Gumbel ($\gamma = 0$), Fréchet ($\gamma = 1$) et Weibull ($\gamma = -1$)

Distribution	$1 - \mathbf{F}(\mathbf{x})$	γ
Burr(β, τ, λ), $\beta > 0, \tau > 0, \lambda > 0$	$\left(\frac{\beta}{\beta+x^\tau}\right)^\lambda$	$\frac{1}{\lambda\tau}$
Fréchet($\frac{1}{\alpha}$), $\alpha > 0$	$1 - \exp(-x^{-\alpha})$	$\frac{1}{\alpha}$
Loggamma($m; \lambda$), $m > 0, \lambda > 0$	$\frac{\lambda^m}{\Gamma(m)} \int_x^\infty (\log(u))^{m-1} u^{-(\beta+1)} du$	$\frac{1}{\lambda}$
Loglogistic(β, α), $\beta > 0, \alpha > 0$	$\frac{1}{1+\beta x^\alpha}$	$\frac{1}{\alpha}$
Pareto(α), $\alpha > 0$	$x^{-\alpha}$	$\frac{1}{\alpha}$

TAB. 1.1 – Quelques distributions appartenant au domaine d’attraction de Fréchet

Distribution	$1 - \mathbf{F}(\mathbf{x})$	γ
Gamma(m, λ), $m \in \mathbb{N}, \lambda > 0$	$\frac{\lambda^m}{\Gamma(m)} \int_x^\infty u^{m-1} \exp(-\lambda u) du$	0
Gumble(μ, β), $\mu \in \mathbb{R}, \beta > 0$	$\exp(-\exp(-\frac{x-\mu}{\beta}))$	0
Logistic	$\frac{2}{1+\exp(x)}$	0
Log normale(μ, σ), $\mu \in \mathbb{R}, \sigma > 0$	$\frac{1}{\sqrt{2\pi}} \int_x^\infty \frac{1}{u} \exp(-\frac{1}{2\sigma^2}(\log u - \mu)^2) du$	0
Weibull (λ, τ), $\lambda > 0, \tau > 0$	$\exp(-\lambda x^\tau)$	0

TAB. 1.2 – Quelques distributions appartenant au domaine d’attraction de Gumble

Distribution	$1 - \mathbf{F}(\mathbf{x})$	γ
Uniforme (0, 1)	$1 - x$	-1
Reverse Burr($\beta, \tau, \lambda, \tau_F$), $\beta > 0, \tau > 0, \lambda > 0$	$\left(\frac{\beta}{\beta+(\tau_F-x)^{-\tau}}\right)^\lambda$	$-\frac{1}{\lambda\tau}$

TAB. 1.3 – Quelques distributions appartenant au domaine d’attraction de Weibull

Proposition 1.2.2 [10] (Conditions du premier ordre, de Haan et Ferreira (2006)). Les assertions suivantes sont équivalentes :

1. F est à queue lourde

$$F \in D.A(\text{fréchet}), \gamma > 0$$

2. $1 - F$ est une fonction à variation régulière à l'infini d'indice $-1/\gamma$

$$\lim_{t \rightarrow \infty} \frac{1 - F(tx)}{1 - F(t)} = x^{-1/\gamma}, \quad x > 0 \quad (1.6)$$

3. $Q(1 - s)$ est une fonction à variations régulières à zéro d'indice $-\gamma$

$$\lim_{s \rightarrow 0} \frac{Q(1 - sx)}{Q(1 - s)} = x^{-\gamma}, \quad x > 0$$

4. U est une fonction à variation régulière à l'infini d'indice γ

$$\lim_{t \rightarrow \infty} \frac{U(tx)}{U(t)} = x^\gamma, \quad x > 0$$

Proposition 1.2.3 [10] (Conditions du second ordre de Haan et Ferreira (2006)). Une fonction de répartition $F(\cdot) \in D.A(\text{fréchet}), \gamma > 0$, admet une condition du second ordre à l'infini si elle satisfait à l'une des assertions suivantes :

1. Il existe un paramètre $\rho \leq 0$, et une fonction $A_1(\cdot)$ qui tend vers 0 (ne change pas de signe à l'infini) définie par, $\forall x > 0$

$$\lim_{t \rightarrow \infty} \frac{(1 - F(tx))/(1 - F(t)) - x^{-1/\gamma}}{A_1(t)} = x^{-1/\gamma} \frac{x^\rho - 1}{\rho}$$

2. S'il existe un paramètre $\rho \leq 0$ et une fonction $A_2(\cdot)$ qui tend vers 0 (ne change pas de signe à zéro) définie par, $\forall x > 0$

$$\lim_{s \rightarrow 0} \frac{Q(1 - sx)/Q(1 - s) - x^{-\gamma}}{A_2(s)} = x^{-\gamma} \frac{x^\rho - 1}{\rho},$$

3. S'il existe un paramètre $\rho \leq 0$, et une fonction $A(\cdot)$ qui tend vers 0 (ne change pas de signe à l'infini) définie par, $\forall x > 0$

$$\lim_{t \rightarrow \infty} \frac{U(tx)/U(t) - x^\gamma}{A(t)} = x^\gamma \frac{x^\rho - 1}{\rho}$$

si $\rho = 0$, on remplace $(x^\rho - 1)/\rho$ par $\log x$

Les fonctions $A(\cdot)$, $A_1(\cdot)$, $A_2(\cdot)$ sont à variations régulières à l'infini d'indices respectifs ρ , ρ/γ , et $-\rho$, avec $A_1(t) = A(1/(1 - F(t)))$ et $A_2(s) = A(1/s)$.

1.2.3 Estimation de l'indice des valeurs extrêmes γ

Dans la littérature de la TVE, il existe plusieurs méthodes et techniques pour l'estimation de l'IVE, dans cette partie on reste limiter à trois méthodes.

Soit X_1, \dots, X_n une suite de variable aléatoire indépendantes et identiquement distribuées et $X, \dots, X_{n,n}$ les statistique d'ordre associées. $k = k_n$ une suite d'entier satisfaisant :

$$1 < k < n, \quad k \longrightarrow \infty \quad \text{et} \quad \frac{k}{n} \longrightarrow 0 \quad \text{quand} \quad n \longrightarrow \infty \quad (1.7)$$

Estimateur de Pickands

L'estimateur de Pickands (1975) est le premier estimateur suggéré pour le paramètre $\gamma \in \mathbb{R}$ et il est donné par la formule suivante :

$$\hat{\gamma}_{k,n}^p = \frac{1}{\log(2)} \frac{X_{n-k,n} - X_{n-2k,n}}{X_{n-2k,n} - X_{n-4k,n}}$$

Théorème 1.2.6 (*Propriétés asymptotiques de $\hat{\gamma}_n^p$, [3]*)

Pour $\gamma > 0$, $F \in DA(G_\gamma(x))$ et k que vérifie (1.7) on a :

(a) *consistance faible*

$$\hat{\gamma}_{k,n}^p \xrightarrow{p} \gamma \quad \text{quand } n \rightarrow \infty$$

(b) *consistance forte* $k/\log \log(n) \rightarrow \infty$, pour $n \rightarrow \infty$, alors

$$\hat{\gamma}_{k,n}^H \xrightarrow{p.s} \gamma \quad \text{quand } n \rightarrow \infty$$

(c) *Normalité asymptotique* sous certaines conditions sur k et F on a

$$\sqrt{k}(\hat{\gamma}_{k,n}^p - \gamma) \xrightarrow{D} N(0, \sigma^2) \quad \text{quand } n \rightarrow \infty$$

où

$$\sigma^2 = \frac{\gamma \sqrt{(2^{2\gamma+1} + 1)}}{(2(2^\gamma - 1) \log(2))}$$

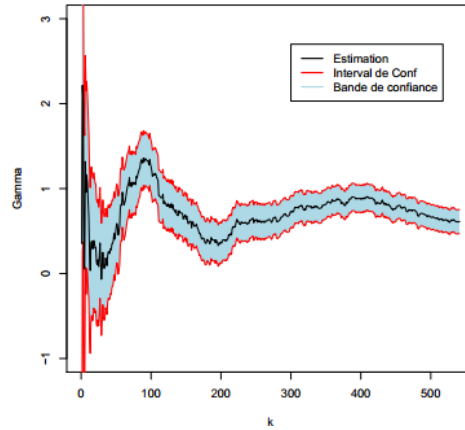


FIG. 1.3 – Comportement asymptotique de l'estimateur de Pickands avec des intervalles de confiance au niveau de confiance 0.95 sur les données d'assurance incendie Danoise pour la période allant de 01/01/1980 au 31/12/1990

Estimateur de Hill $\hat{\gamma}$

Cet estimateur qui a été présenté par *Hill* [39] en 1975 et l'estimateur le plus célèbre parmi tous les estimateurs de l'indice de queue. Il s'applique seulement dans le cas où l'indice de queue est positif ($\gamma > 0$), qui correspond aux distributions appartenant au domaine d'attraction de Fréchet. Cet estimateur donné par la définition suivante :

Définition 1.2.4 (*Estimateur de Hill*) : Pour $\gamma < 0$ l'estimateur de Hill est défini par :

$$\begin{aligned}\hat{\gamma}_{k_n}^H &= \frac{1}{k} \sum_{i=1}^k \log X_{n-i+1,n} - \log X_{n-k,n} \\ &= \frac{1}{k} \sum_{i=1}^k i (\log X_{n-i+1,n} - \log X_{n-i,n})\end{aligned}\quad (1.8)$$

1. Pour obtenir l'estimateur de $\hat{\gamma}_{k_n}^H$, en effet la condition (1.6) a une forme équivalente

$$\lim_{t \rightarrow \infty} \frac{1}{\overline{F}(t)} \int_t^\infty \frac{F(x)}{x} dx = \gamma$$

Par l'intégration par partie, on obtient

$$\lim_{t \rightarrow \infty} \frac{1}{\overline{F}(t)} \int_t^\infty \log \frac{x}{t} dF(x) = \gamma$$

On remplace F par F_n et $t = X_{n-k,n}$, l'estimateur de Hill $\hat{\gamma}_{k_n}^H$ défini par

$$\hat{\gamma}_{k_n}^H = \frac{1}{\overline{F}(X_{n-k,n})} \int_{X_{n-k,n}}^\infty \log \frac{x}{X_{n-k,n}} dF(x)$$

où

$$\hat{\gamma}_{k_n}^H = \frac{1}{k} \sum_{i=1}^k \log X_{n-i+1,n} - \log X_{n-k,n}$$

Théoreme (Propriétés asymptotiques de $\hat{\gamma}_{k_n}^H$ [40],)

Pour $\gamma > 0$, $F \in DA(\Phi_\gamma)$ et k que vérifie (1.7) on a :

(a) *consistance faible*

$$\hat{\gamma}_{k_n}^H \xrightarrow{p} \gamma \quad \text{quand } n \rightarrow \infty$$

(b) *consistance forte*

$$\hat{\gamma}_{k_n}^H \xrightarrow{p.s} \gamma \quad \text{quand } n \rightarrow \infty$$

(c) *Normalité asymptotique*

Supposons que $F(\cdot)$ vérifie la condition (c) de la proposition 1.2.3 et $\sqrt{k}A(n/k) \rightarrow \lambda$ si $n \rightarrow \infty$ alors

$$\sqrt{k}(\hat{\gamma}_{k_n}^H - \gamma) \xrightarrow{d} N\left(\frac{\lambda}{1-\rho}\right) \quad \text{quand } n \rightarrow \infty$$

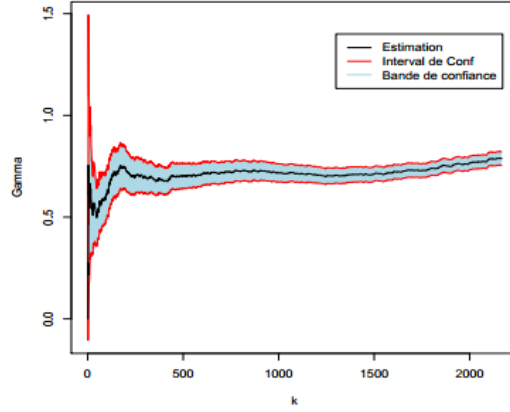


FIG. 1.4 – Comportement asymptotique de l’estimateur de Hill avec des intervalles de confiance au niveau de confiance 0.95 sur les données d’assurance incendie Danoise pour la période allant de 01/01/1980 au 31/12/1990.

Estimateur des Moments

Un autre estimateur qui peut être considéré comme une adaptation de l’estimateur de Hill, pour obtenir la consistance pour quelque soit le signe de l’indice, a été proposé par *Dekkers et al* c’est l’estimateur de moment.

Définition 1.2.5 (*Estimateur des Moments de Dekkers et al*)(1989,[13])

pour $\gamma \in \mathbb{R}$, l’estimateur de moment est défini par :

$$\begin{aligned}\hat{\gamma}_n^M &= M_{k,n}^{(1)} + 1 + \frac{1}{2} \left(1 - \frac{(M_{k,n}^{(1)})^2}{M_{k,n}^{(2)}} \right)^{-1} \\ &= \hat{\gamma}_{k,n}^H + 1 + \frac{1}{2} \left(1 - \frac{(\hat{\gamma}_{k,n}^H)^2}{M_{k,n}^{(2)}} \right)^{-1} \\ M_{k,n}^{(r)} &= \frac{1}{k} \sum_{ii=1}^k (\log X_{n-i+1,n} - \log X_{n-k,n})^r \quad r = 1, 2\end{aligned}$$

Théorème 1.2.7 (*Propriétés asymptotiques de $\hat{\gamma}_n^M$, [40]*)

Pour $\gamma > 0$, $F \in DA(\Phi_{1/\gamma})$ et k que vérifie (1.7) on a :

(a) *consistance faible*

$$\hat{\gamma}_n^M \xrightarrow{p} \gamma \quad \text{quand } n \rightarrow \infty$$

(b) *consistance forte* : si $(k/\log(n))^\delta \rightarrow \infty$ quand $n \rightarrow \infty$ avec $\delta > 0$

$$\hat{\gamma}_n^M \xrightarrow{p.s} \gamma \quad \text{quand } n \rightarrow \infty$$

(c) *Normalité asymptotique*

Supposons que $F(\cdot)$ vérifie la condition (3) de la proposition **1.2.4** et $\sqrt{k}A(n/k) \rightarrow \lambda$ si $n \rightarrow \infty$ alors

$$\sqrt{k}(\hat{\gamma}_n^M - \gamma) \xrightarrow{D} N(0, \eta^2) \quad \text{quand } n \rightarrow \infty$$

où

$$\eta \begin{cases} 1 + \gamma^2 & \text{si } \gamma > 0 \\ (1 - \gamma^2)(1 - 2\gamma) \left[4 - 8 \frac{(1-2\gamma)}{(1-3\gamma)} + \frac{(5-11\gamma)(1-2\gamma)}{(1-3\gamma)(1-4\gamma)} \right] & \text{si } \gamma < 0 \end{cases}$$

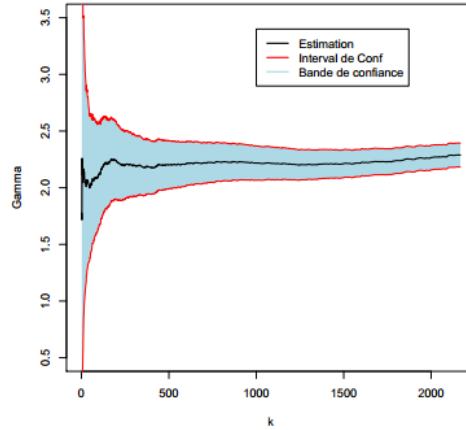


FIG. 1.5 – Comportement asymptotique de l'estimateur des moments avec des intervalles de confiance au niveau de confiance 0.95 sur les données d'assurance incendie Danoise pour la période allant de 01/01/1980 au 31/12/1990

Estimateur UH de Beirlant et al

Définition 1.2.6 (*Estimateur UH de Beirlant et al*)(1996,[40]) :

pour $\gamma \in \mathbb{R}$, l'estimateur de moment est défini par :

$$\hat{\gamma}_{k,n}^H = \frac{1}{k} \sum_{i=1}^k \log UH_{(i)} - \log UH_{(k+1)}$$

$$\frac{1}{k} \sum_{i=1}^k \log(X_{(n-i,n)} \hat{\gamma}_{i,n}^H) - \log(X_{(n-k-1,n)} \hat{\gamma}_{k+1,n}^H)$$

où

$$UH_i = Y_{n-j} \left(\frac{1}{j} \sum_{i=1}^k \log X_{(n-i+1,n)} - \log X_{(n-j,n)} \right) = X_{(n-j,n)} \hat{\gamma}_{j,n}^H$$

Théorème 1.2.8 (*Propriétés asymptotiques de $\hat{\gamma}_n^{UH}$* , [40])

Supposons que $F(\cdot)$ vérifie la condition (3) de la proposition **1.2.4** $\sqrt{k}(A(n/k) \rightarrow \lambda$ si $n \rightarrow \infty$ alors

$$\sqrt{k}(\hat{\gamma}_n^{UH} - \gamma) \xrightarrow{D} N\left(\frac{\lambda}{1-\rho}, \sigma^2\right) \quad n \rightarrow \infty$$

où

$$\sigma^2 = \begin{cases} 1 + \gamma^2 & \text{si } \gamma \geq 0 \\ \frac{(1-\gamma)(1+\gamma+2\gamma^2)}{1-2\gamma} & \text{si } \gamma < 0 \end{cases}$$

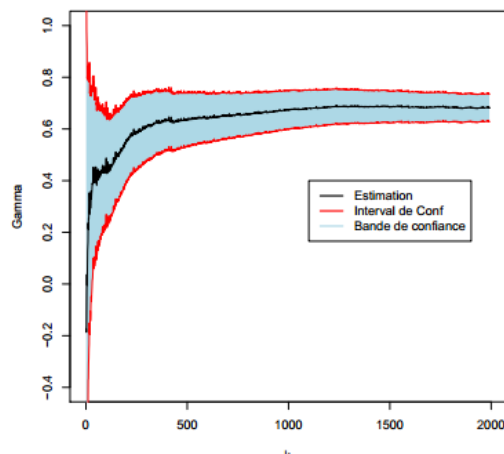


FIG. 1.6 – Comportement asymptotique de l’estimateur UH avec des intervalles de confiance au niveau de confiance 0.95 sur les données d’assurance incendie Danoise pour la période allant de 01/01/1980 au 31/12/1990.

1.3 Quelques rappelles sur la modélisation des valeurs extrêmes conditionnelles

1.3.1 Fonction de répartition conditionnelle

Soit Y une variable aléatoire réelle mesurée conjointement avec une covariable non-aléatoire $X \in X$, où X est un espace métrique muni d’une distance d . Soit $\{(X_i, Y_i), i = 1, \dots, n\}$ un échantillon d’observations indépendantes et identiquement distribuées du couple $(X, Y) \in X \times R$. la fonction de répartition conditionnelle $F(x, y)$ de Y_i sachant X_i , donnée par

$$F(y | x) = p(Y \leq y | X = x)$$

et la fonction de survie conditionnelle donnée par

$$\bar{F}(y | x) = p(Y > y | X = x)$$

dans le domaine de Fréchet où les queues de distribution sont lourdes : la fonction de distribution conditionnelle de Y sachant X se modélise alors de la façon suivante

$$F(y | x) = y^{-1/\gamma(x)} \ell(y, x)$$

où x est fixée dans un espace métrique muni d’une distance d , et $\ell(y; x)$ est une fonction à variations lentes.

1.3.2 Estimation de la fonction répartition conditionnelle

Méthode d’estimation indirecte

On dénombre deux techniques d’estimation indirectes de la fonction de répartition conditionnelle :

- Les estimateurs à noyaux qui peuvent être construits suivant :
 - (a) la méthode d’estimation du simple noyau ;
 - (b) la méthode d’estimation par noyau produit ou de Roussas (1969) ;
 - (c) la méthode du médianogramme ou de la fenêtre mobile.

- Les estimateurs au sens des plus proches voisins.
Nous nous intéressons à la première technique

Les estimateurs à noyaux

la méthode d'estimation du simple noyau :[1] l'estimateurs par simple noyau de la fonction de répartition conditionnelle défini par :

$$\widehat{F}_n(y | x) = \sum_{i=1}^n K\left(\frac{d(X_i, x)}{h_n}\right) \mathbb{I}\{Y_i \leq y\} / \sum_{i=1}^n K\left(\frac{d(X_i, x)}{h_n}\right) \quad (1.9)$$

où K est une densité de probabilité appelée « noyau » et h_n un paramètre qui converge vers zéro lorsque n tend vers l'infini

la méthode d'estimation par noyau produit ou de Roussas (1969) :[1] L'estimateur de la densité conditionnelle étant défini par le rapport entre les estimateurs de la densité du couple $f(x, y)$ et de la densité marginale $g(x)$, il en découle que l'estimateur à noyau de la fonction de répartition conditionnelle[2] est :

$$\widehat{F}_n(y | x) = \int_{-\infty}^y \widehat{f}_n(u | x) du = \frac{\int_{-\infty}^y \widehat{f}_n(x, u) du}{\widehat{g}_n(x)} \quad (1.10)$$

ou

$$\widehat{f}_n(x, y) = \frac{1}{nh_n^{d+1}} \sum_{i=1}^n K_0\left(\frac{d(X_i, x)}{h_n}\right) K_1\left(\frac{d(Y_i, y)}{h_n}\right) \quad (1.11)$$

$$\widehat{g}_n(x) = \int_{\mathbb{R}^P} \widehat{f}_n(x, y) dy = \frac{1}{nh_n^d} \sum_{i=1}^n K_0\left(\frac{d(X_i, x)}{h_n}\right) \quad (1.12)$$

avec K_0 et K_1 des noyaux de probabilité

la méthode du médianogramme ou de la fenêtre mobile :[1] Cet estimateur est valable quel que soit le design. Pour un réel fixé $h_n > 0$, on se place en un point fixé x et on sélectionne les seuls Y_i pour lesquels les points d'observations X_i (ou x_i) appartiennent à la boule centrée en x et de rayon h_n . Ceci revient alors à estimer la fonction de répartition conditionnelle par (1.9) en posant :

$$\begin{aligned} w_{in}(x) &= \mathbb{I}\{X_i \in B(x, h_n)\} / \sum_{i=1}^n \mathbb{I}\{X_i \in B(x, h_n)\}, \quad (\text{design aléatoire}) \\ w_{in}(x) &= \mathbb{I}\{x_i \in B(x, h_n)\} / \sum_{i=1}^n \mathbb{I}\{x_i \in B(x, h_n)\}, \quad (\text{design fixe}) \end{aligned} \quad (1.13)$$

où

$$w_{in}(x) = K\left(\frac{d(X_i, x)}{h_n}\right) / \sum_{i=1}^n K\left(\frac{d(X_i, x)}{h_n}\right) \quad (1.14)$$

et $B(x, h_n)$ est une boule centrée en x et de rayon $h_n \rightarrow 0$ quand $n \rightarrow \infty$.

1.3.3 Estimation des quantiles conditionnels

Estimation paramétrique des quantiles conditionnels

Elle est généralement utilisée quand l'on dispose d'un échantillon de petite taille. Afin d'estimer le quantile conditionnel, une façon de procéder consiste à supposer la fonction de répartition conditionnelle gaussienne. L'estimateur correspondant du quantile conditionnel $\hat{q}(\alpha | x)$ est défini par [1] :

$$\hat{q}(\alpha | x) = \hat{m}_n(x) + z_\alpha \hat{\sigma}_n(x)$$

où $\hat{m}_n(x)$ (resp $\hat{\sigma}_n(x)$) désigne l'estimateur de l'espérance conditionnelle (resp. l'écart type conditionnel) de Y sachant X = x et z_α le quantile d'ordre α de la loi normal centrée réduite.

Estimation non paramétrique des quantiles conditionnels

Dans le cadre de l'approche non-paramétrique, on distingue deux méthodes d'estimation. La première consiste à estimer au préalable la fonction de répartition conditionnelle puis à l'inverser pour en obtenir un estimateur du quantile conditionnel. La seconde consiste quant à elle en une estimation directe basée sur le principe des moindres carrés.

Nous ne discuterons que de la première méthode

Estimateurs à noyaux L'expression de l'estimateur du quantile conditionnel construit en inversant la fonction de répartition conditionnelle est donnée par :

$$\hat{q}(\alpha | x) = \inf\{y | \hat{F}_n(y | x) \geq \alpha\}$$

où

$\hat{F}_n(y | x)$ est Les estimateurs par simple noyau de la fonction de répartition conditionnelle qui nous donneront donc l'expression (1.9) ou l'estimateur de Roussas (1969) l'expression (1.10) ou l'estimateur de la fenêtre mobile l'expression (1.13)

1.4 Quelques généralités sur la censure

1.4.1 données de survie

L'analyse de survie, autrement dit la modélisation du temps de survenue d'un événement, apporte un outil principal d'évaluation théorique et pratique. L'analyse de ce type de données possède deux particularités intrinsèques, d'une part, celle-ci ne concerne que des variables aléatoires positives et d'autre part, la présence de données non complètement observées comme nous l'expliquons ci-dessous.

Désignons par X une variable d'intérêt, c'est à dire une variable aléatoire positive décrivant le temps qui s'écoule entre deux événements par exemple [4]

- *En fiabilité* : durée de vie d'une lampe, durée de vie d'un matériel...
- *En biologie* : en culture de cellules les durées d'apparition de parasites...
- *En médecine* : durée de guérison d'un patient, durée de rémission d'un malade...
- *En économie* : durée de chômage...
- *En éducation* : durée d'apprentissage d'une langue...
- *En assurance* : durée de vie d'un contrat qui peut être définie comme la différence entre la date de résiliation et la date de création du contrat.

Définition 1.6. (Fonction de survie). La fonction de survie qu'on note par $S(t)$ ou

$\bar{F}(t)$ est définie sur R_+ par

$$\begin{aligned} S(T) &= \bar{F}(t) \\ &= p(x > t), \quad t \geq 0 \end{aligned}$$

Pour t fixé c'est la probabilités de survivre jusqu'à l'instant t :

Définition 1.7. La fonction de risque instantané, pour t fixé représente la probabilité de mourir dans un petit intervalle de temps après t , conditionnellement au fait d'avoir survécu jusqu'au temps t (c'est-à-dire le risque de mort instantané pour ceux qui ont survécu) :

$$\begin{aligned} \lambda(t) &= \lim_{h \rightarrow 0} \frac{\mathbb{P}(t \leq X < t + h \mid X \geq t)}{h} \\ &= \frac{f(t)}{S(t)} \end{aligned}$$

où f est la densité de probabilité de X

Définition 1.8. (Fonctions empiriques de répartition et de survie) Soit $X_1; \dots, X_n$ un échantillon de taille $n \geq 1$ d'une variable positive X de fonction répartition F et de fonction de survie S . Les fonctions empiriques de répartition et de survie, F_n et S_n sont respectivement définies par

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{X_i \leq t\}, \quad \forall t \geq 0,$$

$$S_n(t) = 1 - F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{X_i \geq t\}, \quad \forall t \geq 0,$$

où $\mathbb{I}\{A\}$ est la fonction indicatrice de l'ensemble A

1.4.2 données censurées

Une des caractéristiques des données de survie est l'existence d'observations incomplètes. En effet, les données sont souvent recueillies partiellement, notamment, à cause des processus de censure et de troncature. Les données censurées proviennent du fait qu'on n'a pas accès à toute l'information. Au lieu d'observer des réalisations indépendantes et identiquement distribuées de durée X , on observe la réalisation de la variable X soumise à diverses perturbations indépendantes ou non de l'événement étudié.

Définition 1.10. La variable de censure C est définie par la non-observation de l'événement étudié. Si au lieu d'observer X , on observe C , et que l'on sait que $X > C$ (respectivement $X < C$, $C_1 < X < C_2$, on dit qu'il y a censure à droite (respectivement censure à gauche, censure par intervalle).

Caractéristiques

La censure est le phénomène le plus couramment rencontré lors du recueil de données en statistique. Pour un individu donné i , nous allons considérer [40]

- . son temps de survie X_i de fonction de répartition F
- . sa variable de censure C_i de fonction de répartition G
- . sa variable réellement observée Z_i de fonction de répartition H .

Dans la littérature on distingue trois types de censure [40] :

Censure à droite : La variable d'intérêt est dite censurée à droite si l'individu concerné n'a aucune information sur sa dernière observation. Ainsi, en présence de censure à droite les variables d'intérêt ne sont pas toutes observées. 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 trouve aussi ce genre de phénomène dans les 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 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.

Censure à gauche : Il y a censure à gauche lorsque l'individu a déjà subi l'événement avant qu'il soit observé. On sait uniquement que la variable d'intérêt est inférieure ou égale à une variable connue. Par exemple si nous voulons étudier en fiabilité un certain composant électronique qui est branché en parallèle avec un ou plusieurs autres composants : 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). Ainsi donc, la durée observée pour ce composant est censurée à gauche.

Censure par intervalle : Dans ce cas, comme son nom l'indique, on observe à la fois une borne inférieure et une borne supérieure de la variable d'intérêt. On retrouve ce modèle en général 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. Nous avons aussi ce genre de données qui sont censurées à droite ou, plus rarement, à gauche. Un avantage de ce type est qu'il permet de présenter les données censurées à droite ou à gauche par des intervalles du type $[c, +\infty[$ et $[0, c]$ respectivement.

Ces trois catégories de censure décrites ci-dessus peuvent se présenter en fonction du mode ou mécanisme de censure. Ainsi, dans la littérature on retrouve les types suivants :

La censure de type I

Définition 1.11. *La censure est dite non-aléatoire du type I si, étant donné un nombre positif fixé C et un n -échantillons X_1, \dots, X_n les observations consistent en (Z_i, δ_i) où*

$$\begin{cases} Z_i = X_i \wedge C \\ \delta_i = \mathbb{I}_{\{X_i \leq C\}} \end{cases}$$

La censure de type II

Elle est présente quand on décide d'observer les durées de survie des n patients jusqu'à ce que r d'entre eux soient décédés et d'arrêter l'étude à ce moment là. Soient $X_{(i)}$ et $Z_{(i)}$ les statistiques d'ordre des variables X_i et Z_i : La date de censure est donc $X_{(r)}$ et on observe les variables suivantes

$$\begin{aligned}
Z_{(1)} &= X_{(1)} \\
Z_{(2)} &= X_{(2)} \\
&\vdots \\
Z_{(r)} &= X_{(r)} \\
Z_{(r+1)} &= X_{(r)} \\
&\vdots \\
Z_{(n)} &= X_{(r)}
\end{aligned}$$

Définition 1.12. La censure est dite non-aléatoire du type II si, étant donné un nombre positif fixé r et un n -échantillons X_1, \dots, X_n les observations consistent en (Z_i, δ_i) où

$$\begin{cases} Z_i = X_i \wedge X_{(r)} \\ \delta_i = \mathbb{I}_{\{X_i \leq X_{(r)}\}} \end{cases}$$

La censure de type III : C'est la version aléatoire du type I

Définition 1.13. La censure est dite aléatoire du type I si, étant donné un n -échantillons X_1, \dots, X_n , il existe un v.a. n -dimensionnelle (C_1, \dots, C_n) de $(\mathbb{R}^+)^n$ telle que les observations consistent en (Z_i, δ_i) où

$$\begin{cases} Z_i = X_i \wedge C_i \\ \delta_i = \mathbb{I}_{\{X_i \leq C_i\}} \end{cases}$$

1.4.3 Estimation de la fonction de survie

Dans la littérature plusieurs auteurs se sont intéressés pour l'estimation de la fonction de survie dans le cas où les données sont censurées. Parmi ces derniers nous pouvons citer *Kaplan et Meier* (1958) ont proposé un estimateur de la fonction de survie.

Estimateur de Kaplan-Meier

Soit $(Z_i, \delta_i)_{1 \leq i \leq n}$ l'échantillon réellement observé et soit $(Z_{i,n}, \delta_{i,n})_{1 \leq i \leq n}$ sa statistique d'ordre croissant. L'estimateur de *Kaplan-Meier* est défini par :

$$\begin{aligned}
\hat{S}_n(t) &= 1 - \hat{F}_n(t) = \prod_{i=1}^n \left(\frac{n-i}{n-i+1} \right)^{\delta_{i,n} \mathbb{I}_{\{Z_{i,n} \leq t\}}} \\
&= \prod_{i=1}^n \left[1 - \frac{\delta_{i,n} \mathbb{I}_{\{Z_{i,n} \leq t\}}}{n-i+1} \right]
\end{aligned}$$

Il est aussi appelé « **produit limite** » car il s'obtient comme la limite d'un produit.

- . Cet estimateur de Kaplan-Meier est une fonction étagée avec des sauts seulement aux observations non-censurées.
- . La hauteur des sauts de cet estimateur est aléatoire.
- . Quand toutes les observations sont non-censurées alors on obtient la fonction de répartition empirique.

L'estimateur de Kaplan-Meier est asymptotiquement gaussien, précisément on a le résultat suivant :

Théorème 1.7 (Dreesbeke et Saporta (2011)). Si les fonctions de répartition de la survie et de la censure n'ont aucune discontinuité commune, alors :

$$\sup_{t \geq 1} \left| \hat{S}_n(t) - S(t) \right| \xrightarrow{p.s} 0$$

et pour tout $t \geq 0$,

$$\sqrt{n} \left(\hat{S}_n(t) - S(t) \right) \xrightarrow{d} W_t$$

où $(W_t)_{t \geq 0}$ est un processus gaussien centré qui vérifie pour tous t et s strictement positifs,

$$\text{cov}(W_s, W_t) = S(t)S(s) \int_0^{s \wedge t} \frac{dF(u)}{(1 - F(u))^2(1 - G(u))}.$$

Chapitre 2

Estimation de l'indice des valeurs extrêmes conditionnelles en présence de données Complètes et de données censurées

Dans ce chapitre, nous considérons le problème de l'estimation de l'indice des valeurs extrêmes conditionnels en présence de données complètes et de données censurées aléatoirement à droite. Les domaines privilégiés dans ce chapitre concernent ceux dont les fonctions de répartition sont à queues lourdes telles que l'hydrologie avec les inondations, la biologie avec la présence de données manquantes dans les bases de données, la finance avec les cracks boursiers, la télécommunication avec les trafics routiers, etc.

2.1 Estimation de l'indice des valeurs extrêmes conditionnelles en design fixe en présence de données Complètes

La connaissance de γ est nécessaire pour résoudre un certain nombre de problèmes l'analyse des valeurs extrêmes. L'estimation de l'indice de queue a été largement étudiée dans la littérature et plusieurs estimateurs proposés. L'estimateur le plus populaire a été proposé par Hill (1975) dans le contexte de distributions à queue lourde. Diebolt et al. (2008) ont examiné le cas des distributions de queue de Weibull et le cas général a été étudié par Dekkers et al (1989).

Dans les applications pratiques, il arrive souvent que la variable d'intérêt Y puisse être lié à une covariable X . Dans cette situation, l'indice de valeur extrême de la distribution conditionnelle de Y dépend de la valeur x observée de la covariable X et doit être mentionné dans le suivant, comme indice de queue conditionnel. son estimation a été abordée dans la littérature récente sur les valeurs extrême, principalement dans le cas « *design fixe* », c'est-à-dire lorsque les covariables ne sont pas aléatoires. Smith [48] et Davison et Smith, [18] ont considéré un modèle de régression tandis que Hall et Tajvidi [8] ont utilisé une approche semi-paramétrique pour estimer l'indice de queue conditionnel. Des méthodes entièrement non paramétriques ont été envisagées en utilisant des splines, polynômes locaux, une approche de fenêtre mobile ou une approche de voisine.

2.1.1 Définitions des estimateurs

Soient $(x_i)_{1 \leq i \leq n}$, un échantillon déterministe de X et $(y_i)_{1 \leq i \leq n}$ des réalisations des variables aléatoires Y dans l'espace (Ω, A, p) et la fonction de répartition conditionnelle de Y sachant x et notée $F(y, x)$ est à queue lourde. On veut définir quelque Estimateur de « indice de queue conditionnel » ou « indice des valeurs extrêmes conditionnel » et pour tout x non aléatoire

Estimateur de Hill adapté(1975,[38])

$$\gamma_{k_x m_{n,x}}^H = M_{k_x m_{n,x}}^1 = \frac{1}{k_x} \sum_{i=1}^{k_x} \left(\frac{y_{(m_{n,x}-i+1)}^x}{y_{(m_{n,x}-i)}^x} \right), i = 1 \dots m_{n,x} \quad (2.1)$$

Estimateur de Dekkers-Einmahl-de Haan adapté (1989,[40])

$$\hat{\gamma}_{k_x m_{n,x}}^M(x) = M_{k_x m_{n,x}}^{(1)} + 1 - \frac{1}{2} \left(1 - \frac{(M_{k_x m_{n,x}}^{(1)})^2}{M_{k_x m_{n,x}}^{(2)}} \right)^{-1}, i = 1 \dots m_{n,x} \quad (2.2)$$

où

$$M_{k_x m_{n,x}}^{(2)} = \frac{1}{k_x} \sum_{i=1}^{k_x} \left(i \log \left(\frac{y_{(m_{n,x}-i+1)}^x}{y_{(m_{n,x}-i)}^x} \right) \right)^2$$

Estimateur UH adapté (1996, [8])

$$\hat{\gamma}_{k_x m_{n,x}}^{UH}(x) = \frac{1}{k} \sum_{i=1}^k \log(y_{(m_{n,x}-i)}^x \hat{\gamma}_{i, m_{n,x}}^H) - \log(y_{(m_{n,x}-k_x)}^x \hat{\gamma}_{k_x, m_{n,x}}^H) \quad (2.3)$$

2.1.2 Propriétés asymptotiques de l'estimateur de l'indice

Pour tout $x \in X$, nous supposons les assertions suivantes :

- a $F(\cdot | x)$ et $G(\cdot | x)$ sont absolument continues,
- b Il existe une fonction $\rho(x) < 0$ et une fonction à variations régulières $b(\cdot | x)$ d'indice $\rho(x)$ telles que pour tout $u > 0$,

$$\lim_{t \rightarrow \infty} = \frac{H^{\leftarrow} \left(1 - \frac{1}{tu} | x \right) / H^{\leftarrow} \left(1 - \frac{1}{t} | x \right) - u^{\rho(x)}}{b(t, x)} = u^{\rho(x)} \frac{u^{\rho(x)} - 1}{\rho(x)}$$

- c Il existe des constantes positives c_u , z_u et $\alpha_u \leq 1$ telles que pour tout $x \in B(x', 1)$

$$\sup_{z \geq z_u} \left| \frac{\log u(z | x)}{\log u(z | x')} - 1 \right| \leq c_u d^{\alpha_u}(x, x')$$

où $u(z | x) = \inf \left\{ s; H(s | x) \geq 1 - \frac{1}{z} \right\}$

Pour tout $x \in X$ nous supposons les assertions suivantes si $n \rightarrow \infty$, $k_x \rightarrow \infty$, $k_x / m_{n,x} \rightarrow 0$ et:

- d $\sqrt{k_x} b \left(\frac{m_{n,x}}{k_x}, x \right) \rightarrow \lambda(x) < \infty$

- c** $\frac{1}{\sqrt{k_x}} \sum_{i=1}^{k_x} \left[p_x \left(H^{\leftarrow} \left(1 - \frac{i}{m_{n,x}} \mid x \right) \right) - p_x \right] \longrightarrow \varepsilon(x) < \infty$. Soient $c > 0$ et $A(s, t) := \{1 - k_x/m_{n,x} \leq t < 1, |t - s| \leq c\sqrt{k_x}/m_{n,x}, s < 1\}$ Nous assumons que si $n \rightarrow \infty$
- e** $\sqrt{k_x} \sup_{A(s,t)} |p_x(H^{\leftarrow}(t \mid x)) - p_x(H^{\leftarrow}(s \mid x))| \longrightarrow 0$.

Sous ces conditions, nous avons les résultats asymptotiques de les estimateure

Théorème 2.1.1 *soit $x \in X$. Sous les conditions **c-e** et s'il existe des fonctions $m(\cdot)$ et $\sigma(\cdot)$ telles que $\sqrt{k_x} \left(\hat{\gamma}_{k_x m_{n,x}}^{(\cdot)}(x) - \gamma(x) \right) \xrightarrow{d} N(m(x)\lambda(x), \sigma^2(x))$, alors, nous avons :*

$$\sqrt{k_x} \left(\hat{\gamma}_{k_x m_{n,x}}^{(c,\cdot)}(x) - \gamma_1(x) \right) \xrightarrow{d} N \left(\frac{\lambda(x)m(x) - \gamma_1(x)\varepsilon(x)}{p_x}, \frac{\sigma^2(x) + \gamma_1(x)^2 p_x(1 - p_x)}{p_x^2} \right)$$

Corollaire 2.1.1 *Sous les hypothèses **a-e** et $k_x^{1/2} h_{n,x}^{\alpha_u} \longrightarrow 0$, nous avons*

$$\sqrt{k_x} \left(\hat{\gamma}_{k_x m_{n,x}}^{(c,H)}(x) - \gamma_1(x) \right) \xrightarrow{d} N \left(\frac{-\gamma_1(x)\varepsilon(x)}{p_x} + \frac{\lambda(x)}{p_x(1 - \rho(x))}, \frac{\gamma_1^3(x)}{\gamma(x)} \right)$$

$$\sqrt{k_x} \left(\hat{\gamma}_{k_x m_{n,x}}^{(c,UH)}(x) - \gamma_1(x) \right) \xrightarrow{d} N \left(\frac{-\gamma_1(x)\varepsilon(x)}{p_x} + \frac{\lambda(x)}{p_x(1 - \rho(x))}, \frac{\gamma_1^2(x)}{\gamma^2(x)}(1 + \gamma_1(x)\gamma(x)) \right)$$

$$\sqrt{k_x} \left(\hat{\gamma}_{k_x m_{n,x}}^{(c,M)}(x) - \gamma_1(x) \right) \xrightarrow{d} N \left(\frac{-\gamma_1(x)\varepsilon(x)}{p_x} + \frac{\lambda(x)}{p_x(1 - \rho(x))}, \frac{\gamma_1^2(x)}{\gamma^2(x)}(1 + \gamma_1(x)\gamma(x)) \right)$$

2.2 Estimation de l'indice des valeurs extrêmes conditionnelles en design aléatoire en présence de données Complètes

Soit Y une variable aléatoire ayant comme fonction de répartition F et Y des copies indépendantes de Y . Lorsqu'une covariable X est disponible et la distribution de Y dépend de X , le problème est d'estimer l'indice de queue conditionnel de la fonction de répartition conditionnelle $F(\cdot \mid x)$ de Y sachant $X = x$. Motivés par les pluies extrêmes. Motivés par les pluies extrêmes mesurées conjointement avec leurs positions géographiques (longitude, altitude, latitude) considérées comme covariables, Gardes et Girard (2010, [27]) ont proposé un estimateur de l'indice de queue conditionnel et son quantile extrême conditionnel. En suivant la même logique, Ferrez et al. (2011, [25]) ont fait l'analyse des températures extrêmes avec des paramètres topologiques et Pisarenko et Sornette (2003, [42]) ont étudié les séismes extrêmes avec comme covariable leurs emplacements. L'estimation de l'indice des valeurs extrêmes conditionnel et quantiles extrêmes conditionnels avec covariables fixes (ou non aléatoires) a été étudiée assez largement dans la littérature récente de la statistique des valeurs extrêmes.

Nous nous référons à Beirlant et al. (2004, [6]), à Gardes et Girard (2008, [20]), à Gardes et al. (2010, [29]), à Stupfler (2013, [49]) pour un aperçu de la méthodologie disponible, y compris le cas où la variable explicative est fonctionnelle (Gardes et Girard, 2012, [28]). A ce jour, beaucoup d'auteurs s'attardent sur le cas à covariable aléatoire, en dépit de son intérêt pratique. Gardes et Stupfler (2014, [30]) et Goegebeur et al. (2014, [32]) adaptent l'estimateur

de Hill de l'indice de queue d'une distribution à queue lourde en présence d'une covariable aléatoire. L'estimateur des moments introduit par Dekkers et al. (1989,[19]) a été adapté aux covariables aléatoires par Goegebeur et al. (2014,[31]). Daouia et al. (2011, [16]) ont proposé un estimateur basé sur le noyau de quantiles extrêmes conditionnels avec des covariables aléatoires. Delafosse et Guillou (2002,[20]), Gomes et Oliveira (2003, [34]), Beirlant et al. (2007, [7]), Einmahl et al. (2008, [17]), Beirlant et al. (2010, [8]), Gomes et Neves (2011, [60]), Brahimy et al. (2013, [19]) et Worms et Worms (2014,[40]) ont proposé des estimateurs de l'indice des valeurs extrêmes conditionnel et parfois leurs quantiles extrêmes conditionnels. Par contre, Ndao et al. (2014,[40]), voir le chapitre précédent) ont proposé un estimateur de l'indice de queue conditionnel et un estimateur de quantile extrême conditionnel basé sur une adaptation de l'estimateur de Weissman (1978, [52]) et de la méthode de la fenêtre mobile

2.2.1 Définitions des estimateurs

Soient $(x_i)_{1 \leq i \leq n}$, un échantillon déterministe de X et $(y_i)_{1 \leq i \leq n}$ des réalisations des variables aléatoires Y dans l'espace (Ω, A, p) et la fonction de répartition conditionnelle de Y sachant x et notée $F(y, x)$ est à queue lourde. ,on veut définirquelque Estimateur de « indice de queue conditionnel » ou « indice des valeurs extrêmes conditionnel » et pour tout x aléatoire

(2014b, [59]) ont proposé un estimateur de l'indice des valeurs extrêmes conditionnel par la version du noyau de l'estimateur de Hill :

$$\gamma_{t_n}^H(x) = \sum_{i=1}^n K_h(x - X_i)(\log Z_i - \log t_n)\mathbb{I}_{\{z_i > t_n\}} / \sum_{i=1}^n K_h(x - X_i)\mathbb{I}_{\{z_i > t_n\}} \quad (2.4)$$

où $K_h(x) = h^{-p}K(x/h)$, k est une densité de probabilité dans R^p , $h = h_n$ est une suite positive non aléatoire appelée fenêtre telle que $h \rightarrow 0$ quand $n \rightarrow \infty$ et t_n est une suite positive non aléatoire appelée seuil telle que $t_n \rightarrow \infty$ si $n \rightarrow \infty$

2.2.2 Propriétés asymptotiques de l'estimateur de l'indice

Hypothèses lipschitziennes[40]

Pour tout $(x, x_0) \in X \times X$, nous supposons :

- a Il existe $c_\gamma > 0$ tel que $\left| \frac{1}{\gamma(x)} - \frac{1}{\gamma(x')} \right| \leq c_\gamma d(x, x')$
- b Il existe $c_L > 0$ et u_0 tel que

$$\left| \frac{\log L(u \mid x)}{\log(u)} - \frac{\log L(u \mid x')}{\log(u)} \right| \leq c_L d(x, x')$$

- c Il existe $c_g > 0$ tel que $|g(x) - g(x')| \leq c_g d(x, x')$.

Hypothèse classique (condition du second ordre)[10]

- b Pour tout $x \in R$, il existe une fonction $\rho(x) < 0$ et une fonction à variations régulières à l'infini $b(\cdot \mid x)$ d'indice $\rho(x)$ telles que pour tout $u > 0$,

$$\lim_{t \rightarrow \infty} \frac{H^{\leftarrow} \left(1 - \frac{1}{tu} \mid x\right) / H^{\leftarrow} \left(1 - \frac{1}{t} \mid x\right) - u^{\rho(x)}}{b(t \mid x)} = u^{\rho(x)} \frac{u^{\rho(x)} - 1}{\rho(x)}$$

Hypothèse sur le noyau[40]

Proposition 2.2.1 [10]Supposons les hypothèses **a-d** vérifiées. Soit t_n une suite non aléatoire telle que $t_n \rightarrow \infty, nh^p \bar{H}(t \mid x) \rightarrow \infty, nh^{p+2} \bar{H}(t \mid x)(\log t_n)^2 \rightarrow 0$ et $\sqrt{nh^p \bar{H}(t \mid x)b \left(\frac{1}{\bar{H}(t_n \mid x)} \right)} \rightarrow 0$ quand $n \rightarrow \infty$. Alors pour tout $x \in X$

$$\sqrt{nh^p \bar{H}(t \mid x)}(\widehat{p}_{t_n}(x) - p_x) \xrightarrow{d} N\left(0, \frac{p_x(1-p_x) \|K\|_2^2}{g(x)}\right) \quad \text{quand } n \rightarrow \infty$$

$$\text{ou } p_x = \gamma_2(x) / (\gamma_1(x) + \gamma_2(x)) \text{ et } \|K\|_2^2 = \int K^2(u) du$$

Théorème 2.2.1 Supposons les hypothèses **a-d** vérifiées. Soit t_n une suite non aléatoire telle que $t_n \rightarrow \infty, nh^p \bar{H}(t \mid x) \rightarrow \infty, nh^{p+2} \bar{H}(t \mid x)(\log t_n)^2 \rightarrow 0$ et $\sqrt{nh^p \bar{H}(t \mid x)b \left(\frac{1}{\bar{H}(t_n \mid x)} \right)} \rightarrow 0$ et $\sqrt{nh^p \bar{H}(t \mid x)b(t_n \mid x)} \rightarrow 0$ quand $n \rightarrow \infty$. Alors, pour tout $x \in X$ tel que $g(x) > 0$,

$$\sqrt{nh^p \bar{H}(t_n \mid x)}(\widehat{\gamma}_{t_n}^{c.H}(x) - \gamma_1(x)) \xrightarrow{d} N\left(0, \frac{\gamma_1^3(x) \|K\|_2^2}{\gamma(x)g(x)}\right) \quad \text{quand } n \rightarrow \infty$$

2.3 Estimation de l'indice des valeurs extrêmes conditionnelles en présence de données censurées

Dans cette partie on va intéresser au problème de l'estimation de l'IVE mais cette fois c'est en présence de données censurées aléatoirement à droite. Ce problème, est très récent dans la littérature, les premiers qui ont mentionné ce sujet sont Reiss et Thomas[43] en (2007) mais sans résultat asymptotique. En (2001) ; Beirlant et Guillou [5]ils ont proposé un estimateur mais pour les données tronquées. En 2007, Beirlant et al[4] ils ont introduit une méthode pour les estimateurs de Hill et de moment,... de plus, ils ont proposé les estimateurs des quantiles extrêmes et ont discuté leurs propriétés asymptotiques lorsque les données sont censurées pour un seuil déterministe et l'année suivante Einmahl et al.[23] ils ont adapté différents estimateurs de l'IVE au cas où les données sont censurées par un seuil aléatoire et ils ont proposé une méthode unifiée pour établir leur normalité asymptotique.

Soit F et G sont supposées être dans le domaine d'attraction maximales $D(G_{\gamma_1})$ et $D(G_{\gamma_2})$ respectivement où $\gamma_1, \gamma_2 \in \mathbb{R}$, l'échantillon $X_1 \dots X_n$ et $Y_1 \dots Y_n$. soit $\{(Z_i, \delta_i), 1 \leq i \leq n\}$ l'échantillon réellement observé où $Z_i = X_i \wedge Y_i$ et $\delta_i = \mathbb{I}_{\{X_i \leq Y_i\}}$, avec points terminales τ_F et τ_G , où $\tau_F = \sup \{x, F(x) < 1\}$, alors cela signifie que $H \in D(G_\gamma)$. Einmahl et al.(2008, [7]) ont examiné les trois cas les plus intéressants suivants:

$$\begin{cases} \text{cas 1} & : \gamma_1 > 0, \gamma_2 > 0 & \tau_F = \tau_G & \gamma = \frac{\gamma_1 \gamma_2}{\gamma_1 + \gamma_2} \\ \text{cas 2} & : \gamma_1 < 0, \gamma_2 < 0 & \tau_F = \tau_G & \gamma = \frac{\gamma_1 \gamma_2}{\gamma_1 + \gamma_2} \\ \text{cas 3} & : \gamma_1 = \gamma_2 = 0 & \tau_F = \tau_G = \infty & \gamma = 0 \end{cases} \quad (2.5)$$

les estimateurs de l'IVE de données censuré sont basés sur un estimateur standard de l'indice de queue divisé par l'estimateur de la proportion de données non censurées dans le plus grand k de Z

$$\widehat{\gamma}_{X,k,n}^{(c,\cdot)} = \frac{\widehat{\gamma}_{Z,k,n}^{(\cdot)}}{\widehat{p}} \quad (2.6)$$

ou

$$\widehat{p}_x = \frac{1}{k_x} \sum_{i=1}^{k_x} \delta_{[m_{n,x-i+1}]^x} \quad \text{pour la covariable } x \text{ fixe}$$

et

$$\hat{p}_{t_n}(x) = \frac{\overline{H}_n^1(t_n | x)}{\overline{H}(t_n | x)} \text{ pour la covariable } x \text{ aléatoire}$$

avec $\overline{H}(t_n | x) = \sum_{i=1}^n B_i(x) \mathbb{I}_{\{z_i > t_n\}}$, $\overline{H}^1(t_n | x) = \sum_{i=1}^n B_i(x) \mathbb{I}_{\{z_i > t_n, \delta_i = 1\}}$ et les poids de Nadarya-Watson définis par

$$B_i(x) = K\left(\frac{x - X_i}{h}\right) / \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right)$$

donc on va définir l'estimateur de l'indice de valeur extrême conditionnelle en présence de données censuré

Soit $(x_i)_{1 \leq i \leq n}$, un échantillon déterministe de X dans l'espace métrique (x, d) , $(y_i)_{1 \leq i \leq n}$ et $(c_i)_{1 \leq i \leq n}$ des réalisations respectives des variables aléatoires Y et C dans l'espace probabilisé $(\mathcal{X}, \mathcal{A}, P)$ avec F et G respectivement les fonctions de répartition conditionnelles. Soient $(Z_i)_{1 \leq i \leq n} = (Y_i \wedge C_i)_{1 \leq i \leq n}$ et $\delta_i = \mathbb{I}_{\{Y_i \leq C_i\}}$, $\mathbb{I}_{\{\cdot\}}$ est une fonction indicatrice avec Y et C sont indépendantes conditionnellement à X

– l'estimateur de l'indice de valeur extrême conditionnelle en design fixe en présence de données censuré défini par :

$$\hat{\gamma}_{X,k,n}^{(c,\cdot)} = \frac{\hat{\gamma}_{k_x m_{n,x}}^{(\cdot)}}{\hat{p}_x} \quad (2.7)$$

où

$$\hat{p}_x = \frac{1}{k_x} \sum_{i=1}^{k_x} \delta_{[m_{n,x} - i + 1]}^x$$

et $\hat{\gamma}_{k_x m_{n,x}}^{(\cdot)}$ est l'estimateur de l'indice de valeur extrême conditionnelle qui définit pour l'estimateur de Hill dans l'équation (2.1) ou l'estimateur de Dekkers-Einmahl-de Haan dans l'équation (2.2) ou l'estimateur UH dans l'équation (2.3)

– l'estimateur de l'indice de valeur extrême conditionnelle en design aléatoire en présence de données censuré défini par :

$$\hat{\gamma}_{X,k,n}^{(c,\cdot)} = \frac{\hat{\gamma}_{k_x m_{n,x}}^{(\cdot)}}{\hat{p}_x}$$

$$\hat{p}_{t_n}(x) = \frac{\overline{H}_n^1(t_n | x)}{\overline{H}(t_n | x)}$$

2.4 Construction d'un nouvel estimateur de l'indice des valeurs extrêmes conditionnelles

Avant de parler de nouvelle estimateur nous allons voir l'estimateur d'indice de queue de moment harmonique qui nous allons définir sous

Définition 2.4.1 (*Estimateurs de l'indice de queue de moment harmonique*) soit X_1, \dots, X_n suite de variable aléatoire i.i.d avec F la distributions de paréto

$$H^{(\beta)} = \frac{1}{\beta - 1} \left\{ \left[k^{-1} \sum_{i=1}^n \left(\frac{X_{n-k,n}}{X_{n-i+1,n}} \right)^{\beta-1} \right] - 1 \right\}$$

où $1 \leq k \leq n - 1$ et $\beta > 0$ est paramètre de réglage .pour $\beta = 1$, $H_{n,k}^{(\beta)}$ est interprété comme limite pour $\beta \rightarrow 1$,i.e

$$H_{n,k}^{(1)} = \lim_{\beta \rightarrow 1} H_{n,k}^{(\beta)} = k^{-1} \sum_{i=1}^k \log(X_{n-i+1,n}/X_{n-k,n})$$

dans cette parté nous nous proposons d'étudier un nouvel estimateur de l'indice de valeur extrême conditinelle en présence de données complete et de données censuré. Pour cela nous'allons pas utiliser toutes les données mais simplement les observations dont leurs covariables sont proches de x Ainsi, l'approche que nous allons utiliser est la méthode de la **fenêtre mobile** comme dans Gardes et Girard (2008,), définie par une boule, $B(x, h_{n,x})$ où x est le centre de la boule et $h_{n,x}$ est le rayon de la boule :

$$B(x, h_{n,x}) = \{t \in X, d(x, t) \leq h_{n,x}\}$$

Cette méthode renferme trois étapes :

- **Etape 1** : On fixe une covariable x déterministe, centre de la boule,
- **Etape 2** : Puis on fixe le rayon, $h_{n,x}$,
- **Etape 3** : Enfn on sélectionne uniquement les observations (y_i) pour lesquelles $x_i \in B(x, h_{n,x})$:

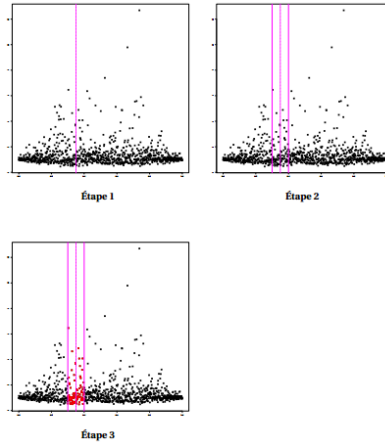


FIG. 2.1 – Les différentes étapes de sélection des données. En ordonnée on a la variable Y et en abscisse la covariable X.

2.4.1 Sélection de l'échantillon conditionnelle

Soit $Y \in R$ une variable aléatoire associée à une covariable non-aléatoire $x \in E$, où E désigne un espace métrique, non nécessairement de dimension finie. Soient $\{(X_i, Y_i), i = 1, \dots, n\}$ des observations indépendantes du couple $(x, Y) \in E \times R$

Nous nous plaçons dans le domaine de Fréchet où les queues de distribution sont lourdes, la fonction de distribution conditionnelle de Y sachant X se modélise alors de la façon suivante

$$F(y, x) = 1 - x^{-1/\gamma(x)} L(y, x)$$

où $\gamma(x)$ est de paramètre fonctionnel positif de la covariable x fixé, $\ell(\cdot, x)$ est une fonction à variations lentes qui vérifie, pour $\lambda > 0$

$$\lim_{n \rightarrow \infty} \frac{L(\lambda y, x)}{L(\lambda, x)}$$

nous utilisons une méthode d'estimation dite de la fenêtre mobile pour construire nos estimateurs. Pour cela, on introduit une boule centrée en x , de rayon $h_{n,x} > 0$, notée $B(x, h_{n,x})$ et définie par :

$$B(x, h_{n,x}) = \{t \in X, d(x, t) \leq h_{n,x}\}$$

Étant donné $(h_{n,x}, x)_{n \geq 1}$ une suite positive convergeant vers zéro quand n tend vers l'infini, on se propose de sélectionner que les observations Y_i pour lesquelles les covariables x_i sont dans la boule $B(x, h_{n,x})$. La proportion de tels points est ainsi donnée par :

$$\varphi(h_{n,x}) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{B(x, h_{n,x})}(x_i)$$

Le résultat de l'observation sample dans $(0, \infty) \times B(x, h_{n,x})$ donné par

$$\{Z_i(x), i = 1, \dots, m_{n,x}\}$$

où $m_{n,x} = n\varphi(h_{n,x})$ et $Z_{1,m_{n,x}} \leq \dots \leq Z_{m_{n,x},m_{n,x}}$ les statistiques ordonnées correspondantes

2.4.2 Estimateur pour des données complètes

pour le même échantillon considéré dans la sous-section ci-dessus où $l(y, x) = 1$ et pour un seuil arbitraire $z_0 > 1$ soit

$$T = \frac{Z}{z_0} \mathbb{I}_{\{Z > z_0\}}$$

puis conditionnellement sur $Z > z_0$,

$$p(T \leq t / Z > z_0) = 1 - t^{-1/\gamma(x)}$$

De plus conditionnellement $Z > z_0$, on a

$$T^{-1/\gamma(x)} = 1 - F(T) \stackrel{d}{=} U$$

où U uniformément distribué sur $[0, 1]$. Étant donné une séquence iid $T_1 \dots T_k$ des dépassements relatifs sur un seuil z_0 , la loi forte des grands nombres implique que

$$\frac{1}{k_{n,x}} \sum_{i=1}^{k_{n,x}} T_i^{1-\beta} = \frac{1}{k_{n,x}} \sum_{i=1}^{k_{n,x}} \left(T_i^{-1/\gamma(x)}\right)^{\gamma(x)(\beta-1)} \stackrel{d}{=} \frac{1}{k_{n,x}} \sum_{i=1}^{k_{n,x}} U_i^{\gamma(x)(\beta-1)}$$

$$\xrightarrow{\alpha.s} E[U^{\gamma(x)(\beta-1)}] = \int_0^1 u^{\gamma(x)(\beta-1)} du = \frac{1}{\gamma(x)(\beta-1)+1}$$

à condition que $\gamma(x)(\beta-1) + 1 \neq 0$

$$\frac{1}{\beta - 1} \left(\frac{1}{k \sum_{i=1}^k T_i^{1-\beta}} - 1 \right) \xrightarrow{\alpha.s} \gamma(x)$$

$$\hat{\gamma} = \frac{1}{\beta - 1} \left\{ \frac{1}{k_{n,x}} \sum_{i=1}^{k_{n,x}} \left[\frac{Z_{m_{n,x}-k_{n,x},m_{n,x}}(x)}{Z_{m_{n,x}-i+1,m_{n,x}}(x)} \right]^{\beta-1} - 1 \right\} \quad (2.8)$$

où $1 \leq k \leq m_{n,x} - 1$ et $\beta > 0$ est paramètre de réglage.

2.4.3 Estimateur pour des données censuré

Modèle et notation

On considère l'espace probabilisé (Ω, A, P) et l'espace métrique (x, d) , soient l'échantillon $(x_i)_{1 \leq i \leq n}$ un échantillon déterministe de X , $(Y_i)_{1 \leq i \leq n}$ et $(C_i)_{1 \leq i \leq n}$ réalisations respectives des variables aléatoires Y et C dans l'espace (Ω, A, P) . Nous supposons que Y et C ont respectivement comme fonctions de répartition conditionnelles F et G . soit $(Z_i)_{1 \leq i \leq n} = (Y_i \wedge C_i)_{1 \leq i \leq n}$ et $\delta_i = \mathbb{I}_{\{Y_i \leq C_i\}}$, où $\mathbb{I}_{\{\cdot\}}$ est une fonction indicatrice. Nous supposons aussi que Y et C sont indépendantes conditionnellement à X . Nous observons, $(z_1, \delta_1, x_1) \dots (z_n, \delta_n, x_n)$ des copies indépendantes de (Z, δ, X) . Dans la suite nous supposons que la variable Z a comme fonction de répartition conditionnelle H . nous considérons que les fonctions de répartition conditionnelles de Y et C pour tout $x \in X$ sont ainsi définies :

$$F(u | x) = 1 - u^{-1/\gamma_1(x)} L_1(u, x)$$

$$G(u | x) = 1 - u^{-1/\gamma_2(x)} L_2(u, x)$$

où $\gamma_1(\cdot)$ et $\gamma_2(\cdot)$ sont des paramètres fonctionnels positifs de la covariable x fixé appelés indices de queue conditionnels ou indices des valeurs extrêmes conditionnels, $L_1(\cdot, x)$ et $L_2(\cdot, x)$ sont des fonctions à variations lentes. pour tout $\lambda > 0$,

$$\lim_{u \rightarrow \infty} \frac{L_i(\lambda u, x)}{L_i(u, x)} = 1, i = 1, 2.$$

Si Y et C sont indépendantes, alors la fonction de répartition conditionnelle $H(\cdot | x)$ de Z pour $X = x$ fixé, $\gamma(x) = \gamma_1(x)\gamma_2(x) / (\gamma_1(x) + \gamma_2(x))$. Ce qui donne pour tout u et x

$$1 - H(u | x) = (1 - F(u | x))(1 - G(u | x)).$$

$$= u^{-1/\gamma_1(x)} L_1(u, x) u^{-1/\gamma_2(x)} L_2(u, x)$$

$$u^{-\left(\frac{1}{\gamma_1(x)} + \frac{1}{\gamma_2(x)}\right)} L_1(u, x) L_2(u, x)$$

$$u^{-\left(\frac{\gamma_1(x) + \gamma_2(x)}{\gamma_1(x)\gamma_2(x)}\right)} L(u, x)$$

$$u^{-1/\gamma(x)} L(u, x)$$

où $\gamma(x) = (\gamma_1(x)\gamma_2(x))/(\gamma_1(x)+\gamma_2(x))$ et $L(x) = L_1(x) L_2(x)$. Donc H est une fonction de répartition appartenant au domaine d'attraction de Fréchet

$$1 - H(u | x) \in RV(-1/\gamma(x)), \quad \text{avec } \gamma(x) = (\gamma_1(x)\gamma_2(x))/((\gamma_1(x) + \gamma_2(x)) .$$

pour obtenu le niveau l'estimateure pour de donne consuré on devisé 2.8 sur p_x

$$\hat{\gamma}_{X,k,n}^{(c,.)} = \frac{\hat{\gamma}_{k_x m_n, x}^{(.)}}{\hat{p}_x}$$

ou

$$\hat{p}_x = \frac{1}{k_x} \sum_{i=1}^{k_x} \delta_{[m_n, x-i+1]}^x \quad \text{pour la covariable } x \text{ fixe}$$

$$\hat{p}_{t_n}(x) = \frac{\overline{H}_n^1(t_n | x)}{\overline{H}(t_n | x)} \text{ pour la covariable } x \text{ aléatoire}$$

Chapitre 3

Etude de simulation

Dans ce chapitre nous considérons la simulation de notre nouveau estimateur $\hat{\gamma}$ dans le cas des données complètes et censurées et nous appliquons la méthode du bootstrap pour une étude asymptotique

3.1 introduction du logiciel R

3.1.1 Qu'est-ce que le logiciel R ?

R est un logiciel permettant de faire des analyses statistiques et de produire des graphiques. Mais R est également un langage de programmation complet, c'est cet aspect qui fait que R est différent des autres logiciels statistiques. Les informations sur R sont disponibles sur la homepage du projet

<http://www.r-project.org/>

R est un clone gratuit du logiciel S-Plus commercialisé par MathSoft et développé par Statistical Sciences autour du langage

3.1.2 les avantages de logiciel R

- permet de faire des graphiques très flexibles et d'une qualité exceptionnelle
- permet de construire facilement vos propres fonctions
- est facile à interfacer avec d'autres langages
- est aussi un logiciel mathématique (calcul matriciel, intégration numérique, optimisation, ..)
- les graphes sont bien plus jolis
- R est très facile à utiliser
- R est gratuit

3.2 Princip du Bootstrap

Le principe fondamental de cette technique de ré-échantillonnage est de substituer à la distribution de probabilité inconnue F , dont est issu l'échantillon d'apprentissage, la distribution empirique F_n qui donne un poids $\frac{1}{n}$ à chaque réalisation. Ainsi on obtient un échantillon de taille n dit échantillon bootstrap selon la distribution empirique F_n par n tirages aléatoires avec remise parmi les n observations initiales. Si n est grand, la distribution empirique \hat{F}_n est proche de F , on aura donc une bonne approximation de la loi de X en utilisant \hat{F}_n à la place de F . La méthode de bootstrap consiste à construire un nombre B (B entier)

d'échantillons bootstrap notée X^* (images de l'échantillon initial), afin de les utiliser pour faire des inférences, et plus le nombre d'images simulées est grand, plus la statistique est précise, pour chaque nouvel échantillon, on calcule de la même façon un nouvel estimateur (image simulée de l'estimateur initial). L'ensemble des images simulées de l'estimateur initial est considéré comme un modèle de sa distribution sur la population de l'échantillon initial.

3.3 Simulation pour des données complètes

3.3.1 Échantillon initial et paramètres de simulation

La loi de simulation utilisée dans ce cas est une loi de Pareto de paramètre γ de fonction de répartition conditionnel,

$$F(y) = 1 - y^{-1/\gamma(x)}$$

Nous avons généré un échantillon $(Y_i)_{1 \leq i \leq n} \sim \text{Pareto}(\gamma_1)$ de taille $n = 1000$, à partir d'une variable u de $U([0, 1])$, le modèle ajusté sera :

$$F^{-1}(u) = (1 - u)^{-\gamma_1}$$

les covariables $(X_i)_{1 \leq i \leq n} \sim \text{Pareto}(\gamma_2)$ de taille $n = 1000$, à partir d'une variable v de $v([0, 1])$, le modèle ajusté sera :

$$H^{-1}(v) = (1 - v)^{-\gamma_2}$$

Par la méthode de fenêtre mobile on obtient l'échantillon conditionnel Z

3.3.2 Sélectionnement de l'échantillon conditionnel

pour sélectionner l'échantillon conditionnel on utilise la méthode de fenêtre mobile

Algorithme 1

1. Nous avons généré un échantillon Y et une covariable X suit la loi de Pareto de taille $n=1000$ de paramètre $\gamma_1 = 1.5$ (resp $\gamma_2 = 0.6$), à partir d'une variable u de $U([0; 1])$ (resp v de $U([0; 1])$), le modèle ajusté sera :

$$F^{-1}(u) = 1 - u^{-1/\gamma}$$

$$H^{-1}(v) = 1 - v^{-1/\gamma}$$

1. On fixe une valeur x_1 déterministe, à partir de la covariable x
2. Puis on choisit les variables x_i qui vérifient

$$d(x, x_1) \leq h_{n,x} \tag{3.1}$$

où $h_{n,x} = 0.3$

et on détermine la taille de x_i , $m_{n,x} = n\varphi(h_{n,x})$ où

$$\varphi(h_{n,x}) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{B(x, h_{n,x})}(x_i)$$

Enfin on sélectionne uniquement les observations (y_i) pour lesquelles x_i vérifient la condition 3.1

Programme sous R

```
n=?
gamma1=?
gamma2=0.6
u<-runif(n)
v<-runif(n)
y<-(1-u)^(-gamma1)
x<-(1-v)^(-gamma2)
x1<-sample(x,1)
d<-numeric(n)
b<-numeric(n)
h<-(0.3)
for (i in 1 :n)
{
d[i]<-abs(x[i]-x1)
if (d[i]<h)
b[i]<-1 else b[i]<-0
}
Q<-(sum(b)/n)
m<-(n*Q)
z<-numeric(m)
for (i in 1 :n)
{
d[i]<-abs(x[i]-x1)
if (d[i]<h)
z[i]<-y[i] else z[i]<-0
}
z
z<-z[z!=0]
```

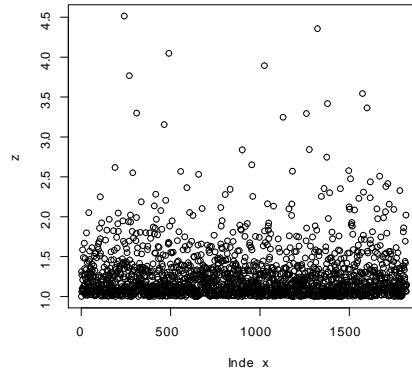


FIG. 3.1 – Échantillon conditionel Z pour ($\gamma_1 = 0.2$) et ($\gamma_2 = 0.8$)

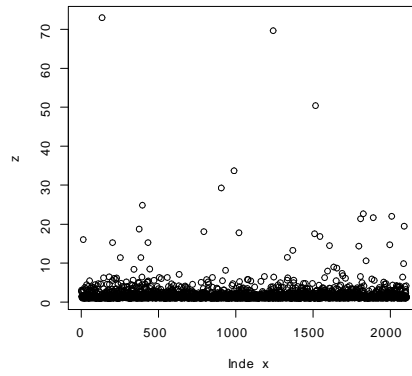


FIG. 3.2 – Échantillon conditionel Z pour ($\gamma_1 = 0.5$) et ($\gamma_2 = 0.8$)

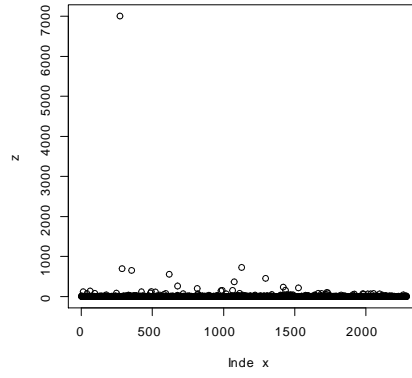


FIG. 3.3 – Échantillon conditionel Z pour ($\gamma_1 = 1$) et ($\gamma_2 = 0.8$)

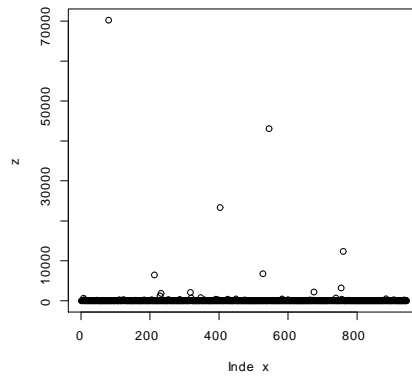


FIG. 3.4 – Échantillon conditionel Z pour ($\gamma_1 = 1.5$) et ($\gamma_2 = 0.8$)

3.3.3 Simulation de nouvel estimateur γ^h en fonction k

Algorithme2

1. on sélectionne l'échantillon conditionnel z pour la même méthode du l'algorithme1
2. on fixe $beta = 2$ et k varé de 1 á $m - 1$, où m la taill de l' échantillon Z , et en simuler l'estimateur de queue $gama[k]$ où

$$gama[k] = \frac{1}{beta - 1} \left\{ \frac{1}{k} \sum_{i=1}^k \left[\frac{Z_{m_{n,x}-k_{n,x},m_{n,x}}(x)}{Z_{m_{n,x}-i+1,m_{n,x}}(x)} \right]^{beta-1} - 1 \right\}$$

3. Nous dessinons la courbe de gama en fonction "k"
4. Nous choisissons la valeur appropriée de "k" pour le corépondant a la mieur valeur du gama que l'on note k_{opt}

Programme sous R

#paramétra de simulation

n=?

gamma=0.5

gamma2=0.8

#échantillon initial

u<-runif(n)

v<-runif(n)

y<-(1-u)^(-gamma1)

x<-(1-v)^(-gamma2)

y<-sort(y)

x1<-sample(x,1)

d<-numeric(n)

b<-numeric(n)

h<-(0.3)

for (i in 1 :n)

{

d[i]<-abs(x[i]-x1)

if (d[i]<h)

b[i]<-1 else b[i]<-0

}

Q<-(sum(b)/n)

m<-(n*Q)

z<-numeric(m)

for (i in 1 :n)

{

d[i]<-abs(x[i]-x1)

if (d[i]<h)

z[i]<-y[i] else z[i]<-0

}

z

z<-z[z!=0]

z

z<-sort(z)

```

z
#"graphe gama vs k"
beta=?
gama<-numeric(m-1)
for (k in 1 :m-1)
{
i<-seq(1,k)
gama[k]<-abs((1/(beta-1))*((1/k)*sum((z[m-k]/z[m-i+1])^(beta-1))-1))
}
k<-seq(1,m-1)
plot(k,gama[k],col="red",type="l",lwd="2",main="graphe de estimateure"
abline(h=gamma,col="blue",lwd="2")
# k optimal graphique
kopt<-k[which.min(abs(gamma-gama))]
kopt
gamaopt<-gama[kopt]
gamaopt

```

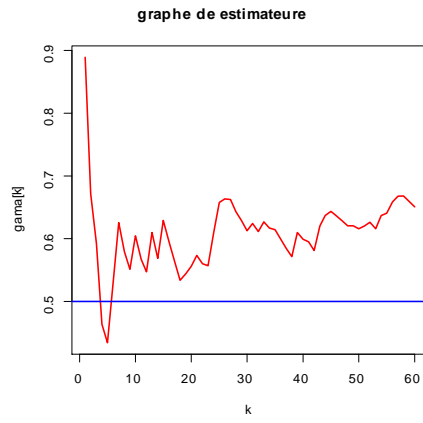


FIG. 3.5 – Comportement graphique de $\hat{\gamma}$ vs k , pour $n=200$ et $\beta=2$

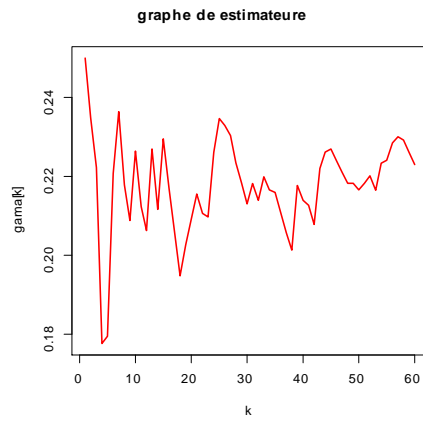


FIG. 3.6 – Comportement graphique de $\hat{\gamma}$ vs k , pour $n=200$ et $\beta=5$

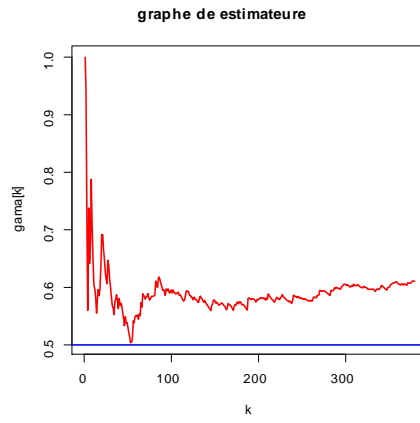


FIG. 3.7 – Comportement graphique de $\hat{\gamma}$ vs k , pour $n=200$ et $\beta=10$

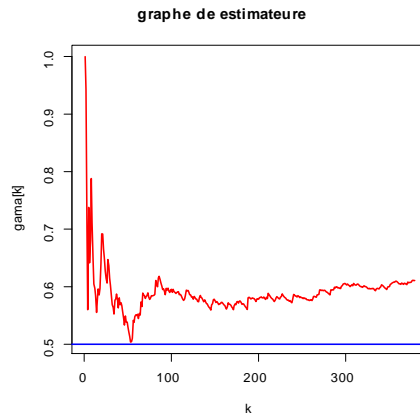


FIG. 3.8 – Comportement graphique de $\hat{\gamma}$ vs k , pour $n=1000$ et $\beta=2$

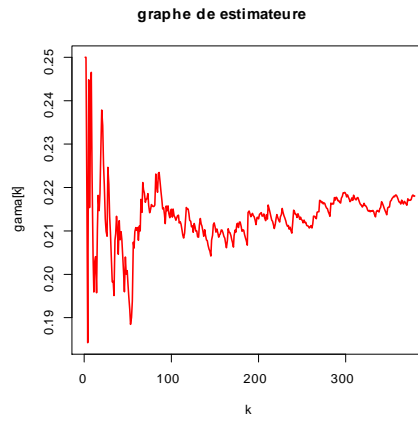


FIG. 3.9 – Comportement graphique de $\hat{\gamma}$ vs k , pour $n=1000$ et $\beta=5$

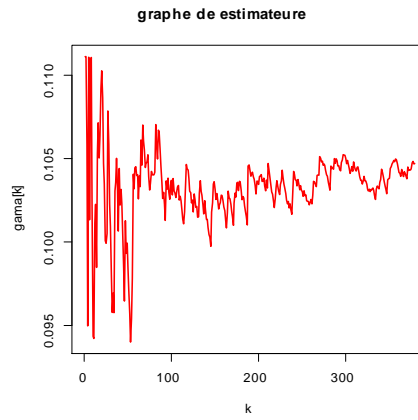


FIG. 3.10 – Comportement graphique de $\hat{\gamma}$ vs k , pour $n=1000$ et $\beta=10$

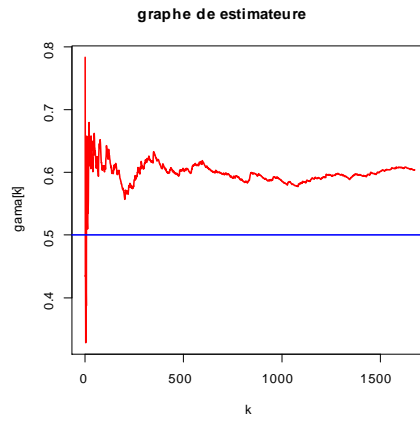


FIG. 3.11 – Comportement graphique de $\hat{\gamma}$ vs k , pour $n=5000$ et $\beta=2$

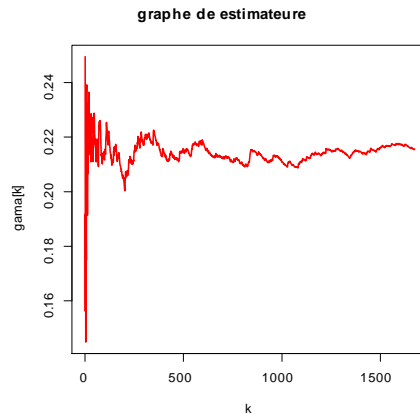


FIG. 3.12 – Comportement graphique de $\hat{\gamma}$ vs k , pour $n=5000$ et $\beta=5$

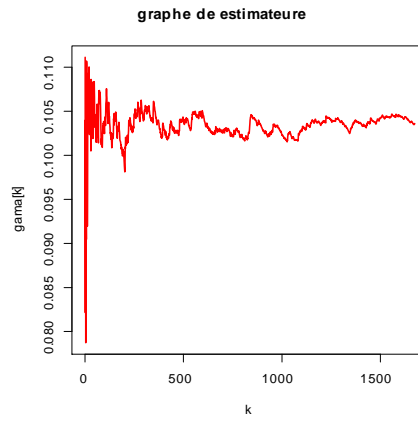


FIG. 3.13 – Comportement graphique de $\hat{\gamma}$ vs k , pour $n=5000$ et $\beta=10$

3.3.4 En ploit de bootsrap pour de $\hat{\gamma}$

Algorithme3

1. On sélectionne l'échantillon conditionnel de la même méthode que l'algorithme1
2. On choisit le valeur de $kopt$ de le même méthode que de l'algorithme2
3. On simule l'estimateur de queue $gama1$ inétial pour $kopt$

$$gama1 = \frac{1}{beta - 1} \left\{ \frac{1}{kopt} \sum_{i=1}^{kopt} \left[\frac{Z_{m_{n,x}-k_{n,x},m_{n,x}}(x)}{Z_{m_{n,x}-i+1,m_{n,x}}(x)} \right]^{beta-1} - 1 \right\}$$

4. On crée B échantillons indépendants $Z^{*1}, Z^{*2}, \dots, Z^{*B}$ où chaque échantillon Z^{*b} est obtenu en tirant n observations avec remise dans l'échantillon $Zboot = \{Zboot_1, \dots, Zboot_n\}$
5. Pour chaque échantillon b tel que, $1 \leq b \leq B$, on calcule l'estimateur de queue $:gmma[b]$,

$$gama[b] = \frac{1}{beta - 1} \left\{ \frac{1}{kopt} \sum_{i=1}^{kopt} \left[\frac{Zboot_{m_{n,x}-k_{n,x},m_{n,x}}(x)}{Zboot_{m_{n,x}-i+1,m_{n,x}}(x)} \right]^{beta-1} - 1 \right\}$$

On obtient alors un échantillon de B valeurs de l'estimateur de queue $gama$

6. On estime alors la moyenne $E(gama)$ et l'erreur standard $se_F(gama)$ et $QQnorme$ et du $biais$

Programme sous R

```
#paramétra de simulation
n=1000
gamma1=1.5
gamma2=0.6
beta=2
#échantillon initial
u<-runif(n)
v<-runif(n)
y<-(1-u)^(-gamma1)
x<-(1-v)^(-gamma2)
x1<-sample(x,1)
d<-numeric(n)
b<-numeric(n)
h<-(0.3)
for (i in 1 :n)
{
d[i]<-abs(x[i]-x1)
if (d[i]<h)
b[i]<-1 else b[i]<-0
}
Q<-(sum(b)/n)
m<-(n*Q)
z<-numeric(m)
for (i in 1 :n)
{
```

```

d[i]<-abs(x[i]-x1)
if (d[i]<h)
z[i]<-y[i] else z[i]<-0
}
z
z<-z[z!=0]
z
z<-(sort(z))
z
#choit kopt et gama inétial
gama<-numeric(m-1)
for (k in 1 :m-1)
{
i<-seq(1,k)
gama[k]<-abs((1/(beta-1))*((1/k)*sum((z[m-k]/z[m-i+1])^(beta-1))-1))
}
kopt<-which.min(abs(gamma-gama))
kopt
gama1<-gama[kopt]
#gama bootstrap
r<-2000
gama<-numeric(r)
for (b in 1 :r)
{
i<-seq(1,kopt)
zboot<-sample(z,length(z),TRUE)
zboot<-(sort(zboot))
gama[b]<-abs((1/(beta-1))*((1/kopt)*sum((zboot[m-kopt]/zboot[m-i+1])^(beta-1))-1))
}
#moyenne gama
meangama<-mean(gama)
meangama
#graph de probabilié
qqnorm(gama)
#estimation bootstrap de l'ereure standar
sd<-sd(gama)
sd
#Estémation bootstrap de bias
bias<-mean(gama)-gama1

```

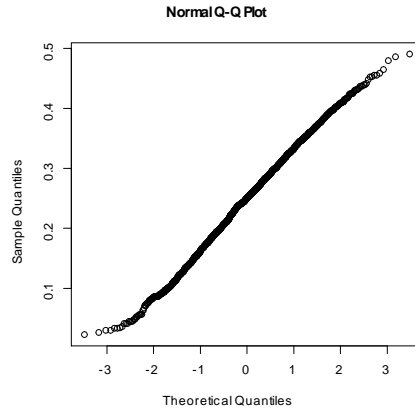


FIG. 3.14 – QQ-norm de la distribution limite bootstrap de 2000 répétition de Paréto ($\gamma_1 = 0.2$), $n = 1000$

	$R = 2000$	$R = 1000$	$R = 500$
k_{opt}	6	6	6
$\hat{\gamma}_1^{(h)}$	0.3248784	0.3248784	0.3248784
$\hat{\gamma}_{1^{boot}}^{(h)}$	0.1403471	0.2132133	0.1830524
$E(\hat{\gamma}^{(h)})$	0.2482064	0.246846	0.2481463
sd	0.8330169	0.08091615	0.08562753
$biais$	-0.07667207	-0.07803245	-0.07673214

TAB. 3.1 – Résultats de simulation pour $n=1000$

3.4 Simulation pour des données censurées

3.4.1 Échantillon initial et paramètres de simulation

La loi de simulation utilisée dans ce cas est une loi de Pareto de paramètre γ de fonction de répartition conditionnel,

$$F(y) = 1 - y^{-1/\gamma(x)}$$

les covariables $(X_i)_{1 \leq i \leq n} \sim \text{Pareto}(\gamma_3)$ de taille $n = 1000$, à partir d'une variable w de $w \in [0, 1]$, le modèle ajusté sera :

$$H^{-1}(w) = (1 - w)^{-\gamma_3}$$

Nous avons généré un échantillon $(Z_{1_i})_{1 \leq i \leq n} \sim \text{Pareto}(\gamma_1)$ de taille $n = 1000$, à partir d'une variable u de $U([0, 1])$, le modèle ajusté sera :

$$F^{-1}(u) = (1 - u)^{-\gamma_1}$$

L'échantillon $(Z_{1_i})_{1 \leq i \leq n}$ est censuré par un deuxième échantillon $(Z_{2_i})_{1 \leq i \leq n} \sim \text{Pareto}(\gamma_2)$ à partir d'une variable v de $U([0, 1])$:

$$G^{-1}(v) = (1 - v)^{\gamma_2}$$

Les variables que nous observons sont d'une part les $Z_i \sim \text{Pareto}(\gamma)$ définies par :

$$Y_i = Z_{1_i} \wedge Z_{2_i}$$

les indicateurs de censure sont,

$$\delta_i = \mathbb{I}_{\{Z_{1_i} \leq Z_{2_i}\}}$$

3.4.2 Sélectionnement de l'échantillon conditionnel censure

Algorithme4

1. Nous avons généré deux échantillon y_1 et y_2 suivant la loi de Pareto de taille $n=100$ de paramètre $\gamma_1 = 1.5$ (resp $\gamma_2 = 0.6$), à partir d'une variable u de $U([0; 1])$ (resp v de $U([0; 1])$), et une covariable x suivant la loi de Pareto de taille $n=100$ de paramètre $\gamma_3 = 0.8$, à partir d'une variable w de $w \in [0; 1]$, le modèle ajusté sera :

$$F^{-1}(u) = 1 - u^{-1/\gamma_1}$$

$$F^{-1}(v) = 1 - v^{-1/\gamma_2}$$

$$F^{-1}(w) = 1 - w^{-1/\gamma_3}$$

1. on sélectionne deux échantillons conditionnels z_1 et z_2 de même covariable que x , pour cela on utilise la même méthode de l'algorithme1
- 3 On sélectionne l'échantillon initial où $Y_i = Z_{1_i} \wedge Z_{2_i}$
- 4 On définit δ où $\delta_i = \mathbb{I}_{\{Z_{1_i} \leq Z_{2_i}\}}$
- 5 On ordonne les données pour obtenir l'échantillon censuré $Y = \{(Y_1, \delta_1), \dots, (Y_n, \delta_n)\}$,

Programme sous R

#paramètre de simulation

n=1000

gamma1=1.5

gamma2=0.6

gamma3=0.8

u<-runif(n)

v<-runif(n)

w<-runif(n)

y1<-(1-u)^(-gamma1)

x<-(1-v)^(-gamma2)

y2<-(1-w)^(-gamma2)

#sélectionner l'échantion condetionel z

x1<-sample(x,1)

d<-numeric(n)

b<-numeric(n)

h<-(0.3)

for (i in 1 :n)

{

d[i]<-abs(x[i]-x1)

if (d[i]<h)

b[i]<-1 else b[i]<-0

}

Q<-(sum(b)/n)

m<-(n*Q)

#selectioner l'échantillon z1

z1<-numeric(m)

for (i in 1 :n)

{

d[i]<-abs(x[i]-x1)

if (d[i]<h)

z1[i]<-y1[i] else z1[i]<-0

}

z1

z1<-z1[z1!=0]

z1

seléctoner l'échantillon z2

z2<-numeric(m)

for (i in 1 :n)

{

d[i]<-abs(x[i]-x1)

if (d[i]<h)

z2[i]<-y2[i] else z2[i]<-0

}

z2

z2<-z2[z2!=0]

z2

#sélectionner l'échantillon initial

delta<-numeric(m)

```

y<-pmin(z1,z2)
for (i in 1 :m)
{
if (z1[i]<=z2[i])
delta[i]<-1 else delta[i]<-0
}
delta
r<-rank(y)
#ordonner des donner
y<-sort(y)
r1<-sort(r)
L<-numeric(length(y))
B<-numeric(length(y))
Delta<-numeric(length(y))
for (i in 1 :length(y))
{
L[i]<-(delta[i]+length(y))*r[i]
}
L
A<-sort(L)
A
for (i in 1 :length(y))
{
B[i]<-A[i]/r1[i]
}
B
Delta<-B-length(y)
Delta
plot(y

```

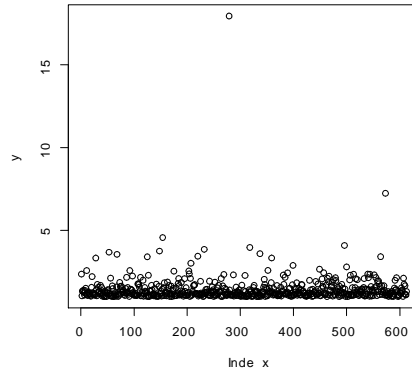


FIG. 3.15 – Échantillon conditionel Y, ($\gamma_1 = 0.35$) et ($\gamma_2 = 2.5$) (10% censurées)

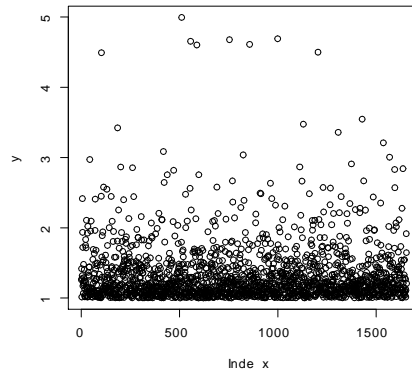


FIG. 3.16 – Échantillon conditionel Y, ($\gamma_1 = 0.35$) et ($\gamma_2 = 1$) (25% censurées)

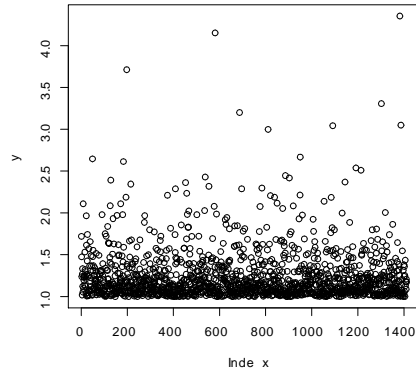


FIG. 3.17 – Échantillon conditionel Y , ($\gamma_1 = 0.35$) et ($\gamma_2 = 0.5$) (40% censurées)

3.4.3 Simulation de nouvel estimateur de l'indice en fonction de k

Algorithme5

1. On sélectionne l'échantillon conditionnel censuré Y de la même méthode que l'algorithme4
2. On fixe $beta = 2$ et k vari de 1 á $m - 1$, où m la taill de léchantillon Y , et on simule l'estimateur de queue $gama[k]$ où

$$gama[k] = \left(\frac{1}{beta - 1} \left\{ \frac{1}{k} \sum_{i=1}^k \left[\frac{Y_{m_n, x-k_n, x, m_n, x}(x)}{Y_{m_n, x-i+1, m_n, x}(x)} \right]^{beta-1} - 1 \right\} \right) / P(x)$$

où

$$P(x) = \frac{1}{kopt} \sum_{i=1}^{kopt} \delta_{[m_n, x-i+1]}^x$$

3. Nous dessinons la courbe de gama en fonction " k "
4. Nous choisissons la valeur appropriée de " k " pour la meilleure valeur de gama corespondante gama et on noté $kopt$

Programme sous R

```
#paramètre de simulation
n=1000
gamma1=0.35
gamma2=
gamma3=0.8
u<-runif(n)
v<-runif(n)
w<-runif(n)
y1<-(1-u)^(-gamma1)
x<-(1-v)^(-gamma2)
y2<-(1-w)^(-gamma2)
#sélectionner l'écheantion condetionel z
x1<-sample(x,1)
d<-numeric(n)
b<-numeric(n)
h<-(0.3)
for (i in 1 :n)
{
d[i]<-abs(x[i]-x1)
if (d[i]<h)
b[i]<-1 else b[i]<-0
}
Q<-(sum(b)/n)
m<-(n*Q)
#sélectionner l'échantillon z1
z1<-numeric(m)
for (i in 1 :n)
{
```

```

d[i]<-abs(x[i]-x1)
if (d[i]<h)
z1[i]<-y1[i] else z1[i]<-0
}
z1
z1<-z1[z1 !=0]
z1
#sélectionner l'échantillon z2
z2<-numeric(m)
for (i in 1 :n)
{
d[i]<-abs(x[i]-x1)
if (d[i]<h)
z2[i]<-y2[i] else z2[i]<-0
}
z2
z2<-z2[z2 !=0]
z2
#sélectionner l'échantillon initial
delta<-numeric(m)
y<-pmin(z1,z2)
plot(y)
for (i in 1 :m)
{
if (z1[i]<=z2[i])
delta[i]<-1 else delta[i]<-0
}
delta
r<-rank(y)
#ordonner des donner
y<-sort(y)
r1<-sort(r)
L<-numeric(length(y))
B<-numeric(length(y))
Delta<-numeric(length(y))
for (i in 1 :length(y))
{
L[i]<-(delta[i]+length(y))*r[i]
}
L
A<-sort(L)
A
for (i in 1 :length(y))
{
B[i]<-A[i]/r1[i]
}
B
Delta<-B-length(y)
Delta
#"graphe gama vs k"

```

```

beta<-2
gama<-numeric(m-1)
for (k in 1 :m-1)
{
i<-seq(1,k)
gama[k]<-abs((1/(beta-1))*((1/k)*sum((y[m-k]/y[m-i+1])^(beta-1))-1))/((1/k)*sum(Delta[m-
i+1]))
}
k<-seq(1,m-1)
plot(k,gama[k],col="red",type="l",lwd="2",main="graphe de estimateur" )
abline(h=gama1,col="blue",lwd="2")
# k optimal graphique
kopt<-k[which.min(abs(gama-gama1))]
kopt
gamaopt<-gama[kopt]
gamaopt

```

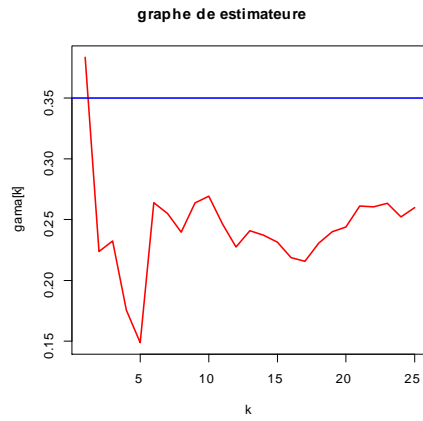


FIG. 3.18 – Comportement graphique de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 2.5$), (10% de censure), pour $n=200$

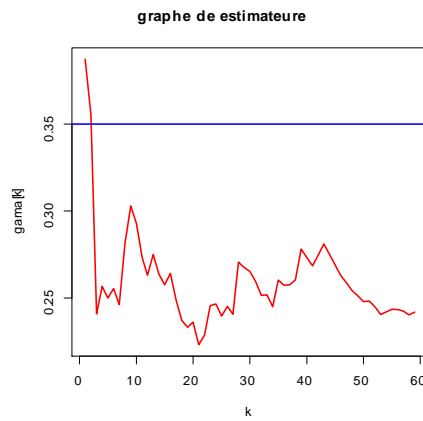


FIG. 3.19 – Comportement graphique de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 2.5$), (10% de censure), pour $n=1000$

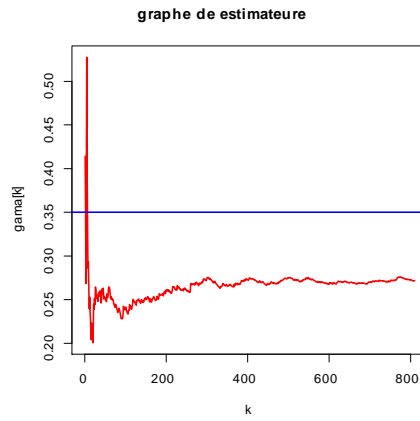


FIG. 3.20 – Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 2.5$), (10% de censure), pour n=5000

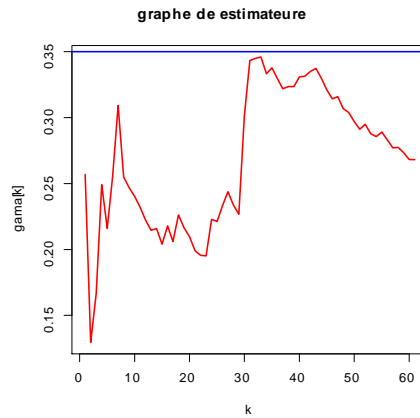


FIG. 3.21 – Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 1$), (25% de censure), pour n=200

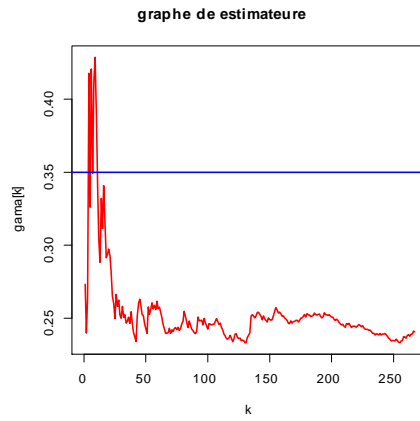


FIG. 3.22 – Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 1$), (25% de censure), pour $n=1000$

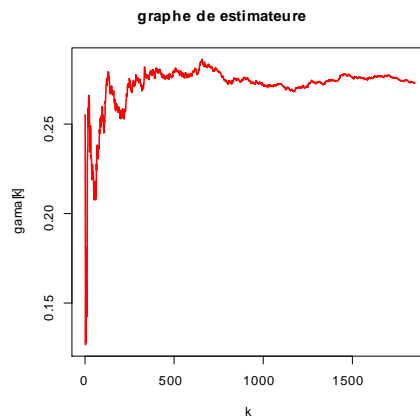


FIG. 3.23 – Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 1$), (25% de censure), pour $n=5000$

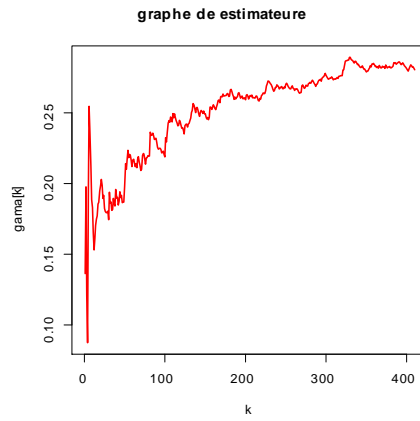


FIG. 3.24 – Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 0.5$), (40% de censure), pour n=200

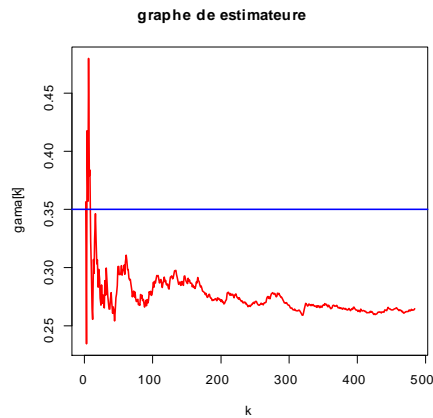


FIG. 3.25 – Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 0.5$), (40% de censure), pour n=1000

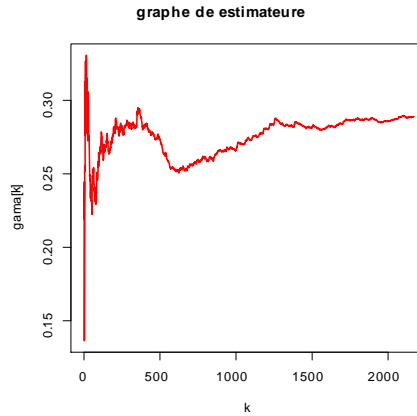


FIG. 3.26 – Comportement graphic de $\hat{\gamma}$ vs k issu de la distribution de pareto ($\gamma_1 = 0.35$) censurée par pareto ($\gamma_2 = 0.5$), (40% de censure), pour n=5000

Comontaire

3.4.4 emploi du bootsrap

Algorithme

1. On sélectionne l'échantillon conditionnel Y de la même méthode que l'algorithme 4
2. On choisit la valeur de k_{opt} et $gama1$ (estimateur initial) pour la même méthode que l'algorithme 5
3. On crée B échantillons indépendants $Z^{*1}, Z^{*2}, \dots, Z^{*B}$ où chaque échantillon Z^{*b} est obtenu en tirant n observations avec remise dans l'échantillon $Y_{boot} = \{(Y_{boot1}, \delta_1), \dots, (Y_{bootn}, \delta_n)\}$
4. Pour chaque échantillon b tel que, $1 \leq b \leq B$, on calcule l'estimateur de queue $gama[b]$,

$$gama[b] = \frac{1}{beta - 1} \left\{ \frac{1}{k_{opt}} \sum_{i=1}^{k_{opt}} \left[\frac{Y_{boot_{m_{n,x}-k_{n,x}, m_{n,x}}}(x)}{Y_{boot_{m_{n,x}-i+1, m_{n,x}}}(x)} \right]^{beta-1} - 1 \right\} / P(x)$$

$$P(x) = \frac{1}{k_{opt}} \sum_{i=1}^{k_{opt}} \delta_{[m_{n,x}-i+1]}^x$$

On obtient alors un échantillon de B valeurs de l'estimateur de queue $gama$

5. On estime alors la moyenne $E(gama)$ et l'erreur standard $se_F(gama)$ et $QQnorme$ et du biais

Programme sous R

```
#paramètre de simulation
#paramètre de simulation
n=100
gamma1=1.5
gamma2=0.6
gamma3=0.8
beta=2
#sélectionner l'échantillon conditionnel z
u<-runif(n)
```



```

v<-runif(n)
w<-runif(n)
y1<-(1-u)^(-gamma1)
x<-(1-v)^(-gamma2)
y2<-(1-w)^(-gamma2)
x1<-sample(x,1)
d<-numeric(n)
b<-numeric(n)
h<-(0.3)
for (i in 1 :n)
{
d[i]<-abs(x[i]-x1)
if (d[i]<h)
b[i]<-1 else b[i]<-0
}
Q<-(sum(b)/n)
m<-(n*Q)
#selectioner l'échantillon z1
z1<-numeric(m)
for (i in 1 :n)
{
d[i]<-abs(x[i]-x1)
if (d[i]<h)
z1[i]<-y1[i] else z1[i]<-0
}
z1
z1<-z1[z1!=0]
z1
sélectionner l'échantillon z2
z2<-numeric(m)
for (i in 1 :n)
{
d[i]<-abs(x[i]-x1)
if (d[i]<h)
z2[i]<-y2[i] else z2[i]<-0
}
z2
z2<-z2[z2!=0]
z2
#sélectionner l'échantillon initial
delta<-numeric(m)
y<-pmin(z1,z2)
for (i in 1 :m)
{
if (z1[i]<=z2[i])
delta[i]<-1 else delta[i]<-0
}
delta
r<-rank(y)
#ordonner des données

```

```

y<-(sort(y))
r1<-(sort(r))
L<-numeric(length(y))
B<-numeric(length(y))
Delta<-numeric(length(y))
for (i in 1 :length(y))
{
L[i]<-(delta[i]+length(y))*r[i]
}
L
A<-(sort(L))
A
for (i in 1 :length(y))
{
B[i]<-A[i]/r1[i]
}
B
Delta<-B-length(y)
Delta
#chiosir de kopt et gama inétial
beta<-2
for (k in 1 :m-1)
{
i<-seq(1,k)
gama[k]<-abs((1/(beta-1))*((1/k)*sum((y[m-k]/y[m-i+1])^(beta-1))-1))/((1/k)*sum(Delta[m-
i+1]))
}
kopt<-which.min(abs(gama1-gama))
kopt
gama1<-gama[kopt]
gama1
#gama bootstrap
R<-2000
gama<-numeric(R)
for (j in 1 :R)
{
indice <-sample(1 :length(y),length(y),replace=TRUE)
yboot<-y[indice]
deltaboot<-delta[indice]
yboot
deltaboot
r<-rank(y)
#ordonner des donner
yboot<-sort(yboot)
r1<-(sort(r))
L<-numeric(length(yboot))
B<-numeric(length(yboot))
Delta<-numeric(length(y))
for (i in 1 :length(yboot))
{

```

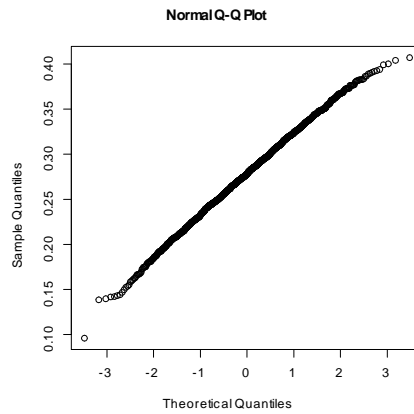
```

L[i]<-(deltaboot[i]+length(yboot))*r[i]
}
L
A<-(sort(L))
A
for (i in 1 :length(yboot))
{
B[i]<-A[i]/r1[i]
}
B
Deltaboot<-B-length(y)
Deltaboot
s<-seq(1,kopt)
gama[j]<-abs((1/(beta-1))*((1/kopt)*sum((yboot[m-kopt]/yboot[m-s+1])^(beta-1))-1))/((1/kopt)*sum
s+1]))
}
#moyenne gama
meangama<-mean(gama)
meangama
#graph de probabilié
qqnorm(gama)
#estimation bootstrap de l'ereure standar
sd<-sd(gama)
sd
#Estémation bootstrap de bias
bias<-mean(gama)-gama1
bias

```

	$c = 10\%$	$c = 25\%$	$c = 40\%$
k_{opt}	20	47	41
$\hat{\gamma}_1^{(h)}$	0.3035514	0.3376852	0.3495238
$\hat{\gamma}_1^{(h)}$	0.3073708	0.3397037	0.3374236
$E(\hat{\gamma}_1^{(h)})$	0.2779916	0.3258021	0.3414811
sd	0.04537116	0.03366893	0.04272734
$biais$	-0.02555979	-0.01188314	-0.008042705

TAB. 3.2 – Résultats de simulation pour $n=1000$



QQ-norm de la distribution
limite bootstrap de 2000
répétition de Paréto
($\gamma_1 = 0.35$), censurée
par ($\gamma_2 = 2.5$), (10% de
censure), $n=1000$

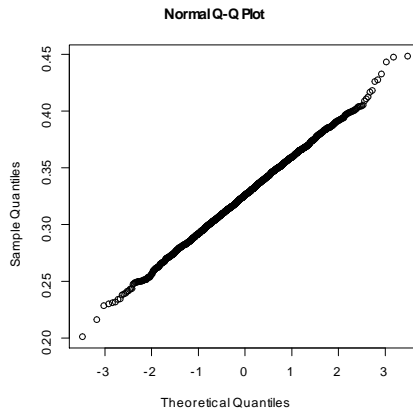


FIG. 3.27 – QQ-norm de la distribution limite bootstrap de 2000 répétition de Paréto ($\gamma_1 = 0.35$), censurée par ($\gamma_2 = 1$), (25% de censure), $n=1000$

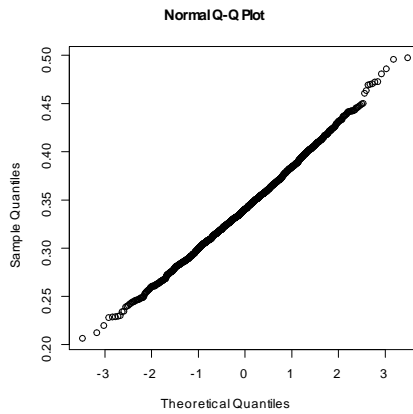


FIG. 3.28 – QQ-norm de la distribution limite bootstrap de 2000 répétition de Paréto ($\gamma_1 = 0.35$), censurée par ($\gamma_2 = 0.5$), (40% de censure), $n=1000$

Conclusion

Nous avons calculé dans notre travail un nouvel estimateur de l'indice des valeurs extrêmes, en adaptant un estimateur harmonique au cas des données conditionnel, par la méthode du fenêtre mobile, pour une distribution à queue lourde dans le cas du domaine d'attraction de Fréchet.

Notre nouvel estimateur présente une forme bien explicite, q on a simulé en fixant certain paramètres initiaux.

Notre estimateur à été présenté en deux versions, dans le cas des données complètes ainsi que des données censurées

Nous avons consacré deux objectifs dans nos études de simulation dans la première partie en étudiant l'estimateur en fonction du nombre de valeurs extrêmes, ensuite en variant le paramètre de réglage pour avoir une idée sur une meilleur convergence empirique dans un premier temps, ainsi pour maîtriser l'influence du censure selon l'échantillon étudié.

Dans la deuxième partie nous avons considéré le bootsrap pour étudier la normalité asymptotique empirique. Nous avons constaté que notre estimateur présente bel et bien une convergence en loi disant empiriquement, ce qui reste comme perspectifs à démontrer théoriquement ainsi que la fixation du nombre des valeurs extrêmes correspondant en minimisant l'erreur quadratique moyenne théorique correspondante.

Nous envisageons à présenter une application réelles pour des données de crédit majores des banques où les echanciers du remboursements des benifisseurs presentent de la censure

Bibliographie

- [1] Alexandre Lekina. Estimation non-paramétrique des quantiles extrêmes conditionnels. Mathématiques [math]. Université Joseph-Fourier - Grenoble I; Université de Grenoble, 2010. Français. tel-00529476v2ffAlvarado, E., Sandberg, D. et Pickford, S. (1998). Modeling large forest fires as extreme events. Northwest Science, 72 :66-75.
- [2] Alvarado, E., Sandberg, D. et Pickford, S. (1998). Modeling large forest fires as extreme events. Northwest Science, 72 :66-75
- [3] Balkema, A., de Haan, L. Residual life at great age. Ann. Probab., 2 :792-804, 1974. (Cit  en pages xxi, 121, 124, 128 et 135.)
- [4] Beirlant, J., Goegebeur, Y., Segers, J., and Teugels, J. (2006). Statistics of Extremes : Theory and Applications. John Wiley.
- [5] Beirlant, J., and Guillou, A. (2001). Pareto index estimation under moderate right censoring. Scand. Actuar. J., 111 125
- [6] Beirlant, J., Goegebeur, Y., Segers, J., Teugels, J. Statistics of extremes, theory and applications. Wiley, New York, 2004. (Cit  en pages xix, 22, 81 et 82.)-
- [7] Beirlant, J., Guillou, A., Dierckx, G., Fils-Villetard, A. . Estimation of the extreme value index and extreme quantiles under random censoring. Springer Science+Business Media, LLC 2007, 10 :151-174, 2007. (Cit  en pages xx, 73, 74, 82, 86, 114 et 123.)-
- [8] Beirlant, J., Vjncikier, P., Teugels, J.L. Tail index estimation, pareto quantile plots, and regression diagnostics. J. Amer Statist. Assoc., 91 :1659-1667, 1996. (Cit  en pages xix, 14, 22, 40 et 123.)-
- [9] Beirlant, J., Goegebeur, Y., Segers, J., Teugels, J., de Waal, D. et Ferro, C. (2004b). Statistics of Extremes : Theory and Applications. John Wiley Sons.
- [10] Bingham, N., Goldie, C. et Teugels, J. (1987). Regular Variation. Cambridge University Press. 22, 23, 24, 68, 111
- [11] Bouleau, N. (1991). Splendeurs et mis res des lois de valeurs extr mes. Revue Risques, 4 :85-92.
- [12] Brahim, B., Meraghni, D., Necir, A. On the asymptotic normality of hill's estimator of the tail index under random censoring. Preprint : arXiv-1302.1666, 2013. (Cit  en pages 44, 82 et 107.)
- [13] Censure, th se de doctorat, universit  Gaston Berger de Saint-Louis.
- [14] Ceresetti, D., Ursu, E., Carreau, J., Anquetin, S., Creutin, J., Gardes, L., Girard, S. et Moli ni , G. (2012). Evaluation of classical spatial-analysis schemes of extreme rainfall. Natural Hazards and Earth System Sciences, 12 :3229-3240.
- [15] Coles, S. et Tawn, J. (1996). A bayesian analysis of extreme rainfall data. Applied Statistics, pages 463-478.
- [16] Daouia, A., Gardes, L., Girard, S., Lekina, A. Kernel estimators of extreme level curves. Test, 20 :311-333, 2011. (Cit  en pages xix, 82, 85, 88, 89 et 108.) Daouia, A., Gardes,

- L., Girard, S. et Lekina, A. (2011). Kernel estimators of extreme level curves. *Test*, 20(2) :311-333.
- [17] Daouia, A., Gardes, L., Girard, S. et Lekina, A. (2011). Kernel estimators of extreme level curves. *Test*, 20(2) :311-333
- [18] Davison, A.C., Smith, R.L. (1990). Models for exceedances over high thresholds, *Journal of the Royal Statistical Society, series B*, 52, 393-442.
- [19] Dekkers, A.L.M., Einmahl, J.H.J., de Haan, L. A moment estimator for the index of an extreme value index. *Annals of statistics*, 17 :1833-1855, 1989. (Cité en pages xix, 13, 14, 18, 40, 45, 82 et 123.)
- [20] Delafosse, E., Guillou, A. Almost sure convergence of a tail index estimator in the presence of censoring. *comptes rendus mathématiques. Académie des Sciences. Paris*, 335 :375-380, 2002. (Cité en page 82.)
- [21] Einmahl, J. et Magnus, J. (2008). Records in athletics through extreme-value theory. *Journal of the American Statistical Association*, 103(484) :1382-1391.
- [22] Einmahl, J. et Sander, S. (2011). Ultimate 100-m world records through extreme-value theory. *Statistica Neerlandica*, 65(1) :32-42
- [23] Einmahl, J. H. J., Fils-Villetard, A., Guillou, A. Statistics of extremes under Random Censoring. *Bernoulli*, 14(1) ; 207-227.
- [24] Einmahl, J. H., Fils-Villetard, A., and Guillou, A. (2008). Statistics of Extremes Under Random Censoring. *Bernoulli*, 14(1) ; 207-227. et *Gestion des Risques. Evénements Rares pour la Gestion des Risques, DESS 203 de l'Université*
- [25] Ferrez, J., Davison, A.C., Rebetez, M. Extreme temperature analysis under forest cover compared to an open field. *Agricultural and Forest Meteorology*, 151 :992-1001, 2011. (Cité en page 81.)
- [26] Gardes, L. et Girard, S. (2010). Conditional extremes from heavy-tailed distributions : an application to the estimation of extreme rainfall return levels. *Extremes*, 13(2) :177-204.
- [27] Gardes, L., Girard, S. Conditional extremes from heavy-tailed distributions : an application to the estimation of extreme rainfall return levels. *Extremes*, 13 :177-204, 2010. (Cité en page 81.)
- [28] Gardes, L., Girard, S. Functional kernel estimators of large conditional quantiles. *Electronic Journal of Statistics*, 6 :1715-1744, 2012. (Cité en page 81.)
- [29] Gardes, L., Girard, S., Lekina, A. Functional nonparametric estimation of conditional extreme quantiles. *Journal of Multivariate Analysis*, 101 :419-433, 2010. (Cité en pages xix, 23, 81 et 134.)
- [30] Gardes, L., Stupfler, G. Estimation of the conditional tail index using a smoothed local hill estimator. *Extremes*, 17 :45-75, 2014. (Cité en page 82.)
- [31] Goegebeur, Y., Guillou, A., Osmann, M. A local moment type estimator for the extreme value index in regression with random covariates. *The Canadian Journal of Statistics*, 42 :487-507, 2014a. (Cité en page 82.)
- [32] Goegebeur, Y., Guillou, A., Schorgen, A. Nonparametric regression estimation of conditional tails : the random covariate case. *Statistics*, 48 :732-755, 2014b. (Cité en pages 82, 84, 85, 88, 91, 111, 117 et 134.)
- [33] Gomes, M.I., Neves, M.M. Estimation of the extreme value index for randomly censored data. *Biometrical Letters*, 48 :1-22, 2011. (Cité en page 82.)

- [34] Gomes, M.I., Oliveira, O. Censoring estimators of a positive tail index. *Statistics. Probability Letter*, 65 :147-159, 2003. (Cité en page 82.)
- [35] Gumbel, E. (1954). *Statistical theory of extreme values and some practical applications : a series of lectures*. Numéro *Applied Mathematics Series*, 33. National Bureau of Standards, Washington.
- [36] Gumbel, E. (1958). *Statistics of extremes*. Columbia University Press, Columbia.-109-
- [37] Hall, P., Tajvidi, N. (2000). Nonparametric analysis of temporal trend when fitting parametric models to extreme-value data, *Statistical Science*, 15, 153-167.
- [38] Hill, B. A simple general approach to inference about the tail of a distribution. *The Annals of Statistics*, 3 :1163-1174, 1975. (Cité en pages xix, 13, 39 et 123.)
- [39] Hill, B. M.(1975). A simple general approach to inference about the tail of a distribution. *Ann. Statist.*, 3(5) ; 1163-1174 in the presence of censoring. *comptes rendus mathématiques*. Académie des Sciences. Paris, 335 :375-380, 2002. (Cité en page 82.)
- [40] Ndao, P., Diop, A., Dupuy, J.-F. Nonparametric estimation of the conditional tail index and extreme quantiles under random censoring. *Computational Statistics Data*, 79 :63-79, 2014. (Cité en page 82.)
- [41] NERC(1975). *Flood studies report. Rapport technique*, London, Natural Environment Research Council. Paris IX Dauphine Marchés Financiers, Marchés des Matières Premières
- [42] Pisarenko, V.F., Sornette, D. Characterization of the frequency of extreme earthquake events by the generalized pareto distribution. *Pure and Applied Geo-physics*, 160 :2343-2364, 2003. (Cité en page 81.)
- [43] Reiss, R.D, and Thomas, M.(2007). *Statistical Analysis of Extreme Values with Applications to Insurance, Finance, Hydrology and Other Fields*. Birkh
- [44] Reiss, R.-D. et Thomas, M. (2001). *Statistical analysis of extreme values : with applications to insurance, finance, hydrology and other fields*. Birkhäuser Verlag.
- [45] Roncalli, T., 2002, *Théorie des Valeurs Extrêmes ou Modélisation des événements Rares pour la Gestion des Risques*, DESS 203 de l'Université Paris IX Dauphine Marchés Financiers, Marchés des Matières Premières et Gestion des Risques
- [46] Rootzén, H. et Tajvidi, N. (2001). Can losses caused by wind storms be predicted from meteorological observations *Scandinavian Actuarial Journal*, 2001(2) :162-175.
- [47] Sandra Plancade. Estimation de la fonction de répartition conditionnelle à partir de données censurées par intervalle, cas 1, par sélection de modèles. 42èmes Journées de Statistique, 2010, Marseille, France. pp.USB-key. inria-00494833ff.
- [48] Stupfler, G. A moment estimator for the conditional extreme-value index. *Electronic Journal of Statistics*, 7 :2298-2343, 2013. (Cité en pages 76, 81 et 129.)
- [49] Smith, R. L. (1989). Extreme value analysis of environmental time series : an application to trend detection in ground-level ozone (with discussion), *Statistical Science*, 4, 367-393.-17-
- [49] Stupfler, G. A moment estimator for the conditional extreme-value index. *Electronic Journal of Statistics*, 7 :2298-2343, 2013. (Cité en pages 76, 81 et 129.)
- [50] Toulemonde, G., 2008, *Estimation et tests en théorie des valeurs extrêmes*, thèse de doctorat de l'université Paris VI
- [51] Worms, J., Worms, R. New estimators of the extreme value index under random right censoring, for heavy-tailed distributions. *Extremes*, 17 :337-358, 2014. (Cité en page 82.)
- [52] Worms, J., Worms, R. New estimators of the extreme value index under random right censoring, for heavy-tailed distributions. *Extremes*, 17 :337-358, 2014. (Cité en page 82.)

Résumé

En 2013 Bearan. J et al on présenté un nouveau estimateur robuste des moments harmoniques de l'indice des valeurs extrêmes à variation régulière avec un paramètre de réglage notamment sous étude asymptotique. Dans notre travail, nous avons adapté ce dernier dans le cas conditionnelle lorsque l'échantillon Y est lié à la covariable X , notamment nous avons calculés notre nouveau estimateur dans le cas des données complètes et censurées. Nous avons présenté quelques simulations simples ainsi que le bootstrap pour l'emploi de la sous-composante éventuelle analyse asymptotique.

Mot clés : la théorie des valeurs extrêmes, les données censurées, variation régulière, indice de valeurs extrêmes, variable conditionnelle, indice de valeurs extrêmes, fenêtre mobile, bootstrap

Abstract

In 2013 Bearan. J et al. presented a new robust estimator of the harmonic moments of the index of the extreme values with regular variation with a tuning parameter, in particular its asymptotic study was established. In our work, we have adapted the latter in the conditional case when the sample Y is linked to the covariate X . Thus we have calculated our new estimator in the case of complete and censored data. We have presented some simple simulations as well as the use of the bootstrap for a possible asymptotic analysis.

Keywords : extreme value theory, censored data, regular variation, extreme value index, conditional variable, extreme value index, moving window, bootstrap

ملخص

في عام 2013 بيران. ج. وآخرون. قدموا مقدرًا جديدًا للعزوم الهندسية لمؤشر القيم القصوى المتغيرة بانتظام ذات وسيط تعديل، ولا سيما تقديم دراسة التقارب اللانهائي المرافق. في عملنا، قمنا بتكييف هذا الأخير في الحالة الشرطية عندما تكون العينة Y مرتبطة بالمتغير المشترك X . وهكذا قمنا بحساب مقدرنا الجديد في حالة البيانات الكاملة والخاضعة للرقابة. لقد قدمنا ببعض عمليات المحاكاة البسيطة بالإضافة إلى استخدام تقنية السحب لتحليل تقارب المقدر اللانهائي.

الكلمات المفتاحية: نظرية القيمة القصوى، البيانات الخاضعة للرقابة، التباين المنتظم، مؤشر القيمة القصوى، المتغير الشرطي، مؤشر القيمة القصوى، النافذة المتحركة، السحب