



## Mémoire d'Actuariat - Promotion 2007

---

# Problématique de seuil dans la modélisation de la sinistralité en Réassurance Non Vie

Emilie Deprey et Alexandre Godzinski

---

**MOTS-CLÉS** : *seuil, sinistralité attritionnelle, sinistralité extrême, théorie des valeurs extrêmes, loi de Pareto généralisée*

**KEYWORDS** : *threshold, attritional losses, extreme losses, Extreme Value Theory, Generalized Pareto Distribution*

**ENCADREMENT** : Jérôme Sarraïl et Célia Vesselinoff (SCOR)

**CORRESPONDANT E.N.S.A.E.** : Christian-Yann Robert

---



## Remerciements

Tout d'abord, nous tenons à remercier Célia Vesselinoff et Jérôme Sarrail de la société SCOR pour nous avoir proposé ce sujet ainsi que pour leur encadrement tout au long de cette étude.

Nous remercions également Christian Robert et Arthur Charpentier pour leur disponibilité et leurs précieux conseils.

# Table des matières

<b>Introduction</b>	<b>5</b>
<b>1 Contexte</b>	<b>6</b>
1.1 Problématique . . . . .	6
1.2 La Réassurance . . . . .	6
1.2.1 Principe de la réassurance . . . . .	6
1.2.2 Les types de réassurance . . . . .	7
1.2.3 Les types de traités . . . . .	7
1.2.4 Les types de sinistres . . . . .	9
1.2.5 La durée de l'engagement . . . . .	9
1.3 Principes de Tarification . . . . .	10
1.3.1 Tarification des traités non proportionnels . . . . .	11
1.3.2 Tarification des traités proportionnels . . . . .	12
1.4 Enjeux du mémoire . . . . .	13
1.4.1 Enjeux pour le Pricing . . . . .	13
1.4.2 Enjeux pour la DFA . . . . .	14
<b>2 Outils mathématiques</b>	<b>16</b>
2.1 Lois usuelles en réassurance . . . . .	16
2.1.1 Loi Lognormale . . . . .	16
2.1.2 Loi Gamma . . . . .	16
2.1.3 Loi de Weibull . . . . .	16
2.1.4 Loi de Pareto . . . . .	17
2.1.5 Loi Pareto Généralisée . . . . .	17
2.2 Tests d'adéquation . . . . .	18
2.2.1 Problématique des données de réassurance et lois tronquées . . . . .	18
2.2.2 Calcul des estimateurs, maximum de vraisemblance . . . . .	18
2.2.3 Le test de Kolmogorov-Smirnov . . . . .	19
2.2.4 Le test d'Anderson Darling . . . . .	20
2.3 Théorie des Valeurs Extrêmes . . . . .	21
2.3.1 Loi du maximum et distribution GEV (Generalized Extreme Value) . . . . .	21
2.3.2 Ajustement de lois GEV . . . . .	22
2.3.3 Loi des excès et distribution GPD (Generalized Pareto Distribution) . . . . .	23
2.3.4 Estimation de l'indice de queue . . . . .	24
<b>3 Approche préliminaire : méthodes graphiques de détermination du seuil</b>	<b>27</b>
3.1 QQ-Plot . . . . .	27
3.2 Mean Excess Function . . . . .	29
3.3 Hill-plot, Pickands-plot . . . . .	31
3.4 Gertensgarbe plot . . . . .	32

<b>4</b>	<b>Modélisation</b>	<b>34</b>
4.1	Cas 1 : Une seule distribution fittée . . . . .	35
4.2	Cas 2 : Deux distributions pour modéliser deux types de sinistralités . . . . .	37
4.3	Application du modèle à deux lois sur données simulées . . . . .	41
4.4	Approches alternatives . . . . .	47
4.4.1	Cas plus élaboré à trois lois . . . . .	47
4.4.2	Modèle de minimisation contraint . . . . .	49
<b>5</b>	<b>Applications</b>	<b>52</b>
5.1	DFA . . . . .	52
5.1.1	Notations . . . . .	53
5.1.2	Revalorisation des sinistres . . . . .	53
5.1.3	Projection à l'ultime de la charge des sinistres . . . . .	55
5.1.4	Etude de la série des ultimes . . . . .	56
5.1.5	Modélisation à l'aide d'une seule distribution . . . . .	58
5.1.6	Modélisation à deux lois . . . . .	59
5.1.7	Estimation de quantiles extrêmes . . . . .	64
5.2	Pricing . . . . .	67
5.2.1	Revalorisation des sinistres . . . . .	67
5.2.2	Obtention des ultimes . . . . .	67
5.2.3	Etude de la série des ultimes . . . . .	67
5.2.4	Modélisation à l'aide d'une seule distribution . . . . .	70
5.2.5	Modélisation à deux lois . . . . .	71
5.2.6	Tarification de traités en excédent de sinistre . . . . .	76
	<b>Conclusion</b>	<b>79</b>
	<b>Références</b>	<b>80</b>
	<b>Annexe 1 : L'outil Excel</b>	<b>82</b>
	<b>Annexe 2 : Rappel sur le modèle de risque collectif</b>	<b>87</b>
	<b>Annexe 3 : Remarque sur l'impact de la méthode du Chain Ladder</b>	<b>88</b>

# Introduction

La sinistralité issue de l'activité de réassurance non vie se caractérise par une nature complexe, ce qui rend difficile sa modélisation. Deux types de sinistralités peuvent généralement être distingués : d'une part la sinistralité attritionnelle, qui correspond à des sinistres de forte fréquence et de faible coût, d'autre part la sinistralité extrême, qui correspond à des sinistres de faible fréquence et de forte sévérité, c'est-à-dire des sinistres survenant rarement mais dont le coût est élevé, qui peuvent être par exemple, selon la branche considérée, des incendies importants en assurance habitation ou des accidents automobiles coûteux en assurance dommages ou responsabilité civile. On souhaite disposer d'un modèle global prenant en compte ces deux types différents de sinistralité. Ce modèle consiste en la proposition d'une distribution pour modéliser la sinistralité au global. La partie centrale de cette distribution modélisera la sinistralité attritionnelle, tandis que la queue de distribution modélisera les sinistres extrêmes. La question sous-jacente est alors la suivante : comment déterminer le seuil départageant un sinistre attritionnel d'un sinistre extrême ? Ce mémoire tentera d'apporter une réponse, en présentant des outils adaptés à la détermination du seuil et en construisant des modèles faisant intervenir ce dernier.

Une première partie sera consacrée à la mise en place du problème et de son contexte : présentation de l'activité de réassurance, des types de traités, des principes de la tarification ainsi que des enjeux du mémoire. Une deuxième partie exposera quant à elle les outils mathématiques nécessaires pour la résolution du problème posé, à savoir principalement la théorie des tests et la théorie des valeurs extrêmes. Dans une troisième partie seront présentées en détail les méthodes graphiques de détermination du seuil, qui seront ensuite implémentées informatiquement et illustrées par une application sur des données simulées. Une quatrième partie permettra de présenter les modèles que nous proposons afin de modéliser dans son ensemble la sinistralité. Enfin, la cinquième et dernière partie mettra en œuvre les outils et modèles proposés sur des données réelles.

# 1 Contexte

## 1.1 Problématique

L'activité de réassurance se caractérise par une sinistralité de nature complexe. Deux types de sinistralités se distinguent : la sinistralité attritionnelle d'une part (les sinistres de forte fréquence et de faible coût) et les sinistres rares mais d'intensité extrême d'autre part. L'utilisation d'une distribution unique pour modéliser ces deux types de sinistralités ne permet pas de rendre compte de ces spécificités. Or on souhaite disposer d'une modélisation globale permettant de prendre en compte ces deux types de sinistralités. Il est important de parvenir à estimer précisément la loi des sinistres fréquents aussi bien que celle des plus rares ; c'est pourquoi il est nécessaire de définir ce qu'est un sinistre de "forte sévérité" en déterminant, pour un portefeuille donné, un seuil de sinistralité critique caractérisant les sinistres de pointe. Ainsi dans cette optique, nous chercherons, à l'aide de différents outils, à déterminer un seuil critique au-delà duquel les sinistres puissent être considérés comme sinistres extrêmes et ce par le biais d'un modèle à établir puis à implémenter.

## 1.2 La Réassurance

### 1.2.1 Principe de la réassurance

*« La réassurance est une opération par laquelle une société d'assurance, ou cédante, s'assure elle-même auprès d'une autre société dénommée réassureur, ou cessionnaire, pour tout ou partie des risques qu'elle a pris en charge ».*

La réassurance ne diffère de l'assurance que par une plus grande complexité due à la diversité plus importante de ses activités et à son caractère international. La réassurance permet à une cédante d'obtenir certains avantages, notamment une réduction de son engagement net sur des risques individuels et une protection contre des pertes multiples ou importantes. La réassurance permet également à une cédante d'obtenir une capacité de souscription supérieure et donc de souscrire des polices portant sur des risques plus importants et plus nombreux, ce qui ne serait pas possible sans une augmentation de ses fonds propres. Le volume du risque, en valeur, que décide de conserver l'assureur pour son propre compte est appelé le « plein de conservation », ou rétention.

La protection apportée par la réassurance a bien évidemment pour contrepartie le paiement d'une prime (la prime de réassurance) par l'assureur à son ou ses réassureur(s). Ces derniers, poursuivant tout comme l'assureur des objectifs de rentabilité et de stabilité des résultats, mutualisent leurs risques à une échelle géographique supérieure à celle de l'assureur direct. Les sociétés de réassurance sont donc généralement par nature des groupes de taille conséquente à dimension internationale. De plus, elles ont, comme l'assureur direct, recours à la réassurance, appelée alors rétrocession, c'est-à-dire *« l'opération par laquelle un réassureur cède à son tour, une partie des risques qu'il a réassurés à un rétrocessionnaire qui peut être une société de réassurance ou une société d'assurance ».*

## 1.2.2 Les types de réassurance

La réassurance comprend trois modes : La réassurance facultative, la réassurance facultative-obligatoire et la réassurance obligatoire. En réassurance facultative, l'assureur et le réassureur sont libres de céder ou d'accepter un risque en totalité ou en partie. Il s'agit de la forme la plus ancienne de réassurance. En réassurance facultative obligatoire, l'assureur a la possibilité de céder ou non, mais le réassureur a l'obligation d'accepter tout ce qui lui est cédé, selon des conditions définies au préalable. Enfin, la réassurance obligatoire établit des obligations réciproques : la société d'assurance s'engage, durant une période donnée, à céder les risques d'une catégorie donnée, la société de réassurance étant obligée de les accepter. Ce mode de réassurance est le plus couramment utilisé. Dans cette étude, nous nous concentrerons sur ce dernier type de réassurance. Le contrat de réassurance obligatoire liant l'assureur et le réassureur s'appelle un traité.

## 1.2.3 Les types de traités

Les traités de réassurance peuvent être soit des traités proportionnels, soit des traités non proportionnels.

### Réassurance proportionnelle

La nature d'un traité de réassurance est dite proportionnelle lorsque la part que le réassureur aura à supporter dans la charge de tout sinistre est égale à la part qu'il a reçue de la prime. Il y a donc égalité entre la proportion de primes reçues par le réassureur et la proportion du coût des sinistres transférée au réassureur.

$$\frac{\text{prime de réassurance}}{\text{primes totales reçues par la cédante}} = \frac{\text{montant des sinistres à la charge du réassureur}}{\text{montant brut des sinistres à la charge de la cédante}}$$

On distingue deux types de traités proportionnels :

- Les traités en quote part :

Le réassureur prend en charge une proportion identique sur tous les risques du portefeuille. Dans ce cas, l'assureur cède la même part sur les risques faibles que sur les risques importants, le profil de portefeuille conservé par le réassureur est semblable au portefeuille initial, seul le niveau des engagements est modifié.

- Les traités en excédent de plein :

Le taux de cession est calculé police par police. Pour chaque police, le réassureur prend en charge uniquement la portion de risque dépassant un niveau de capital appelé plein de rétention.

### Réassurance non proportionnelle

Le traité de réassurance non proportionnelle est défini par une priorité (ou franchise) et un plafond. Le réassureur prend en charge tout ou partie du sinistre qui excède la priorité du traité et dans la limite de la portée (différence entre le plafond et la franchise). En contrepartie, le réassureur perçoit une prime pour compenser le risque qu'il prend. On distingue classiquement trois types de traités non proportionnels : les traités en excédent de sinistre par risque, en excédent de sinistre par événement et les traités en excédent de perte annuelle.

- L'excédent de sinistre par risque :  
Pour chaque risque, on définit une franchise  $f$ , une portée  $p$  et on note " $pXSf$ " l'excédent de sinistre. La cédante obtient ainsi une protection contre les sinistres importants et conserve une partie plus grande des primes. En revanche, le niveau de la prime de réassurance est plus difficile à déterminer et la cédante n'obtient pas de protection contre une accumulation de petits sinistres.
- L'excédent de sinistre par événement :  
Le réassureur couvre un ensemble de sinistres individuels reliés par un même fait générateur et offre ainsi à l'assureur une protection contre l'accumulation de sinistres par événement. L'événement est défini contractuellement dans sa nature, dans l'espace et dans le temps.
- L'excédent de perte annuelle :  
Dans ce cas, l'assureur cherche à se prémunir contre les mauvais résultats puisque ce traité offre à l'assureur une protection contre l'accumulation de sinistres sur une durée. La priorité est donc définie comme la sinistralité annuelle que l'assureur conserve à sa charge.

Par ailleurs, une couverture non proportionnelle est fréquemment composée de plusieurs tranches, chacune étant caractérisée par un couple priorité-portée. Un réassureur assume le risque lorsqu'il dépasse le montant de la rétention jusqu'à concurrence d'une certaine limite. À ce moment-là, un autre réassureur assume la responsabilité jusqu'à un montant supérieur donné, etc... Considérons par exemple le cas où une cédante bénéficie d'une couverture constituée de  $n$  tranches, où les montants  $f_i$  et  $p_i$  décrivent respectivement la franchise et la portée de la tranche  $i$  (pour  $i = 1, \dots, n$ ). Dans la pratique, toute tranche  $i + 1$  sera telle que  $f_{i+1} = f_i + p_i$ , si bien que l'ensemble de la couverture assure à la cédante une protection de  $\sum_{i=1}^n p_i$  en excédent de  $f_1$ , comme l'illustre le graphique ci dessous.

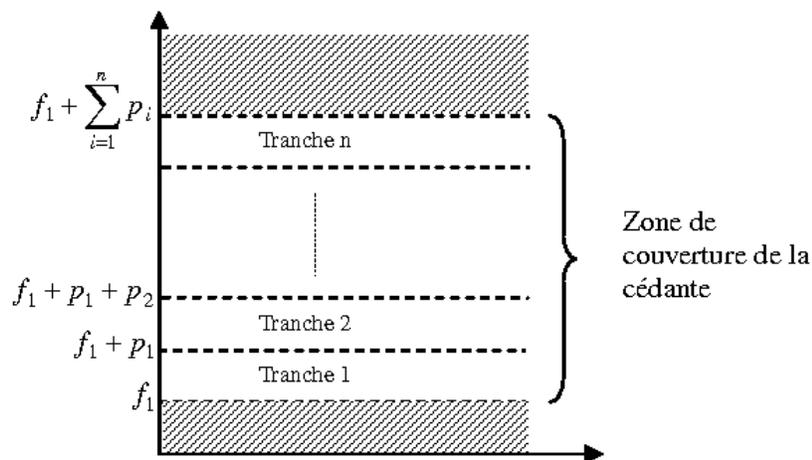


FIG. 1 – Différentes tranches dans un traité non proportionnel.

On dit du réassureur couvrant le risque s'inscrivant immédiatement au-delà de la tranche

de rétention de la cédante qu'il souscrit une tranche basse de réassurance en excédent de sinistre (dans notre exemple, il s'agit du traité  $p_1XSf_1$ ). Un sinistre dont le montant est tout juste supérieur à la rétention de la cédante entraînera des pertes pour le réassureur de la tranche basse, mais pas pour les réassureurs des tranches supérieures. Il est plus facile de prévoir les sinistres dans les tranches basses car ils se produisent plus fréquemment.

#### 1.2.4 Les types de sinistres

En fonction de leur fréquence et de leur sévérité, les sinistres peuvent être classés dans deux grandes catégories :

- La sinistralité attritionnelle d'une part. Il s'agit de la sinistralité classique, récurrente. Elle est composée des sinistres ayant une fréquence importante et une sévérité peu élevée.
- La sinistralité extrême d'autre part. Il s'agit de sinistres rares de forte sévérité.

La sinistralité attritionnelle touche les tranches qui sont alors appelées tranches travaillantes. Elles sont plus faciles à tarifer du fait du plus grand nombre d'observations. A l'inverse, les tranches partiellement travaillantes touchées par la sinistralité extrême et les tranches non travaillantes, pour lesquelles on n'observe aucun sinistre, sont plus compliquées à tarifer. Elles peuvent cependant avoir un impact significatif sur le résultat de l'entreprise, d'où l'intérêt de bien les appréhender.

#### 1.2.5 La durée de l'engagement

Il s'agit du cycle économique complet : souscription, survenance du sinistre, déclaration et enfin règlement. On distingue les branches à développement court (short tail) et les branches à développement long (long tail).

- Les branches courtes :  
La déclaration des sinistres est faite rapidement, l'estimation du montant des dommages se fait dans un délai de quelques mois et le règlement est immédiat. Ce sont des branches qui offrent moins d'imprévus. Les exemples classiques sont l'incendie, le vol, tout ce qui touche à la propriété...
- Les branches longues :  
La déclaration des sinistres peut arriver 10 ou 20 ans après sa survenance. Les estimations des montants des dommages peuvent durer plusieurs années et sont souvent sujettes à de nombreux procès. Ces sinistres sont difficiles à appréhender, les montants des provisions techniques sont plus importants. Un exemple classique est la responsabilité civile.

Ainsi, dans certaines branches comme la Responsabilité Civile (RC), les sinistres se règlent sur plusieurs années. Le coût total de la réparation est alors difficile à évaluer lors de la déclaration du sinistre, puisque le montant dû à la victime lui sera payé le plus souvent en plusieurs fois et sur plusieurs années : Il faut attendre la consolidation de la victime pour pouvoir quantifier les conséquences du sinistre (par exemple, les conséquences des fautes d'un obstétricien peuvent être connues après la puberté de la victime, soit des années après sa naissance). L'indemnisation peut se faire sous forme de rentes viagères ou temporaires.

Enfin, il est fréquent que les instances judiciaires interviennent pour réévaluer le sinistre. Ainsi, pour faciliter l'analyse du montant de sinistre, l'assureur sépare la part déjà payée (Sinistre Payé ou "*SP*", connue), de la part restant à payer (Sinistre à Payer ou "*SAP*", estimée) du montant total du sinistre (sinistre encouru ou "*Incurred*"). On a donc la relation :  $Incurred = SP + SAP$ . Si l'évaluation du montant de sinistre au moment de la déclaration est bien faite, alors à chaque fois que le sinistre sera revu, la part payée augmentera d'autant que la part restant à payer diminuera, laissant l'*incurred* inchangé. Dans la réalité, cela est rarement le cas.

### 1.3 Principes de Tarification

Pour bien tarifier un traité, il est important de pouvoir évaluer au plus juste la prime pure, c'est à dire l'espérance du montant total des sinistres. Pour une année donnée, notons  $N$  le nombre de sinistres (appelé aussi fréquence des sinistres par abus de langage) et  $X_i$  le coût du  $i$ -ème sinistre (également appelé sévérité).  $N$  est une variable aléatoire discrète à valeur dans  $\mathbb{N}$ , alors que  $X_i$  est une variable aléatoire continue à valeur dans  $\mathbb{R}^+$ . On note  $S$  le coût total des sinistres pour une année de cotation :  $S = \sum_{i=1}^N X_i$ .  $S$  représente le montant total des sinistres d'un portefeuille composé de plusieurs polices homogènes où une police peut donner lieu à plusieurs sinistres. L'objectif est de calculer au mieux la prime pure, à savoir  $E(S)$ . Pour cela, différentes méthodes sont envisageables.

- La Tarification sur expérience (méthode du burning cost) :

Il s'agit d'une méthode basée sur la sinistralité observée. C'est une approche par la méthode des moments. Dans un premier temps, on actualise les données (prise en compte de l'inflation, évolution du profil de risque et du coût du risque). On obtient alors une série de données " as if ". Ensuite, on calcule la moyenne pondérée des ratios  $S/P$  revalorisés (le Burning Cost) et on en déduit la prime pure en multipliant le Burning Cost à l'assiette de prime estimée pour l'année considérée. Cette méthode ne donne des résultats acceptables que si l'on se situe dans une tranche totalement travaillante, c'est à dire une tranche traversée totalement par la sinistralité (cf. figure 2). Pour les tranches non travaillantes ou partiellement travaillantes, la tarification s'effectue à l'aide d'une extrapolation (souvent de Pareto).

- Le modèle probabiliste (modèle fréquence / sévérité) :

Cette méthode consiste à modéliser séparément le coût des sinistres ( $X_i$ ) et leur fréquence  $N$ . Sous l'hypothèse que les  $X_i$  sont indépendants et identiquement distribués de même loi qu'une variable aléatoire  $X$ , et que les  $X_i$  et  $N$  sont des variables aléatoires indépendantes, on obtient le résultat fondamental de ce modèle qui est la formule de calcul de prime pure :  $E(S) = E(X) \times E(N)$ <sup>1</sup>.

Il s'agit d'une approche paramétrique. Les lois les plus utilisées pour la fréquence  $N$  sont la distribution de Poisson et la distribution binomiale négative.

Concernant la sévérité, à la SCOR, deux méthodes sont utilisées pour ajuster une loi à la fonction de répartition empirique : le "fitting" et le "blending".

<sup>1</sup>le modèle fréquence / sévérité ou modèle de risque collectif est détaillé plus en détail en annexe 3

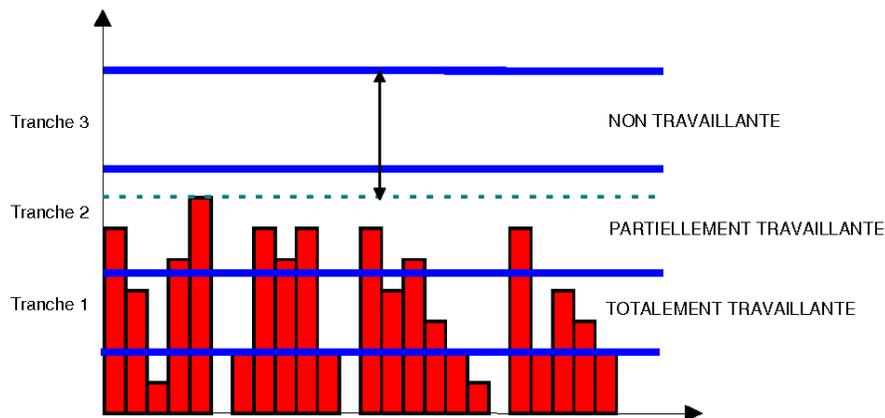


FIG. 2 – Limites de la tarification sur expérience.

Le "fitting" consiste à substituer une loi usuelle à la loi empirique. Les lois usuelles les plus souvent utilisées pour la sévérité sont la loi de Pareto, la loi Gamma, la loi de Weibull et la loi lognormale. La qualité de la substitution est mesurée à l'aide d'un test d'adéquation statistique, qui peut être le test de Kolmogorov-Smirnov ou le test d'Anderson Darling.

Le "blending" est un mélange entre la conservation de la répartition empirique et la substitution par une loi usuelle comme dans le cas du fitting. Pour les sinistres attritionnels, on considère que l'historique est suffisamment représentatif donc on conserve la fonction de répartition empirique. Pour les sinistres extrêmes, on a peu ou pas d'observations, donc on utilise une loi usuelle, par exemple une loi de Pareto. Entre les deux, on raccorde par une fonction de répartition hybride.

– La Tarification par simulation :

C'est une méthode utile lorsque la complexité du modèle envisagé rend difficile toute approche analytique. Elle permet d'obtenir des intervalles de confiance d'une prévision liée au nombre de simulations effectuées.

A noter qu'il existe encore d'autres méthodes de tarification telles que les modèles de régression ou la tarification sur exposition.

### 1.3.1 Tarification des traités non proportionnels

Considérons un traité en excédent de sinistre de franchise  $d$  et de portée infinie<sup>2</sup>. Dans ce cas, la répartition primes / sinistres entre l'assureur et le réassureur est la suivante (on note  $\pi_i$  la prime perçue par l'assureur pour le risque  $i$ ) :

<sup>2</sup>Il est utile de savoir tarifier les traités de portée infinie car on peut ensuite en déduire le tarif des autres traités quelle que soit la portée. On aura :  $\text{prime}(pXSp) = \text{prime}(+\infty XSp) - \text{prime}(+\infty XSp)$

Objet	Assureur	Réassureur
Primes	$\sum_i \pi_i - p$	$p$
Sinistres	$\sum_i \min(X_i, d)$	$\sum_i (X_i - d)_+$

Nous cherchons à déterminer  $p$ , que nous supposons égal à la prime pure, ce qui donne d'après la formule fréquence/sévérité du modèle de tarification probabiliste :

$$p = \sum_{i=1}^N E[(X_i - d)_+] = E(N)E[(X - d)_+]$$

Mais le calcul de  $E[(X - d)_+]$  fait apparaître quelques problèmes.

Par la méthode non paramétrique du burning cost, nous estimons l'espérance par une moyenne :  $E[(X - d)_+] = \frac{1}{n} \sum_{i=1}^n (X_i - d)_+$ . Ceci donne de bons résultats pour les tranches basses, pour lesquelles nous avons beaucoup d'observations, mais pour les tranches élevée, les résultats ne sont pas du tout robustes. D'où l'intérêt d'une méthode paramétrique.

Par la méthode paramétrique, nous sommes amenés à calculer  $E[(X - d)_+] = \int_d^{+\infty} (x - d)dF(x)$  où  $F$  est la fonction de répartition de  $X$ . Il est donc important de bien modéliser la loi de  $X$ . En particulier, lorsque  $d$  est grand, il faut une bonne modélisation de la queue de distribution de  $X$  pour avoir un bon estimateur ou bien une bonne modélisation de  $X$  quand  $X > d$ .

### 1.3.2 Tarification des traités proportionnels

La tarification des traités proportionnels s'effectue en deux étapes. La première étape consiste à tarifier les sinistres attritionnels, tandis que la seconde étape a pour but de tarifier ce que l'on nomme "la tranche de capacité non travaillante". Nous présentons la méthode de tarification utilisée par SCOR dans un cas classique où les montants de sinistres individuels sont connus. Supposons que l'on observe  $n$  sinistres et notons  $max_1$  le plus grand. D'autre part, nous notons  $C$  la capacité du traité.

#### La partie attritionnelle

La partie dite attritionnelle correspond à la tranche pour laquelle nous avons des observations. Elle s'étend donc jusqu'au montant  $max_1$ . Nous utilisons la méthode du burning cost pour estimer le loss ratio moyen espéré pour l'année de cotation. Les étapes sont donc les suivantes :

- Mise à jour et liquidation des primes  $P_i$  : Nous obtenons les  $P_i^{As\ if}$
- Mise à jour et liquidation des sinistres  $S_i$  : Nous obtenons les  $S_i^{As\ if}$
- Calcul des Loss Ratios (LR) annuels As if :  $S_i^{As\ if} / P_i^{As\ if}$
- Estimation de l'espérance du Loss Ratio<sup>3</sup> par la moyenne empirique (ou éventuellement une moyenne pondérée par les assiettes de prime si nous souhaitons accorder plus de poids au passé proche).

#### La tranche de capacité non travaillante

La tranche de capacité non travaillante correspond à la partie non consommée de la capacité du traité, c'est à dire la tranche située entre  $max_1$  et  $C$ . Ainsi, nous nous ramenons à la tarification d'un traité  $(C - max_1)$  XS  $max_1$ , comme l'illustre le graphe suivant.

<sup>3</sup>Généralement, la modélisation du Loss Ratio est étendue à une loi lognormale.

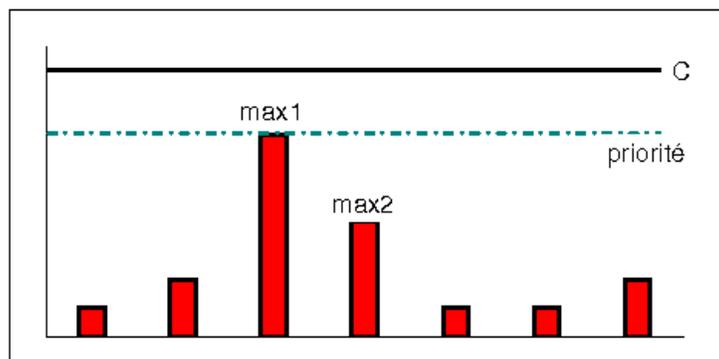


FIG. 3 – Tranche non-travaillante.

Nous voyons dans cette méthode de tarification qu'il est très important de connaître le début de la tranche non travaillante et donc le montant à partir duquel on entre dans la "zone extrême".

## 1.4 Enjeux du mémoire

Dans le cadre d'un besoin à la fois d'amélioration du modèle de tarification des traités de réassurance proportionnelle (approche pricing) et de paramétrage du modèle stochastique de solvabilité (approche Dynamic Financial Analysis (DFA)), ce projet vise à déterminer un seuil de sinistralité critique permettant de caractériser des sinistres dits de pointe et ce par le développement d'un modèle qui sera ensuite implémenté.

### 1.4.1 Enjeux pour le Pricing

L'enjeu du mémoire pour le pricing concerne l'amélioration de la tarification des traités de réassurance proportionnelle.

Nous avons pu voir dans le paragraphe précédent que la tarification des traités proportionnels s'effectuait en deux étapes : dans un premier temps, nous tarifons la partie attritionnelle puis la partie non travaillante. Pour la tranche de capacité non travaillante, le faible nombre (ou même l'absence totale) d'observations ne permet pas une tarification sur expérience. C'est pourquoi à la SCOR, la tranche de capacité non travaillante est tarifée comme s'il s'agissait d'un XS de priorité le plus gros sinistre et de plafond la capacité.

Cependant, cette méthode n'est valable que lorsque nous disposons des montants individuels de sinistres (et donc du montant  $max_1$ ). Or la plupart du temps, la cédante fournit le montant agrégé de ces sinistres au réassureur et l'information individuelle n'est pas disponible. Dans ce cas, comment peut faire le réassureur pour connaître la priorité et la portée de la tranche non travaillante ?

S'il considère que cette portée est égale à la capacité, alors il y a une sur-tarification

puisque la partie attritionnelle a déjà été tarifée à l'expérience, il ne peut donc pas la recompter dans la sinistralité extrême. Tout l'enjeu consiste donc à déterminer une priorité convenable pour cette tranche de capacité non travaillante, correspondant au seuil de sinistralité caractérisant les sinistres de pointe.

In fine, l'objectif est de fournir un seuil de sinistres extrêmes pour diverses branches de la réassurance non-vie en réassurance proportionnelle. Nous nous attacherons également à comparer l'impact de la prise en compte de ce seuil dans la modélisation de la sinistralité sur les primes pures des traités de réassurance non proportionnelle.

## 1.4.2 Enjeux pour la DFA

### Focus sur les principes de la DFA

Le département DFA ("Dynamic Financial Analysis") étudie la solvabilité de la société et en particulier les cash-flows (les primes en entrée, les sinistres en sortie, en simplifiant fortement) qui reflètent les risques que la société encourt. On distingue généralement quatre types de risques :

- Le risque de souscription, lié à la volatilité de la fréquence et de la sévérité des sinistres.
- Le risque de réserve, lié au risque de sous-provisionnement des sinistres.
- Le risque de marché, lié à la volatilité des taux d'intérêt et des actions notamment.
- Le risque de défaut des contreparties de l'assureur (émetteurs obligataires, réassureurs...).

Le département a donc pour objectif de prévoir les flux financiers futurs à l'aide des flux historiques (en supposant que la structure reste la même). Ce qui nous intéresse donc dans cette approche n'est pas le montant total du sinistre mais le coût que le sinistre engendre pour SCOR (c'est à dire après avoir retiré la franchise ou appliqué le plein de conservation suivant le type de tarification).

Concrètement, en notant  $X$  le montant du sinistre, pour un traité en excédent de sinistre  $pXsf$ , on s'intéressera à la variable  $Y$  définie par :

$$Y = \begin{cases} 0 & , X < f \\ X - f & , f < X < f + p \\ p & , f + p < X \end{cases}$$

### Les exigences de solvabilité

Dans le prolongement de la réforme Bâle II pour les banques, l'Union Européenne tente d'établir actuellement un nouveau cadre réglementaire en matière de gestion des risques pour les sociétés d'assurance. La version finalisée de cette réforme baptisée "Solvency II" (ou "Solvabilité II") est attendue courant 2007, pour une mise en application prévue en 2010.

Solvabilité II a pour objectif de fixer des exigences prudentielles qui reflètent au mieux les risques auxquels les sociétés d'assurance sont exposées. La réforme Solvency II s'articule autour de trois piliers :

- Le Pilier I détermine des exigences quantitatives à respecter, notamment sur l'harmonisation des provisions et l'instauration de minima de fonds propres.

- Le Pilier II impose la mise en place de dispositifs de gouvernance des risques (processus, responsabilités, production et suivis d'indicateurs...)
- Le Pilier III, consacré à la discipline de marché, fixe les exigences en termes de reporting et de transparence.

Une des innovations majeures de Solvabilité II provient du Pilier I. Elle consiste en l'introduction de deux niveaux d'exigences de capital :

- Le MCR ("Minimum capital requirement")
- Le SCR ("Solvency capital requirement")

L'adoption de ce mécanisme résulte d'un double constat. D'une part, l'actuelle exigence de marge de solvabilité européenne est trop basse et, sauf cas exceptionnels, une entreprise qui ne la respecte plus n'est pas capable de se redresser par elle-même. D'autre part, on peut observer que les sociétés d'assurance détiennent des fonds propres nettement supérieurs à leur exigence de marge. Par ailleurs, il est apparu nécessaire de mettre en place un système plus adapté au profil de risque de chaque entreprise, qui aboutisse à la fixation d'exigences de capital se rapprochant de la notion de capital économique. Le MCR correspond à un minimum absolu, un "filet de sécurité", en deçà duquel une entreprise ne peut plus opérer : son non-respect entraîne une recapitalisation immédiate sous peine de retrait d'agrément par l'autorité de contrôle. Le SCR correspond quant à lui à un niveau de capital nécessaire pour faire face aux aléas de l'exploitation de l'entreprise, son calcul est basé sur celui d'une Value-at-Risk à 99,5%.

## **Objectifs**

On souhaite anticiper l'ensemble des sorties pour SCOR. On s'attachera donc à modéliser au mieux la loi de sévérité avec une attention particulière pour les comportements des queues de distribution. En effet, ce sont les gros sinistres qui sont les plus dangereux pour la solvabilité de l'entreprise, il est donc important de bien les appréhender. Nous nous attacherons également à fournir des quantiles élevés de la sinistralité afin de mieux connaître les risques extrêmes et de s'adapter à la future directive Solvency II.

## 2 Outils mathématiques

### 2.1 Lois usuelles en réassurance

Dans cette partie, nous présentons rapidement les lois fréquemment utilisées en réassurance pour modéliser la loi du montant des sinistres.

#### 2.1.1 Loi Lognormale

La variable  $X$  suit une loi lognormale de paramètres  $\mu$  et  $\sigma$ , ( $X \sim \ln N(\mu, \sigma)$ ) si le logarithme de  $X$  suit une loi normale ( $\ln(X) \sim N(\mu, \sigma)$ ). Sa densité est donnée par l'expression suivante :

$$f(x) = \frac{1}{x\sigma\sqrt{2\pi}} \exp\left(-\frac{(\ln(x) - \mu)^2}{2\sigma^2}\right)$$

où  $x > 0$ ,  $\mu \in \mathbb{R}$  et  $\sigma > 0$ .

Les moments de la loi lognormale sont donnés par :

$$\begin{aligned} E(X) &= \exp(\mu + \sigma^2/2) \\ Var(X) &= [\exp(2\mu + \sigma^2)][\exp(\sigma^2) - 1] \end{aligned}$$

#### 2.1.2 Loi Gamma

La variable  $X$  suit une loi Gamma de paramètres  $\alpha$  et  $\beta$ , ( $X \sim Ga(\alpha, \beta)$ ) si sa densité est de la forme :

$$f(x) = \frac{\beta^{-\alpha} x^{\alpha-1} \exp(-x/\beta)}{\Gamma(\alpha)} \mathbb{I}_{\mathbb{R}^+}(x)$$

où  $\alpha > 0$ ,  $\beta > 0$  et  $\Gamma(x) = \int_0^{+\infty} t^{x-1} e^{-t} dt$ .

Les moments de la loi Gamma sont donnés par :

$$\begin{aligned} E(X) &= \alpha\beta \\ Var(X) &= \alpha\beta^2 \end{aligned}$$

#### 2.1.3 Loi de Weibull

La variable  $X$  suit une loi de Weibull de paramètres  $\alpha$  et  $\beta$ , ( $X \sim W(\alpha, \beta)$ ) si sa densité est de la forme :

$$f(x) = \alpha\beta^{-\alpha} x^{\alpha-1} e^{-\left(\frac{x}{\beta}\right)^\alpha} \mathbb{I}_{\mathbb{R}^+}(x)$$

et sa fonction de répartition :

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

où  $\alpha > 0$  et  $\beta > 0$ .

Les moments de la loi de Weibull sont donnés par :

$$E(X) = \frac{\beta}{\alpha} \Gamma\left(\frac{1}{\alpha}\right)$$

$$Var(X) = \frac{\beta^2}{\alpha} \left\{ 2\Gamma\left(\frac{2}{\alpha}\right) - \frac{1}{\alpha} \left[ \Gamma\left(\frac{1}{\alpha}\right) \right]^2 \right\}$$

#### 2.1.4 Loi de Pareto

La variable  $X$  suit une loi de Pareto de paramètres  $\theta$  et  $a$ , ( $X \sim Pa(\theta, a)$ ) si sa densité est de la forme :

$$f(x) = \frac{\theta a^\theta}{x^{\theta+1}} \mathbb{I}_{[a, +\infty[}(x)$$

et sa fonction de répartition :

$$F(x) = 1 - \left(\frac{a}{x}\right)^\theta$$

où  $\theta > 0$  et  $a > 0$ .

Les moments de la loi de Pareto sont donnés par :

$$E(X) = \frac{a\theta}{\theta - 1}$$

$$Var(X) = \frac{\theta a^2}{(\theta - 2)(\theta - 1)^2}$$

#### 2.1.5 Loi Pareto Généralisée

La variable  $X$  suit une loi de Pareto Généralisée de paramètres  $\xi$  et  $\sigma > 0$ , ( $X \sim GPD(\xi, \sigma)$ ) si sa densité est de la forme :

$$g_{\xi, \sigma}(x) = \begin{cases} \frac{1}{\sigma} \left(1 + \frac{\xi x}{\sigma}\right)^{-\frac{(\xi+1)}{\xi}} & , \xi \neq 0 \\ \frac{1}{\sigma} \exp(-x/\sigma) & , \xi = 0 \end{cases}$$

et sa fonction de répartition :

$$G_{\xi, \sigma}(x) = \begin{cases} 1 - \left(1 + \frac{\xi x}{\sigma}\right)^{-1/\xi} & , \xi \neq 0 \\ 1 - \exp(-x/\sigma) & , \xi = 0 \end{cases}$$

pour  $x \geq 0$  si  $\xi \geq 0$  et  $0 \leq x \leq -\sigma/\xi$  si  $\xi < 0$ .

L'espérance de la loi de Pareto généralisée est finie si et seulement si  $\xi < 1$ . De manière générale, si  $\xi < 1/r$ , avec  $r \in \mathbb{N}$ , alors on a

$$E(X^r) = \frac{\sigma^r}{\xi^{r+1}} \frac{\Gamma(\xi^{-1} - r)}{\Gamma(1 + \xi^{-1})} r!$$

On remarque que dans le cas  $\xi = 0$ , on retrouve une loi exponentielle. De plus, la loi GPD peut être généralisée à trois paramètres  $\xi, \sigma$  et  $\mu$ , en remplaçant  $x$  par  $(x - \mu)$ .

## 2.2 Tests d'adéquation

### 2.2.1 Problématique des données de réassurance et lois tronquées

La plupart du temps, la cédante fournit une base de données qui ne contient pas la totalité de ses sinistres au réassureur. Elle fournit l'information sur les sinistres dont le montant est supérieur à un seuil de communication. Ce seuil est certes largement inférieur à la priorité mais l'échantillon dont dispose le réassureur afin de faire ses tests d'adéquation est alors naturellement tronqué (à gauche) et l'erreur commise sur la loi ajustée à partir de cet échantillon se retrouve par la suite dans tout le procédé de tarification.

Nous tenterons dans la suite du mémoire de modéliser la sinistralité à l'aide de deux distributions, une pour la sinistralité attritionnelle, l'autre pour la sinistralité extrême. La première sera tronquée à droite au moment du raccord entre les deux distributions. Il faudra en tenir compte lors des tests d'adéquation.

C'est pourquoi, dans toute cette section nous allons généraliser les tests d'adéquation au cas des lois tronquées à gauche et à droite, et ce, en faisant intervenir un terme de normalisation.

Soit  $X$  une variable aléatoire réelle de fonction de répartition  $F_X$  dont on cherche la loi de distribution. Soit  $Y$  une variable aléatoire de fonction de répartition  $F_Y$  qui suit la loi de  $X$  tronquée à gauche en une valeur connue  $d$  et à droite en une valeur connue  $s$ . Nous disposons d'un échantillon  $Y_1, \dots, Y_n$  et des valeurs  $d$  et  $s$ . Par définition de la variable  $Y$ , nous pouvons écrire :

$$F_Y(x) = P(X \leq x/s > X > d) = \frac{P(d < X \leq x)}{P(s > X > d)}$$
$$F_Y(x) = \begin{cases} 1 & \text{si } x > s \\ \frac{F_X(x) - F_X(d)}{F_X(s) - F_X(d)} & \text{si } d < x \leq s \\ 0 & \text{si } x \leq d \end{cases}$$

On obtient alors  $f_Y$ , la densité de  $Y$  suivante :

$$f_Y(x) = \begin{cases} \frac{f_X(x)}{F_X(s) - F_X(d)} & \text{si } d < x \leq s \\ 0 & \text{sinon} \end{cases}$$

Il apparaît donc un coefficient  $F_X(s) - F_X(d)$  au dénominateur qui correspond à une normalisation de la loi tronquée.

### 2.2.2 Calcul des estimateurs, maximum de vraisemblance

Pour effectuer les tests d'adéquation, nous avons besoin d'un échantillon ainsi que d'une loi de référence entièrement définie. Nous testerons ensuite si l'échantillon correspond bien à cette loi. Mais pour cela, nous avons besoin d'estimer les paramètres que nous ignorons a priori. Plusieurs méthodes sont possibles comme la méthode des moments ou celle du maximum de vraisemblance; nous utiliserons celle du maximum de vraisemblance qui donne les

meilleurs résultats.

Cette méthode consiste à maximiser la fonction de vraisemblance qui s'écrit pour un échantillon  $X_1, \dots, X_n$  de densité  $f_X$  :

$$L = \prod_{i=1}^n f_X(x_i)$$

En pratique, on utilise la fonction logvraisemblance  $l$  qui est le logarithme népérien de  $L$  et qui est plus facile à mettre en oeuvre car elle fait intervenir des sommes au lieu de produits. Comme nous ne disposons pas des  $X_1, \dots, X_n$ , mais des  $Y_1, \dots, Y_n$  tronqués, le coefficient correcteur intervient dans la fonction du maximum de vraisemblance :

$$L = \prod_{i=1}^n f_Y(x_i) = \prod_{i=1}^n \frac{f_X(x_i)}{F_X(s) - F_X(d)}$$

et la log vraisemblance s'écrit :

$$l = \sum_{i=1}^n \ln[f_X(x_i)] - n \ln[F_X(s) - F_X(d)]$$

Pour la plupart des lois usuelles, l'estimateur du maximum de vraisemblance est défini de façon unique et se calcule explicitement (quand une détermination explicite est impossible, il faut avoir recours à une détermination numérique, par un algorithme d'optimisation). Sous des hypothèses vérifiées par de nombreux modèles courants, on démontre qu'il est asymptotiquement sans biais et convergent. On démontre de plus que sa variance est minimale. La méthode du maximum de vraisemblance est donc théoriquement la meilleure des méthodes d'estimation.

Une fois les paramètres estimés, on peut effectuer des tests afin de savoir quel ajustement choisir. Pour tester l'adéquation à des lois continues, on a souvent recours aux tests de Kolmogorov-Smirnov et d'Anderson-Darling. Ces tests sont construits à partir de la fonction de répartition empirique. Le test d'Anderson-Darling est particulièrement intéressant car il accorde plus de poids aux queues de distribution que Kolmogorov-Smirnov.

### 2.2.3 Le test de Kolmogorov-Smirnov

Soit  $X_1, \dots, X_n$ , un échantillon de loi  $F_X$  et soit  $F_0$  une loi continue. On souhaite tester :

$$\begin{cases} H_0 : F_X = F_0 \\ H_1 : F_X \neq F_0 \end{cases}$$

Le test de Kolmogorov-Smirnov consiste à mesurer, pour une variable aléatoire continue, la plus grande distance entre la distribution théorique  $F_0$  et la distribution empirique  $F_n$ . La statistique de Kolmogorov-Smirnov est définie par :

$$D_n = \sup_x |F_n(x) - F_0(x)|$$

qui peut aussi s'écrire :

$$D_n = \max\{D_n^+, D_n^-\}$$

avec  $D_n^+ = \max(\frac{i}{n} - F_0(X_i), 0)$  et  $D_n^- = \max(F_0(X_i) - \frac{i-1}{n}, 0)$  où  $X_i$  est la  $i^{\text{ème}}$  observation de l'échantillon et  $n$  le nombre total d'observations.

Le test au seuil  $\alpha$  associé à cette statistique est défini par la région critique de la forme :

$$\{D_n \geq c_\alpha\}$$

où  $c_\alpha$  est le quantile  $(1 - \alpha)$  de la table de Kolmogorov-Smirnov.

Si on travaille avec les lois tronquées, il faut modifier la fonction de répartition théorique :  $F'_0(x) = \frac{F_0(x) - F_0(d)}{F_0(s) - F_0(d)}$  et le test devient alors

$$\begin{cases} H_0 : F_X = F'_0 \\ H_1 : F_X \neq F'_0 \end{cases}$$

#### 2.2.4 Le test d'Anderson Darling

Ce test est également basé sur un calcul de distance entre la fonction de répartition empirique et la fonction de répartition théorique. La statistique d'Anderson-Darling est définie par :

$$A_n^2 = n \int_{-\infty}^{+\infty} \frac{(F_n(x) - F_0(x))^2}{F_0(x)(1 - F_0(x))} dF_0(x)$$

Cette statistique peut également s'écrire :

$$A_n^2 = -n - \frac{1}{n} \sum_{i=1}^n [(2i - 1) \ln(Z_i) + (2n + 1 - 2i) \ln(1 - Z_i)]$$

où  $Z_i = F_0(X_i)$ .

Le test au seuil  $\alpha$  associé à cette statistique est défini par la région critique de la forme :  $\{A_n^2 \geq c_\alpha\}$  où  $c_\alpha$  est le quantile  $(1 - \alpha)$  de la table d'Anderson-Darling.

Si on travaille avec des lois tronquées, il faut modifier la fonction de répartition théorique pour en tenir compte de la même manière que pour le test de Kolmogorov-Smirnov :  $F'_0(x) = \frac{F_0(x) - F_0(d)}{F_0(s) - F_0(d)}$ .

Le test d'Anderson-Darling donne une importance plus grande aux queues de distribution, nous le préférons donc au test de Kolmogorov-Smirnov en cas de conclusions divergentes entre les deux tests.

## 2.3 Théorie des Valeurs Extrêmes

La théorie des valeurs extrêmes permet d'étudier le comportement des queues de distribution. Autrement dit, elle s'intéresse aux événements survenant rarement. Elle est particulièrement utile dans le cadre de l'assurance ou de la réassurance, car elle permet de mieux cerner les sinistres extrêmes, en vue de la tarification de tranches peu ou non travaillantes (approche pricing) ou en vue de constitution de fonds propres suffisants afin d'assurer la solvabilité de l'entreprise (approche DFA).

La théorie des valeurs extrêmes peut se découper en deux parties, cependant liées entre elles. D'une part le théorème de Fischer-Tippett permet de connaître la distribution asymptotique du maximum de  $n$  variables aléatoires. Ce théorème est l'analogue du théorème central limite, qui s'intéresse à la loi asymptotique de la somme : alors que le théorème central limite fait apparaître une loi normale, le théorème de Fischer-Tippett fait apparaître la distribution GEV (Generalized Extreme Value). D'autre part l'étude de "l'excès fonction" (loi des excès) fait apparaître la distribution de Pareto généralisée.

### 2.3.1 Loi du maximum et distribution GEV (Generalized Extreme Value)

Considérons  $n$  variables aléatoires  $X_1, \dots, X_n$  (des coûts de sinistre par exemple) indépendantes et de même loi et notons  $S_n = X_1 + \dots + X_n$  leur somme. Si l'espérance et la variance de ces variables aléatoires sont finies, le théorème central limite permet d'affirmer que :

$$\frac{S_n - n\mathbb{E}(X)}{\sqrt{n}\sigma} \xrightarrow{\text{loi}} N(0, 1)$$

Lorsque l'on relâche l'hypothèse sur le moment d'ordre 2 de  $X$ , le théorème central limite généralisé permet de connaître le comportement de  $\frac{S_n - a_n}{b_n}$  pour des suites  $a_n$  et  $b_n$  bien choisies.

De manière symétrique, le résultat principal de la théorie des valeurs extrêmes va fournir la loi limite du maximum  $M_n = \max(X_1, \dots, X_n)$  correctement normalisé par des suites  $a_n$  et  $b_n$ . Ce résultat est le suivant :

**Théorème de Fischer-Tippett :** Supposons qu'il existe des constantes de normalisation  $a_n \in \mathbb{R}$  et  $b_n > 0$ , et une loi non dégénérée de fonction de répartition  $H$  telles que

$$b_n^{-1}\{M_n - a_n\} \xrightarrow{\text{loi}} H$$

Alors  $H$  est du même type qu'une des trois lois suivantes (données par leur fonction de répartition) :

1. Loi de Fréchet,  $\Phi_\alpha(x) = \exp(-x^{-\alpha})\mathbb{I}_{(x>0)}$ ,  $\alpha > 0$  ;
2. Loi de Weibull,  $\Psi_\alpha(x) = \exp(-x^{-\alpha})$  si  $x \leq 0$ , et 1 sinon,  $\alpha > 0$  ;
3. Loi de Gumbel,  $\Lambda(x) = \exp(-\exp(-x))$ .

Ces trois lois sont en fait les trois cas particuliers de la distribution GEV (Generalized Extreme Value) qui s'écrit :

$$H_{\xi,\mu,\sigma}(x) = \begin{cases} \exp(-[1 - \xi(x - \mu)/\sigma]^{1/\xi}) & , \xi \neq 0 \\ \exp(-\exp[-(x - \mu)/\sigma]) & , \xi = 0 \end{cases}$$

si  $\mu + \xi x/\sigma > 0$ . En effet,  $\xi = \alpha^{-1} > 0$  correspond à la distribution de Fréchet,  $\xi = 0$  correspond à la distribution de Gumbel, et  $\xi = -\alpha^{-1} < 0$  à la distribution de Weibull.

**Remarque 1 :** Max-domain d'attraction.

S'il existe des constantes de normalisation  $a_n \in \mathbb{R}$  et  $b_n > 0$ , et une loi non dégénérée  $GEV(\xi)$  telles que

$$b_n^{-1}\{M_n - a_n\} \xrightarrow{loi} GEV(\xi)$$

on dit que  $F_X$  appartient au max-domain d'attraction de  $GEV(\xi)$  (on écrira par la suite :  $F_X \in MDA(H_\xi)$ ).

**Remarque 2 :** Paramètre de queue.

Le paramètre  $\xi$  de la loi GEV est appelé le paramètre de queue puisqu'il permet de définir l'épaisseur de la queue de distribution de  $F_X$  :

- Si  $F_X \in MDA(H_\xi)$ ,  $\xi > 0$  alors  $F_X \in MDA(\text{Fréchet})$  et la loi de  $X$  est une loi à queue épaisse. Un exemple est la loi de Pareto.
- Si  $F_X \in MDA(H_\xi)$ ,  $\xi = 0$  alors  $F_X \in MDA(\text{Gumbel})$  et la loi de  $X$  est une loi à queue fine ou moyenne. La loi exponentielle, la loi normale et la loi lognormale en sont des exemples.
- Si  $F_X \in MDA(H_\xi)$ ,  $\xi < 0$  alors  $F_X \in MDA(\text{Weibull})$  et la loi de  $X$  est bornée à droite. Nous pouvons citer comme exemple la loi uniforme ou la loi beta.

### 2.3.2 Ajustement de lois GEV

Les lois les plus usuelles en réassurance correspondent au cas  $\xi \geq 0$ . Il peut être intéressant dans un premier temps de savoir si nous sommes dans le cas Gumbel ( $\xi = 0$ ) ou dans le cas Fréchet ( $\xi > 0$ ). Pour cela, une technique consiste à tracer les deux graphes suivants :

#### Graphes pour la loi de Gumbel

La fonction de répartition de la loi de Gumbel s'écrit :  $\Lambda(x) = \exp(-\exp(-(x - \mu)/\sigma))$ . Autrement dit,  $\ln(\ln(1/\Lambda(x)))$  doit être une fonction linéaire en  $x$ . On regarde donc si les points sont alignés puis on ajuste le modèle

$$\ln\left(\ln\left(\frac{1}{\Lambda(x)}\right)\right) = \alpha x + \beta + \varepsilon$$

où on suppose que  $\varepsilon \sim N(0, s^2)$  est un bruit blanc. On teste l'adéquation de la loi de Gumbel via l'étude de la somme des carrés des résidus. Si l'hypothèse est validée, alors  $\sigma = -1/\alpha$  et  $\mu = -\beta/\alpha$ .

### Graphe pour la loi de Fréchet

La fonction de répartition de la loi de Fréchet s'écrit :  $\Phi_\alpha(x) = \exp(-x^{-\alpha})$ , c'est-à-dire que  $\ln(\ln(1/\Phi_\alpha(x)))$  doit être une fonction linéaire en  $\log(x)$ . On regarde donc si les points sont alignés puis on ajuste le modèle

$$\ln\left(\ln\left(\frac{1}{\Phi_\alpha(x)}\right)\right) = \alpha \ln(x) + \beta + \varepsilon$$

où on suppose que  $\varepsilon \sim N(0, s^2)$  est un bruit blanc.

### 2.3.3 Loi des excès et distribution GPD (Generalized Pareto Distribution)

Dans la partie précédente, nous avons étudié les événements rares à l'aide du maximum de la distribution de  $X$ . Néanmoins, une autre méthode est souvent utilisée en réassurance, il s'agit de la méthode des excès (encore appelée POT, *Peaks Over Threshold*) introduite par de Haan et Rootzen. Ainsi, dans cette section, nous allons nous intéresser au comportement de  $X/X > u$  pour des seuils  $u$  suffisamment grand ( $u \rightarrow +\infty$ ).

La fonction de répartition des excès au delà du seuil  $u$  est définie par :

$$F_u(x) = P(X - u \leq x/X > u) = \frac{F(u+x) - F(u)}{1 - F(u)}$$

$x \geq 0$ ,  $u < x_F$ , où  $x_F$  est le point terminal droit de la distribution de  $X$ ,  $x_F = \sup\{x \in \mathbb{R} : F(x) < 1\}$

La loi asymptotique des excès est donnée par le théorème suivant, démontré par Pickands, Balkema et De Hann.

**Théorème (Pickands, Balkema, De Hann) :** *Si  $F$  appartient à l'un des trois domaines d'attraction de la loi des valeurs extrêmes (DA(Fréchet), DA(Gumbel) ou DA(Weibull)), alors il existe une fonction  $\sigma(u)$  positive, définie à une équivalence près quand  $u \rightarrow x_F$ , et un réel  $\xi$  tels que*

$$\lim_{u \rightarrow x_F} \sup_{0 < x < x_F - u} |F_u(x) - G_{\xi, \sigma(u)}(x)| = 0$$

où  $G_{\xi, \sigma}(x)$  est la fonction de répartition de la loi de Pareto généralisée (ou loi GPD, Generalized Pareto Distribution) définie pour  $\sigma > 0$  par

$$G_{\xi, \sigma}(x) = \begin{cases} 1 - \left(1 + \frac{\xi x}{\sigma}\right)^{-1/\xi} & , \xi \neq 0 \\ 1 - \exp(-x/\sigma) & , \xi = 0 \end{cases}$$

pour  $x \geq 0$  si  $\xi \geq 0$  et  $0 \leq x \leq -\sigma/\xi$  si  $\xi < 0$ .

En effet, supposons que l'approximation de la distribution du maximum par une distribution GEV soit satisfaisante :

$$Pr(M_n \leq x) \approx H(x)$$

Alors pour établir le résultat précédent, nous remarquons que

$$F^n(x) \approx \exp\left\{-\left[1 + \xi\left(\frac{x - \mu}{\sigma}\right)\right]^{-1/\xi}\right\}$$

Nous en déduisons que

$$n \ln F(x) \approx -\left[1 + \xi\left(\frac{x - \mu}{\sigma}\right)\right]^{-1/\xi}$$

Or  $\ln F(x) \approx -(1 - F(x))$  pour  $x$  grand. Nous avons donc

$$F(x) \approx 1 - \frac{1}{n} \left[1 + \xi\left(\frac{x - \mu}{\sigma}\right)\right]^{-1/\xi}$$

Nous obtenons alors le résultat final

$$F_u(x) = \frac{F(u + x) - F(u)}{1 - F(u)} = 1 - \left[1 + \xi\left(\frac{x}{\sigma}\right)\right]^{-1/\xi}$$

**Remarque 1 :** Choix du seuil.

Ce théorème stipule donc que pour un seuil  $u$  suffisamment élevé, la fonction de répartition de la loi des excès peut être approchée par  $G_{\xi, \sigma}(x)$  pour des bonnes valeurs de  $\xi$ , le paramètre de queue et de  $\sigma$ , le paramètre d'échelle. Ainsi, pour un seuil suffisamment grand, on s'attend à ce que les données au dessus du seuil aient le comportement d'une pareto généralisée. La principale difficulté est alors de choisir un seuil approprié. Nous sommes face à un arbitrage biais / robustesse. En effet, plus le seuil est élevé, meilleure est la convergence de la loi des excès vers la GPD mais moins bonne est la précision car nous disposons alors de moins d'observations et inversement pour un seuil plus faible. Nous retrouvons ici toute la problématique du sujet de notre mémoire : choisir le meilleur seuil  $u$  possible.

**Remarque 2 :** Lien entre les approches par maximum et par excès.

On peut montrer que même si les lois limites sont dans des familles différentes (GEV et GPD), leur comportement de queue est similaire et ne fait intervenir qu'un unique paramètre  $\xi$  qui est le même que l'on considère l'approche loi du maximum (GEV) ou l'approche loi des excès (GPD). En particulier, ces deux approches permettent de distinguer deux cas,  $\xi = 0$  et  $\xi > 0$ , correspondant respectivement à des queues de type exponentiel ou des queues dites épaisses, de type pareto.

### 2.3.4 Estimation de l'indice de queue

L'indice de queue  $\xi$  nous informe sur l'épaisseur de la queue de distribution de la loi de  $X$ . Il donne donc une indication sur l'importance des risques extrêmes pour une distribution. C'est pourquoi il est important de bien l'estimer. Avec la remarque 2 ci-dessus, on voit que deux approches vont être possibles pour l'estimation de  $\xi$  : une approche par la loi des valeurs extrêmes généralisée (approche GEV) et une approche par la loi de pareto généralisée (approche GPD).

### Approche GEV

En pratique, le maximum des  $n$  variables aléatoires  $X_1, \dots, X_n$  ne constitue qu'une unique observation, ce qui rend l'estimation de la loi délicate. Nous souhaitons disposer d'un échantillon de maxima  $Y_1, \dots, Y_m$ , afin de pouvoir appliquer des méthodes classiques d'ajustement de lois. L'idée proposée par Gumbel (1958) est la suivante : Considérons un échantillon de taille  $n$ , que l'on découpe en  $k$  blocs de taille  $m$  :

$$\underbrace{X_1, \dots, X_m}_{\text{Bloc 1}}, \underbrace{X_{m+1}, \dots, X_{2m}}_{\text{Bloc 2}}, \dots, \underbrace{X_{(k-1)m}, \dots, X_{km}}_{\text{Bloc } k}, \text{ où } mk = n$$

On note alors  $Y_i$ , le maximum obtenu sur le  $i$ ème bloc (i.e.  $Y_i = \max(X_{(i-1)m+1}, \dots, X_{im})$ ). Les  $X_i$  étant indépendants, identiquement distribués et les blocs étant de même taille, on en déduit que les  $Y_i$  sont également iid et  $P(Y_i \leq y) = f_X(y)^m$ . Si  $m$  est suffisamment grand, on a l'approximation  $Y_i \sim GEV(\xi, \mu, \sigma)$  et on peut effectuer la méthode du maximum de vraisemblance sur les  $Y_i$  pour obtenir des estimateurs des paramètres  $(\hat{\xi}, \hat{\mu}, \hat{\sigma})$ .

En particulier, si  $\xi \neq 0$ , la log-vraisemblance de l'échantillon  $Y_1, \dots, Y_m$ , distribué suivant une loi  $GEV(\xi, \mu, \sigma)$  s'écrit :

$$\log \mathcal{L} = -m \log \sigma - (1 + \xi^{-1}) \sum_{i=1}^m \log \left( 1 + \xi \frac{Y_i - \mu}{\sigma} \right) - \sum_{i=1}^m \left( 1 + \xi \frac{Y_i - \mu}{\sigma} \right)^{-1/\xi},$$

et si  $\xi = 0$ ,

$$\log \mathcal{L} = -m \log \sigma - \sum_{i=1}^m \exp \left( 1 + \xi \frac{Y_i - \mu}{\sigma} \right) - \sum_{i=1}^m \left( 1 + \xi \frac{Y_i - \mu}{\sigma} \right).$$

De plus, si  $\xi > -1/2$ , on a alors :  $\sqrt{m}[(\hat{\xi}, \hat{\mu}, \hat{\sigma}) - (\xi, \mu, \sigma)] \xrightarrow{\text{loi}} \mathcal{N}(0, V)$ , où la forme de  $V$  a été donnée par Smith en 1985.

L'estimation du paramètre  $\xi$  en utilisant la distribution GEV est une estimation dite "block componentwise" : à partir d'un échantillon, on construit un échantillon de maxima en formant des blocs de même dimension. Cela implique donc une perte de certaines informations. En particulier, certains blocs peuvent contenir plusieurs extrêmes, alors que d'autres blocs peuvent ne pas en contenir. Nous ne privilégierons donc pas cette approche.

### Approche GPD

Par cette approche, il est également possible d'utiliser des méthodes de maximum de vraisemblance pour estimer le paramètre  $\xi$  de la loi  $GPD$  des excès. En effet, la log-vraisemblance calculée à partir de  $k$  excès  $Y_1, \dots, Y_k$  s'écrit :

$$\log \mathcal{L} = -k \log \sigma - \left( 1 + \frac{1}{\xi} \right) \sum_{i=1}^k \log \left( 1 + \xi \frac{Y_i}{\sigma} \right), \quad \text{si pour tout } i, 1 + \xi \frac{Y_i}{\sigma} > 0$$

Cependant, la plupart des estimateurs du paramètre de queue par l'approche  $GPD$  reposent sur l'utilisation de la statistique d'ordre  $X_{i:n}$ , la  $i$ ème valeur d'un échantillon de taille

$n$ . En particulier, à partir d'un échantillon de taille  $n$ , on s'intéresse aux  $m$  plus grandes valeurs. Les estimateurs les plus utilisés sont celui de Pickands, introduit en 1975, celui de Hill, également introduit en 1975, valable pour  $\xi > 0$ , et celui de Dekkers, Einmalh et De Haan introduit en 1990 qui sont respectivement :

$$\xi_{n,m}^{Pickands} = \frac{1}{\log 2} \log \frac{X_{n-m:n} - X_{n-2m:n}}{X_{n-2m:n} - X_{n-4m:n}}$$

$$\xi_{n,m}^{Hill} = \frac{1}{m} \sum_{i=1}^m \log X_{n-i+1:n} - \log X_{n-m:n}$$

$$\xi_{n,m}^{DEdH} = \xi_{n,m}^{H(1)} + 1 - \frac{1}{2} \left[ 1 - \frac{(\xi_{n,m}^{H(1)})^2}{\xi_{n,m}^{H(2)}} \right]^{-1}$$

où  $\xi_{n,m}^{H(r)} = \frac{1}{m} \sum_{i=1}^{m-1} (\log X_{n-i:n} - X_{n-m:n})^r$ ,  $r = 1, 2, \dots$

Il est possible de montrer que ces estimateurs convergent quand  $n \rightarrow +\infty$ ,  $m \rightarrow +\infty$  et  $m/n \rightarrow 0$ . On a en effet la normalité asymptotique suivante :

$$\sqrt{m}(\xi_{n,m}^{Pickands} - \xi) \xrightarrow{loi} \mathcal{N}\left(0, \frac{\xi^2(2^{\xi+1} + 1)}{(2(2^\xi - 1) \log 2)^2}\right)$$

$$\sqrt{m}(\xi_{n,m}^{Hill} - \xi) \xrightarrow{loi} \mathcal{N}(0, \xi^2) \quad \text{pour } \xi > 0$$

$$\sqrt{m}(\xi_{n,m}^{DEdH} - \xi) \xrightarrow{loi} \mathcal{N}(0, 1 + \xi^2) \quad \text{pour } \xi \geq 0$$

L'estimateur de Hill est généralement le plus utilisé. Il présente une variance asymptotique plus faible que les deux autres. En revanche il n'est valable que si  $\xi > 0$ .

En outre, la principale difficulté dans le calcul de ces trois estimateurs est le choix de  $m$ , le nombre d'excès à considérer. Nous retrouvons l'arbitrage biais robustesse évoqué précédemment : Si  $m$  est trop important<sup>4</sup>, nous n'avons plus la convergence car nous sortons de la queue de distribution et si  $m$  est trop faible, le peu d'observations rend l'estimateur instable.

<sup>4</sup>Si  $m$  est trop important, nous considérons trop d'excès et donc un seuil  $u$  trop faible pour la modélisation de  $X/X > u$ .

### 3 Approche préliminaire : méthodes graphiques de détermination du seuil

La partie précédente, notamment la théorie des valeurs extrêmes, nous permet de proposer plusieurs méthodes graphiques dédiées à l'étude des sinistres extrêmes, des queues de distribution et à la détermination du seuil de sinistralité extrême que nous allons présenter dans cette partie. Mais auparavant, nous présentons le graphique Quantile-Quantile, un autre outil statistique utile dans l'approche exploratoire. Tous ces outils graphiques ont été implémentés par nos soins sous VBA.

#### 3.1 QQ-Plot

Les graphiques Quantiles-Quantiles (Q-Q plot) permettent de tester graphiquement l'adéquation d'une famille de lois à des données. L'idée consiste à regarder si les quantiles de la famille de loi testée ( $Q(p)$ ) et les quantiles de l'échantillon des  $x_i, i = 1, \dots, n, (Q_n(p))$  sont linéairement liés, pour des valeurs de  $p \in [0, 1]$ . Autrement dit, on souhaite vérifier si les points  $(F^{-1}(p), \hat{F}_n^{-1}(p))$  sont alignés pour différentes valeurs de  $p \in [0, 1]$ , où  $F$  est la fonction de répartition de la vraie loi et  $\hat{F}_n$  la fonction de répartition empirique. En pratique, on représente graphiquement les points  $(F^{-1}(\frac{i}{n+1}), x_{(i)})$  pour  $i = 1 \dots n$  où  $x_{(i)}$  est la  $i^{\text{ème}}$  valeur de l'échantillon ordonné.

Le QQ Plot le plus utilisé est vraisemblablement celui de la loi exponentielle, où  $F(x) = 1 - \exp(-x/\lambda)$  pour  $x \geq 0$ . La fonction quantile est donnée par  $F^{-1}(p) = -\lambda \log(1 - p)$ ,  $p \in ]0, 1[$ . Le QQ Plot exponentiel consiste alors à tracer les points

$$\left(-\ln\left(1 - \frac{i}{n+1}\right), x_{(i)}\right) \quad i = 1, \dots, n$$

Son interprétation est simple : Si l'échantillon des  $x_i$  est un échantillon indépendant et identiquement distribué, issu d'une loi exponentielle, alors les points sont alignés selon une droite dont la pente est donnée par  $1/\lambda$ . Si le graphe a une forme concave, alors la distribution des  $x_i$  est à queue plus épaisse. A l'inverse, si le graphe a une forme convexe, la distribution des  $x_i$  est à queue plus fine.

Le graphe suivant (figure 4) permet d'illustrer ces propos. Nous avons tracé le QQ Plot exponentiel pour différents échantillons simulés. Le premier échantillon suit une loi Exponentielle de paramètre 5, et nous obtenons bien une droite dont la pente semble proche de  $1/5$ . Le graphe obtenu avec les deux échantillons suivants a une forme légèrement concave, laissant penser que nous sommes en présence d'échantillons issus de distributions à queue épaisse, ce qui est effectivement le cas, puisque ce sont des échantillons de loi de Pareto et Lognormale. Enfin, le dernier échantillon a été simulé selon une loi de Weibull avec un paramètre  $\alpha > 1$ . Nous sommes donc en présence d'une distribution à queue fine, comme le confirme la forme convexe du graphe.

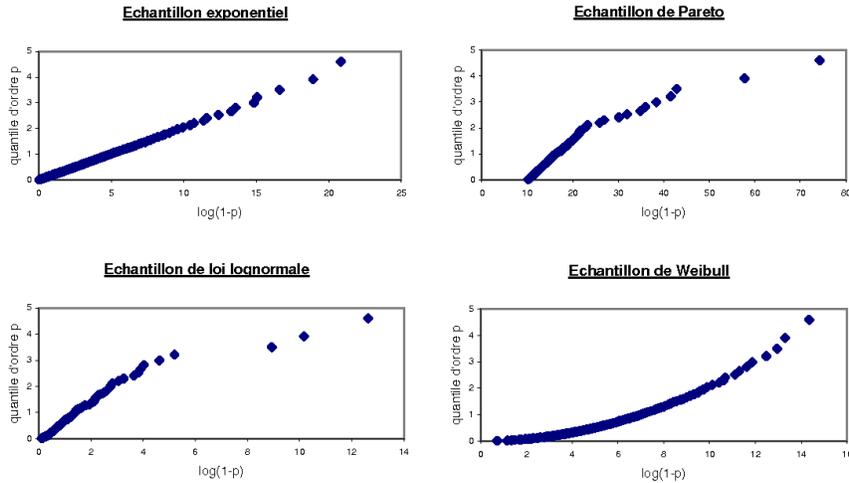


FIG. 4 – QQ Plot exponentiel pour des simulations de loi Exponentielle, de Pareto, Lognormale et de Weibull.

Le QQ Plot peut être adapté à une multitude de lois pour tester graphiquement l'adéquation d'un échantillon à ces lois. Ainsi, dans le cas des lois de Pareto de paramètres  $(\alpha, a)$ , le QQ Pareto plot consiste à tracer les points

$$\left( -\ln\left(1 - \frac{i}{n+1}\right), \ln(x_{(i)}) \right) \quad i = 1, \dots, n$$

On s'attend alors à obtenir une droite de pente  $1/\alpha$ .

Concernant la loi lognormale, on utilise le fait que si l'échantillon suit une loi lognormale, alors son logarithme suit une loi normale. Le QQ Plot Lognormal consiste donc à tracer les points

$$\left( \phi^{-1}\left(\frac{i}{n+1}\right), \ln(x_{(i)}) \right) \quad i = 1, \dots, n$$

Enfin, si on considère la loi de Weibull, le QQ Plot consiste à tracer les points

$$\left( \ln\left(-\ln\left(1 - \frac{i}{n+1}\right)\right), \ln(x_{(i)}) \right) \quad i = 1, \dots, n$$

Les figures suivantes illustrent l'usage des différents QQ-Plots que nous venons de présenter. Nous les avons appliqués à des simulations de loi lognormale.

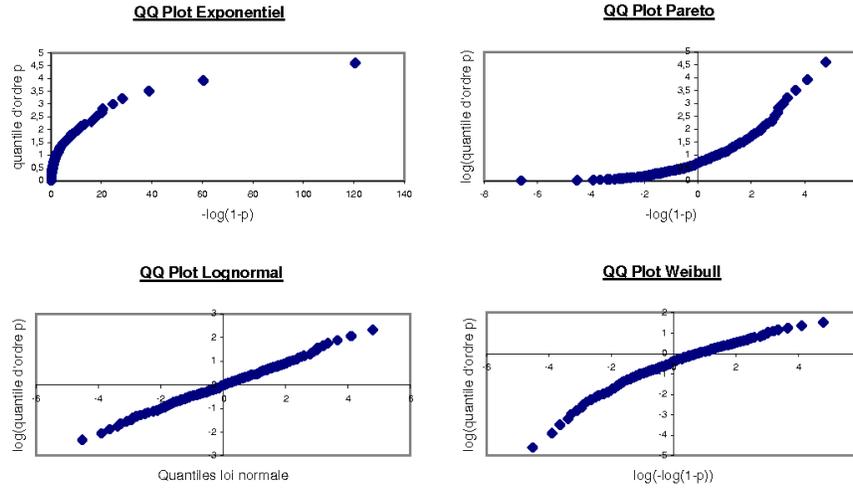


FIG. 5 – QQPlots pour des simulations d’une loi lognormale.

La concavité du QQ Plot exponentiel nous indique que nous sommes en présence d’un échantillon de distribution à queue plus épaisse que la loi Exponentielle. A l’inverse, la convexité du QQ Pareto Plot nous indique que la loi dont est issu l’échantillon est à queue plus fine qu’une loi de Pareto. Finalement, la droite obtenue pour le QQ Plot lognormal laisse peu de doutes quant à l’adéquation de notre échantillon à une loi lognormale.

### 3.2 Mean Excess Function

Un outil graphique utile pour déterminer l’allure de la queue d’une distribution ainsi que pour le choix du seuil à utiliser est la "mean excess function" définie par :

$$e(u) = E(X - u | X > u), \quad 0 \leq u \leq x_F$$

Cette quantité est souvent utilisée en réassurance car elle peut être interprétée comme la perte moyenne espérée dans un traité de priorité  $u$ . On peut faire les remarques suivantes :

- Si  $X \sim Exp(\lambda)$ , alors  $e(u) = \lambda$  pour tout  $u > 0$ , i.e. la mean excess function est horizontale.
- Si  $X \sim GPD(\xi, \sigma)$ , alors  $e(u) = \frac{\sigma}{1-\xi} + \frac{\xi}{1-\xi}u$ , i.e. la mean excess function est affine en  $u$ .

L’estimation des paramètres de la distribution  $GPD$  pose le problème de la détermination du seuil  $u$  qui doit être suffisamment grand pour pouvoir appliquer les résultats de convergence en loi mais pas trop grand afin d’avoir suffisamment de données pour obtenir des estimateurs de bonne qualité. On peut déterminer  $u$  en exploitant le résultat de la deuxième remarque. En effet, on sait que si l’approximation  $GPD$  est valide pour un seuil  $u_0$  alors elle est valide pour tout  $u > u_0$  (par la propriété de stabilité de la  $GPD$ ). Donc pour  $u > u_0$ , la mean excess function  $e(u)$  est linéaire en  $u$  (d’après la remarque précédente). Pour déterminer  $u_0$ ,

on utilise donc le graphique de  $e(u)$ . Plus précisément, on trace  $\hat{e}(u)$  en fonction de  $u$ , avec :

$$\hat{e}(u) = \frac{\sum_{i=1}^n (X_i - u)^+}{\sum_{i=1}^n \mathbb{I}_{X_i > u}}$$

Cela signifie que  $e(u)$  est estimé par la somme des excès au delà du seuil  $u$ , divisé par le nombre de points excédant le seuil  $u$ .

De manière plus générale, le graphe de la mean excess function permet d'avoir une bonne idée du comportement de la queue de distribution de  $X$  : lorsque le graphe de  $e(u)$  est une constante, on a une distribution de type exponentielle, lorsque le graphe de  $e(u)$  est une droite croissante, on a une distribution de type Pareto et lorsque le graphe de  $e(u)$  est entre une constante et une droite croissante, on a une distribution à queue épaisse.

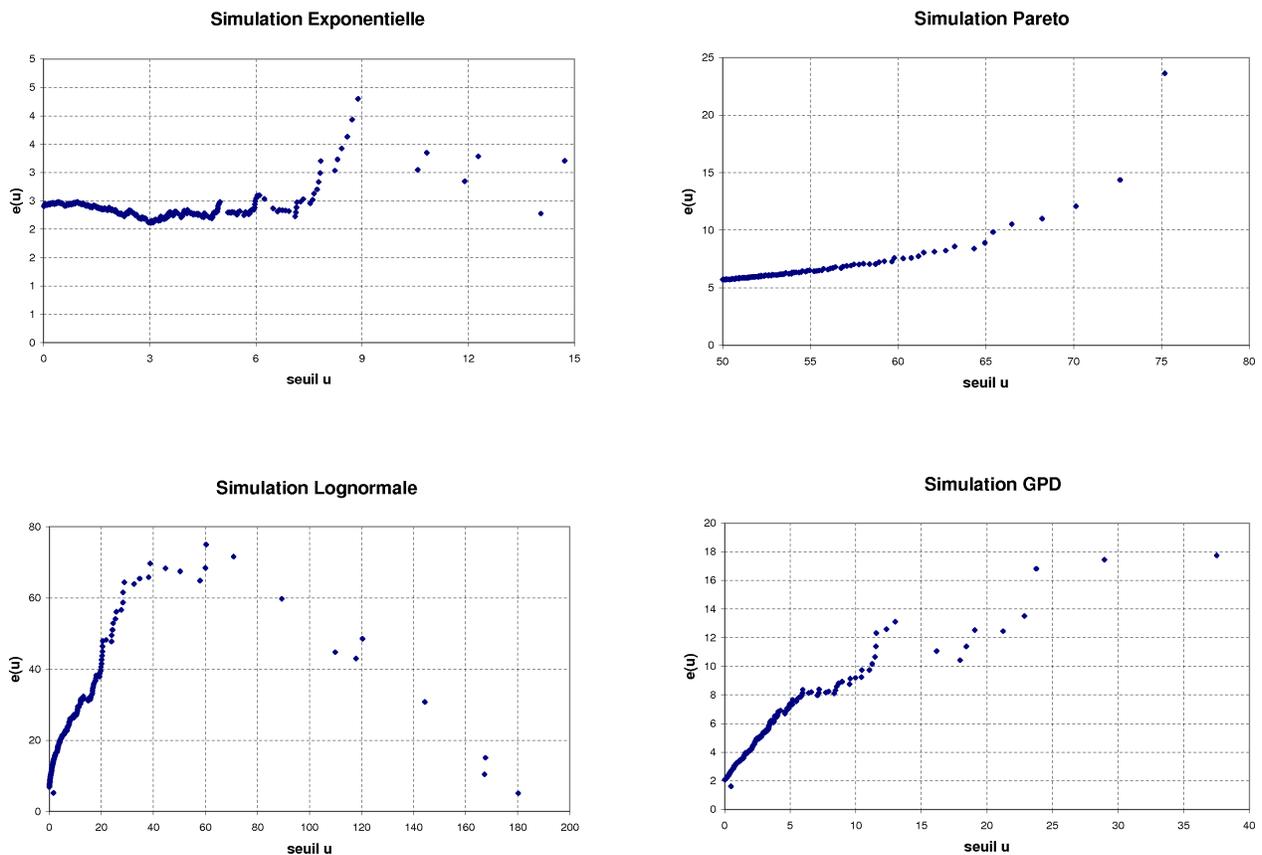


FIG. 6 – Mean excess function pour des simulations de lois paramétriques usuelles.

Dans l'exemple ci-dessus, nous avons utilisé des simulations de quatre lois usuelles effectuées par un logiciel de la SCOR, afin d'observer le comportement de la mean excess function.

Le premier graphique représente la mean excess function d'un échantillon suivant une loi exponentielle de paramètre 2.5. Théoriquement, on doit obtenir une droite horizontale

d'ordonnée 2.5. Empiriquement, on observe bien un alignement horizontal des points, hormis pour les quelques dernières observations.

Concernant les lois de Pareto et GPD, la mean excess function théorique est affine en  $u$ . On constate effectivement une forme affine pour ces deux échantillons. Pour notre simulation de Pareto, l'alignement est très bien vérifié, il est un peu moins bon pour la simulation GPD présentée ici en exemple.

Cet outil graphique est intéressant puisqu'il permet d'avoir une idée sur l'épaisseur de la queue de distribution et sur le seuil d'entrée dans la région extrême. Sur données réelles, le moment où la courbe commence à se comporter comme une droite correspond au seuil, puisqu'une mean excess function affine est caractéristique d'une loi GPD. Cependant, cet outil souffre d'une instabilité pour les dernières observations, ce qui ne facilite pas l'interprétation.

### 3.3 Hill-plot, Pickands-plot

Cette méthode graphique de détermination du seuil repose sur la propriété de stabilité de la loi *GPD*. En effet, si la variable  $[X - u/X > u] \sim GPD(\xi, \sigma)$ , alors  $[X - u'/X > u'] \sim GPD(\xi, \sigma')$ , c'est à dire que le paramètre de queue  $\xi$  est le même pour tout  $u > 0$ . Ainsi, la méthode présentée ici consiste à tracer le graphe des estimateurs  $\hat{\xi}$  obtenus en fonction des seuils  $u$  ou, de manière équivalente, en fonction du nombre  $k$  d'excès considérés. On veut donc tracer le graphe des  $\hat{\xi}_k$  pour tout  $k$ . Or, en fonction du nombre d'excès  $k$ , l'estimateur de Hill s'écrit :

$$\widehat{\xi}_k^{Hill} = \frac{1}{k} \sum_{i=1}^k \log X_{n-i+1:n} - \log X_{n-k:n}$$

et celui de Pickands :

$$\widehat{\xi}_k^{Pick} = \frac{1}{\log 2} \log \frac{X_{n-k:n} - X_{n-2k:n}}{X_{n-2k:n} - X_{n-4k:n}}$$

Plus  $k$  est petit, plus l'intervalle de confiance de  $\xi$  est large (plus il y a de volatilité). Cependant on veut que  $k$  soit le plus petit possible pour avoir le moins de biais. D'où l'idée de choisir le plus petit  $k$  pour lequel l'estimation  $\hat{\xi}_k$  se stabilise.

Nous avons implémenté ces deux estimateurs et nous les avons appliqués dans cette partie à deux échantillons de lois simulées : un échantillon suivant une loi de Pareto de paramètre de queue  $\alpha = 5$  (soit  $\xi = 1/\alpha = 0.2$ ) d'une part et un échantillon suivant une loi GPD de paramètre de queue  $\xi = 1.5$  d'autre part.

L'estimateur de Hill obtenu à partir de ces échantillons (figure 7) donne des valeurs proches des valeurs réelles (autour de 0.2 pour la Pareto et autour de 1.5 pour la GPD). En revanche, l'estimateur de Pickands donne de nettement moins bons résultats. Nous constatons également dans les deux cas que l'estimateur de Pickands est plus volatil que celui de Hill. En effet, l'estimateur de Hill fait intervenir la moyenne des logarithmes des observations, le résultat est donc plus lissé et moins sensible au saut d'une observation. Tout ceci nous amènera à le privilégier par la suite.

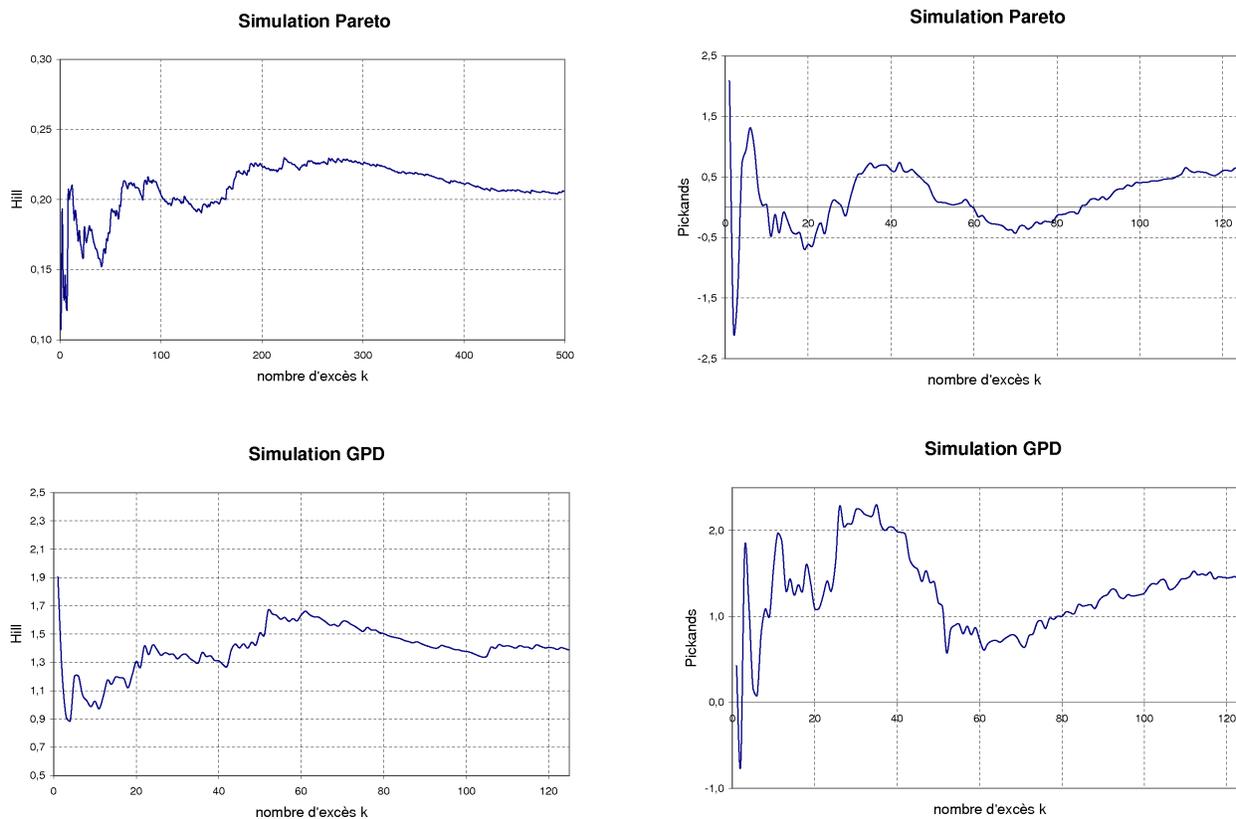


FIG. 7 – Estimateur de Hill (à gauche) et de Pickands (à droite) pour des simulations de lois de Pareto(2;5) et GPD(1.5;1).

De plus, lorsque l'on considère un faible nombre d'excès, l'estimateur de Hill est très volatil, puis il tend à se stabiliser, d'où l'idée de l'utiliser dans le choix du seuil. En effet, on veut le seuil le plus grand possible pour avoir une bonne convergence de la loi des excès vers une GPD (d'après le théorème de Pickands-Balkema-De Haan), mais si le seuil est trop grand, l'estimateur de Hill sera très instable. Le Hill Plot permet donc de voir à partir de quel moment l'estimation devient plus robuste.

### 3.4 Gertensgarbe plot

Cette procédure a été proposée par Gerstengarbe et Werner en 1989. Elle permet de déterminer le point de départ de la région extrême et fournit une estimation du seuil optimal. Plus précisément, pour un échantillon  $x_1, \dots, x_n$  de coûts de sinistres, on considère la série des différences  $\Delta_i = x_{[i]} - x_{[i-1]}$ ,  $i = 2, \dots, n$  de l'échantillon ordonné,  $x_{[1]} \leq \dots \leq x_{[n]}$ . L'idée de la procédure est qu'en entrant dans la zone de sinistres extrêmes, on peut s'attendre à un changement dans le comportement de la série des différences, autrement dit le comportement des  $\Delta_i$  pour les observations extrêmes est certainement différent du comportement des  $\Delta_i$  pour les observations non extrêmes.

Nous cherchons donc à identifier un changement dans une série. Nous allons utiliser pour cela la version séquentielle du test de Mann-Kendall. Pour  $i = 1, \dots, n - 1$ , nous calculons la série  $U_i$  définie par :

$$U_i = \frac{U_i^* - \frac{i(i-1)}{4}}{\sqrt{\frac{i(i-1)(i+5)}{72}}}$$

où  $U_i^* = \sum_{k=2}^i n_k$  et  $n_k$  est le nombre de valeurs  $\Delta_2, \dots, \Delta_k$  inférieures à  $\Delta_k$ . De la même manière, nous calculons une autre série  $U_i^f$  pour la série décroissante des différences,  $\Delta_n, \dots, \Delta_2$ . Le point d'intersection de ces deux séries détermine le seuil d'entrée dans la zone extrême.

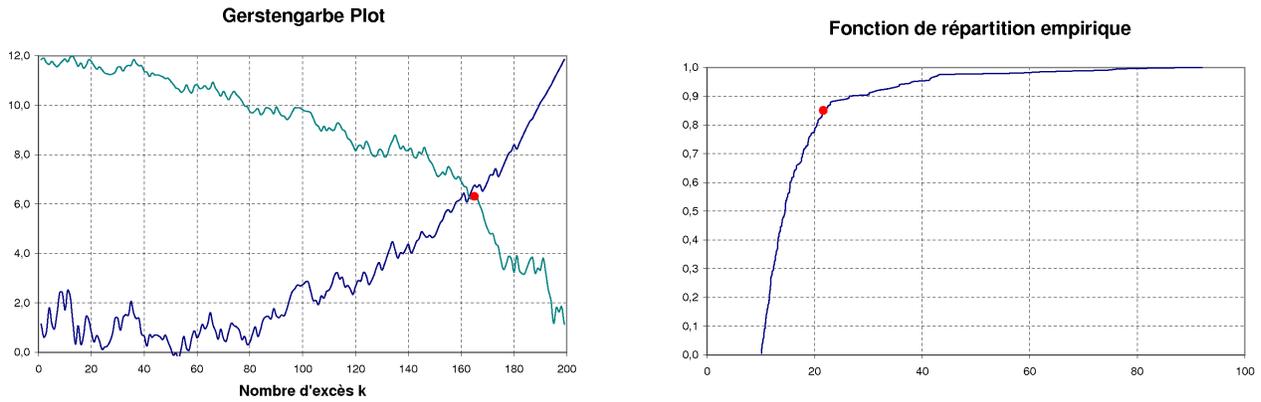


FIG. 8 – Gerstengarbe plot et fonction de répartition empirique pour la simulation d'une loi de Pareto.

Nous avons appliqué cette méthode à un échantillon suivant une loi de Pareto de paramètres 10 et 50. Le point d'intersection des courbes nous indique un nombre d'excès à considérer de 165 observations, ce qui correspond à un seuil de 20,9. Nous constatons de plus en traçant la fonction de répartition empirique de notre échantillon que ce seuil correspond visuellement à un point de cassure de la fonction de répartition. Cela pourrait donc correspondre au seuil d'entrée dans une zone de valeurs extrêmes.

## 4 Modélisation

La deuxième partie a permis de mettre en place les outils mathématiques, notamment la théorie des tests et la théorie des valeurs extrêmes, nécessaires à la résolution de notre problème. La troisième partie a permis de présenter et de mettre en place certains outils graphiques, découlant pour la plupart d'entre eux de la théorie des valeurs extrêmes. Il est donc maintenant possible de présenter des méthodes de modélisation de la sinistralité.

Nous présenterons dans un premier temps un modèle simple à une seule loi. Nous verrons que ce modèle s'adapte souvent mal aux données réelles issues de l'assurance ou de la réassurance, qui rappelons-le comportent souvent deux types de sinistralité : d'une part la sinistralité attritionnelle, qui correspond à des sinistres de faible intensité mais de fréquence importante, c'est-à-dire des sinistres nombreux mais de faible sévérité, d'autre part les sinistres qualifiés d'extrêmes de faible fréquence mais de forte intensité, c'est-à-dire des sinistres survenant rarement voire très rarement mais de forte sévérité.

Dans un deuxième temps, nous présenterons un modèle prenant en compte ces deux types de sinistralité. La sinistralité attritionnelle sera modélisée par une loi classique, une loi lognormale par exemple. La sinistralité extrême, quant à elle, fera l'objet d'une modélisation par une loi de Pareto généralisée, ce qui est cohérent avec les résultats de la théorie des valeurs extrêmes. Le seuil  $u$  entre les deux types de sinistralités sera déterminé par les méthodes graphiques.

Dans un troisième temps, ce modèle à deux lois sera illustré par un exemple. Comme ce modèle ne fait pas partie des modèles classiques et est spécifique à ce mémoire, des données simulées seront utilisées dans cette partie afin de vérifier que l'identification des paramètres sur l'échantillon simulé redonne bien des valeurs proches des valeurs réelles. Ceci permettra de vérifier que notre modèle et notre démarche sont bien robustes.

Enfin, nous présenterons deux extensions du cas précédent. D'abord un modèle où le raccord entre la sinistralité attritionnelle et la sinistralité extrême ne se fait plus en un point, mais sur un intervalle, sur lequel une distribution hybride est utilisée. Le dernier modèle fera intervenir un problème de minimisation contraint. Dans les précédents modèles, le seuil  $u$  est déterminé de manière graphique. Il peut être souhaitable de disposer d'une méthode non graphique, ce qui permet par exemple d'éviter une étape humaine d'interprétation de courbe. Or la littérature actuelle de détermination de seuil est relativement pauvre en ce qui concerne les méthodes non-graphiques et ne fournit pas de solution satisfaisante. Nous formulerons donc une proposition. L'idée générale consiste en la minimisation de la distance entre la sinistralité empirique et la sinistralité modélisée. La distance choisie devra être cohérente avec le problème. Le seuil  $u$  sera celui qui résulte de la minimisation de cette distance.

## 4.1 Cas 1 : Une seule distribution fittée

Dans un premier temps, nous modélisons la sinistralité à l'aide d'une unique distribution. Nous testons ensuite l'adéquation entre la sinistralité empirique et la loi fittée avec le test de Kolmogorov-Smirnov ou d'Anderson-Darling. Ce dernier test est le plus adapté à notre problème car il accorde plus de poids à la queue de distribution que le test de Kolmogorov-Smirnov.

Nous allons considérer deux exemples simples afin d'illustrer ce premier modèle. Le premier échantillon est relatif à des traités proportionnels incendie du marché français et le second est l'analogie pour les marchés de l'Europe du sud. On suppose que les coûts des sinistres composant les échantillons sont bien des réalisations de variables aléatoires indépendantes et identiquement distribuées.

Le premier échantillon comporte 80 observations. Le premier modèle que nous présentons suggère de modéliser la sinistralité à l'aide d'une seule distribution fittée. Nous utilisons pour cela des lois classiques, à savoir la loi lognormale, la loi Gamma, la loi de Weibull, la loi de Pareto, et à titre comparatif et en prévision des modèles suivants, la loi de Pareto généralisée.

Les paramètres de ces différentes lois sont estimés à l'aide de la méthode du maximum de vraisemblance, puis l'adéquation entre la distribution empirique et la loi usuelle est testée à l'aide des tests de Kolmogorov-Smirnov et d'Anderson-Darling. Pour tous les tests, le seuil de significativité est 5%. Les résultats sont synthétisés dans le tableau suivant :

Loi usuelle	Test de Kolmogorov-Smirnov	Test d'Anderson-Darling
Lognormale	Acceptation	Acceptation
Gamma	Rejet	Rejet
Weibull	Acceptation	Acceptation
Pareto	Rejet	Rejet
Pareto généralisée	Rejet	Rejet

On constate que les lois lognormales et de Weibull sont acceptées à la fois par le test de Kolmogorov-Smirnov et d'Anderson-Darling. Les autres lois sont quant à elles rejetées. Ce premier portefeuille de sinistres peut donc être modélisé soit par une loi lognormale, soit par une loi de Weibull.

Graphiquement, les fonctions de répartition empirique et modélisées sont représentées sur la figure suivante.

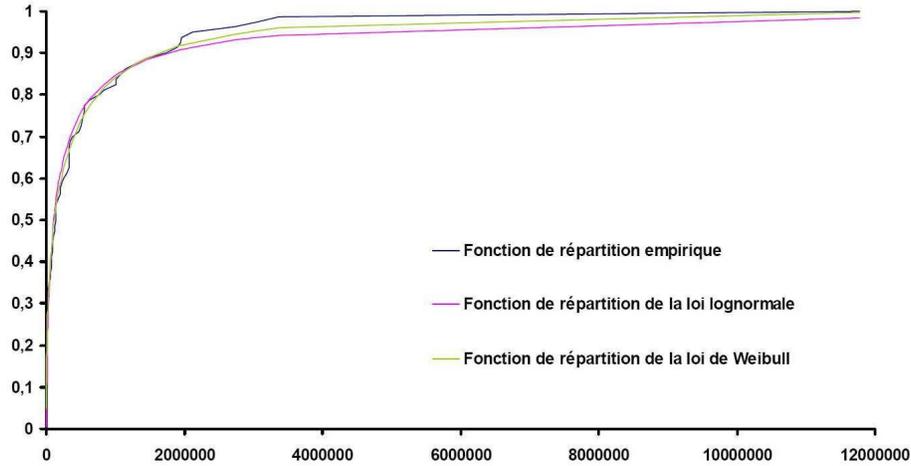


FIG. 9 – Fonctions de répartition empirique et modélisées pour le premier échantillon.

Toutefois, le cas où l’adéquation globale est réalisée par une seule loi est assez rare. Nous allons maintenant considérer le deuxième échantillon, composé de 753 observations. Le même protocole est utilisé que pour le premier échantillon. Les résultats des tests d’adéquation sont les suivants :

Loi usuelle	Test de Kolmogorov-Smirnov	Test d’Anderson-Darling
Lognormale	Rejet	Rejet
Gamma	Rejet	Rejet
Weibull	Acceptation	Rejet
Pareto	Rejet	Rejet
Pareto généralisée	Rejet	Rejet

Ainsi, toutes les lois usuelles présentées ici sont rejetées, à l’exception de la loi de Weibull qui est acceptée par le test de Kolmogorov-Smirnov.

Nous avons tracé la fonction de répartition empirique ainsi que celle de la loi de Weibull (figure 10). Visuellement, nous pouvons constater que pour les petits sinistres, dans cet exemple inférieurs à 100 000, l’adéquation à une loi de Weibull est très bien réalisée. Par contre, au delà de ce seuil, la modélisation par une loi de Weibull ne semble plus justifiée. Le test d’Anderson-Darling accordant plus de poids à la queue de distribution que ne le fait le test de Kolmogoriv-Smirnov, nous comprenons alors pourquoi le test de Kolmorov-Smirnov conduit à accepter l’adéquation, alors que le test d’Anderling-Darling conduit à un rejet.

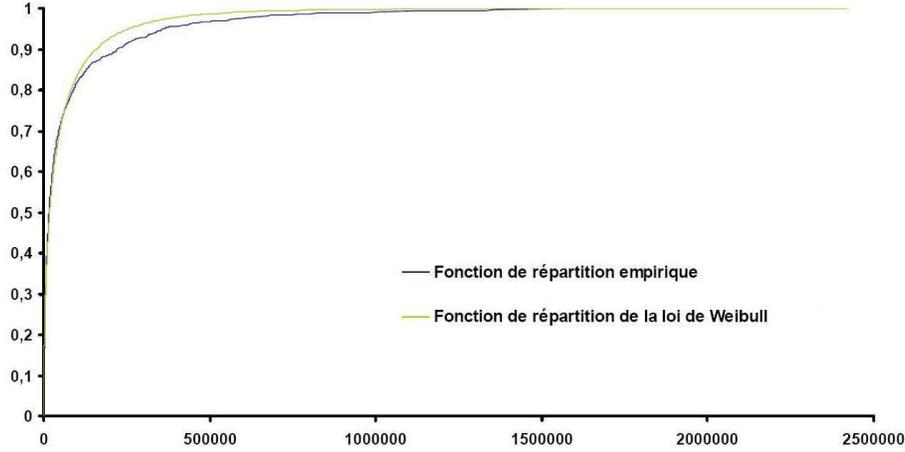


FIG. 10 – Fonctions de répartition empirique et modélisée pour le second échantillon.

Ce deuxième exemple est en fait assez représentatif des données issues de l'assurance ou de la réassurance : la sinistralité est de nature complexe, c'est pourquoi un ajustement par une seule distribution est souvent impossible ou non optimal. Dans ces conditions, nous avons recours à une modélisation mettant en jeu deux distributions.

## 4.2 Cas 2 : Deux distributions pour modéliser deux types de sinistralités

La modélisation par une seule distribution étant généralement peu convaincante avec les données issues de la réassurance, on peut avoir recours à deux distributions. La sinistralité attritionnelle peut être modélisée par une loi classique, par exemple une loi lognormale. Pour la sinistralité extrême, le théorème de Picands-Balkema-De Hann, nous suggère d'utiliser une loi de Pareto généralisée.

Il y a trois types de questions auxquelles il faut répondre. Premièrement, comment choisir le seuil de raccord  $u$  ? Deuxièmement, comment estimer les paramètres des lois, notamment les paramètres de la loi de Pareto généralisée ? Troisièmement, comment à l'aide de deux lois proposer une distribution unique, et essayer si possible que cette loi présente certaines propriétés de régularité, afin de ne pas avoir de saut lors de la tarification des tranches ?

Le seuil  $u$  peut être déterminé par diverses méthodes graphiques présentées en détail dans la troisième partie de ce mémoire. Nous allons à présent proposer une procédure de détermination du seuil.

La détermination peut être faite par le graphe de la mean excess function. Le seuil  $u$  est choisi lorsque la courbe devient linéaire. La détermination peut également être faite par

le graphique de Hill. Le seuil  $u$  est choisi lorsque l'estimateur de Hill se stabilise. Enfin, la détermination peut être faite grâce au Gertensgarbe Plot. Cette méthode, même si elle ne découle pas de la théorie des valeurs extrêmes, présente l'avantage de présenter une réponse généralement non ambiguë, car le seuil est situé à l'intersection de deux courbes. Les méthodes graphiques du Hill Plot et de la mean excess function découlent de la théorie des valeurs extrêmes. Cependant le seuil est à choisir quand l'estimateur de Hill est stable ou quand la fonction moyenne en excès prend l'allure d'une droite, or sur données réelles il peut être difficile de déterminer cette zone et le choix est en fait déterminé par la subjectivité humaine. Un moyen de choisir le seuil peut donc consister en ce protocole provisoire : dans un premier temps, lire le seuil grâce au Gertensgarbe Plot, qui donne une réponse claire car le seuil est donné par l'intersection des deux courbes, et dans un second temps vérifier que ce résultat est cohérent avec le Hill Plot et le graphe de la mean excess function.

Le seuil  $u$  étant déterminé, il faut maintenant estimer les paramètres de deux lois. L'estimation des paramètres de la première loi, par exemple une loi lognormale, relève des méthodes classiques telles que le maximum de vraisemblance ou la méthode des moments. La loi de Pareto généralisée comporte trois paramètres. Le premier paramètre définit le début du support de la distribution. On prend donc le seuil  $u$  précédemment déterminé par les méthodes graphiques. Il reste donc deux paramètres à estimer : le paramètre de queue  $\xi$  et le paramètre d'échelle  $\sigma$ . Il est possible de déterminer conjointement ces deux paramètres soit par la méthode des moments, soit par la méthode du maximum de vraisemblance.

Pour vérifier que ces deux lois conviennent bien, il est nécessaire d'effectuer des tests d'adéquation. Pour la sinistralité attritionnelle, on peut utiliser le test de Kolmogorov-Smirnov ou d'Anderson-Darling. Pour la sinistralité extrême, il faudrait également effectuer un test d'adéquation, cependant la littérature ne nous offre pas de réponse satisfaisante. Le théorème de Picands-Balkema-De Hann est le garant du bien fondé de l'utilisation de la distribution de Pareto généralisée pour la modélisation des excès après le seuil, de la même manière que d'après la théorème central limite, on peut modéliser par une loi normale la moyenne de  $n$  variables aléatoires indépendantes et identiquement distribuées, sous des hypothèses convenables et pour  $n$  assez grand.

Nous avons maintenant deux distributions, une pour la sinistralité attritionnelle et une pour la sinistralité extrême, dont nous avons estimé les paramètres. Il reste donc à créer une distribution globale. La deuxième distribution modélise les sinistres au delà du seuil  $u$ , il faut donc tronquer la première en ce point. Nous avons donc une distribution usuelle tronquée sur  $[0, u]$ , puis la distribution Pareto généralisée à support sur  $[u, +\infty[$ . Le fait de juxtaposer les fonctions de répartition sans précaution conduit à une aberration, puisque en  $u$ , nous passerions d'une valeur strictement positive à 0 et nous n'aurions plus une fonction croissante. Nous voulons une fonction de répartition globale, c'est-à-dire une fonction croissante sur  $[0, +\infty[$ , qui vaut 0 en 0 (ou au point de début de la modélisation s'il y a des sinistres tronqués à gauche) et qui tend vers 1 en  $+\infty$ . Lors de la tarification des tranches, nous aimerions que la prime pure ne connaisse pas de saut lorsque la priorité ou la portée varie très légèrement, donc nous cherchons également une fonction continue. Il semble alors convenable de supposer l'existence d'une densité ; la fonction de répartition, en temps qu'intégrale entre

0 et  $x$  de cette densité, sera donc continue.

Dans un cadre général, si  $f_1$  et  $f_2$  sont deux densités et que l'on cherche à en créer une troisième  $f_3$  à l'aide d'une combinaison linéaire des deux premières, on peut considérer :

$$f_3 = \alpha f_1 + (1 - \alpha) f_2$$

avec  $\alpha$  dans  $[0, 1]$  qui est un coefficient de pondération entre  $f_1$  et  $f_2$ , et on a bien une fonction d'intégrale 1. En effet :

$$\int f_3 = \int (\alpha f_1 + (1 - \alpha) f_2) = \alpha \int f_1 + (1 - \alpha) \int f_2 = \alpha + (1 - \alpha) = 1$$

*Remarque :* Cette formule peut se trouver dans un cadre probabiliste. Supposons, ce qui sera le cas ci-dessous, que  $f_1$  modélise la sinistralité attritionnelle et a pour support  $[0, u]$  et que  $f_2$  modélise la sinistralité extrême et a pour support  $[u, +\infty[$ . Notons respectivement  $F_1$ ,  $F_2$  et  $F_3$  les fonctions de répartition associées aux densités  $f_1$ ,  $f_2$  et  $f_3$ , et  $\mathbb{P}_1$ ,  $\mathbb{P}_2$  et  $\mathbb{P}_3$  les probabilités associées à ces mêmes densités. Nous avons  $\mathbb{P}_3(X \leq x | X < u) = \mathbb{P}_1(X \leq x) = F_1(x)$  et  $\mathbb{P}_3(X \leq x | X \geq u) = \mathbb{P}_2(X \leq x) = F_2(x)$ . Nous avons alors :

$$\begin{aligned} F_3(x) = \mathbb{P}_3(X \leq x) &= \mathbb{P}_3(X \leq x \cap X < u) + \mathbb{P}_3(X \leq x \cap X \geq u) \\ &= \mathbb{P}_3(X < u) \mathbb{P}_3(X \leq x | X < u) + (1 - \mathbb{P}_3(X < u)) \mathbb{P}_3(X \leq x | X \geq u) \\ &= \mathbb{P}_3(X < u) F_1(x) + (1 - \mathbb{P}_3(X < u)) F_2(x) \end{aligned}$$

Notons  $\alpha = \mathbb{P}_3(X < u)$ , nous obtenons :

$$F_3(x) = \alpha F_1(x) + (1 - \alpha) F_2(x)$$

Et enfin par dérivation :

$$f_3(x) = \alpha f_1(x) + (1 - \alpha) f_2(x)$$

Nous retrouvons bien la formule définissant  $f_3$  à partir de  $f_1$  et de  $f_2$ .

Dans le cas qui nous intéresse, nous allons modéliser la sinistralité globale par une loi usuelle tronquée, la loi usuelle retenue étant soit une loi lognormale soit une loi de Weibull, puis par une loi de Pareto généralisée après le seuil  $u$ .  $f_1$  sera donc la densité d'une loi usuelle tronquée et  $f_2$  la densité d'une loi de Pareto généralisée.

Notons  $f^{Usuelle}$  la densité de la loi usuelle choisie (loi lognormale ou loi de Weibull),  $f^{UsuelleTronque}$  la densité de la loi usuelle tronquée dont le support est  $[0, u]$ ,  $f^{GPD}$  la densité de la loi de Pareto généralisée de seuil  $u$ , donc dont le support est  $[u, +\infty[$  et  $f^{Globale}$  la densité utilisée pour modéliser la sinistralité globale. Les fonctions de répartition associées sont notées de manière analogue  $F^{Usuelle}$ ,  $F^{UsuelleTronque}$ ,  $F^{GPD}$ , et  $F^{Globale}$ .

Nous utilisons maintenant la formule de pondération définie ci-dessus, avec  $f_1 = f^{UsuelleTronque}$ ,  $f_2 = f^{GPD}$  et  $f_3 = f^{Globale}$ . La densité modélisant la sinistralité globale a donc pour expression :

$$f^{Globale} = \alpha f^{UsuelleTronque} + (1 - \alpha) f^{GPD}$$

Il est également possible d'exprimer la densité de la loi usuelle tronquée en fonction de la densité de la loi usuelle. En effet, sur  $[0, u]$ , les deux densités sont proportionnelles :

$$f^{UsuelleTronque} = \beta f^{Usuelle} \cdot \mathbb{I}_{[0,u]}$$

En prenant l'intégrale sur  $[0, +\infty[$ , on obtient :

$$1 = \beta \int_0^u f^{Usuelle}$$

D'où :

$$\beta = \frac{1}{\int_0^u f^{Usuelle}} = \frac{1}{F^{Usuelle}(u)}$$

Ainsi :

$$f^{UsuelleTronque} = \frac{f^{Usuelle} \cdot \mathbb{I}_{[0,u]}}{F^{Usuelle}(u)}$$

L'expression de la densité modélisant la sinistralité globale est donc maintenant :

$$f^{Globale} = \alpha \frac{f^{Usuelle} \cdot \mathbb{I}_{[0,u]}}{F^{Usuelle}(u)} + (1 - \alpha) f^{GPD}$$

Ainsi, en choisissant  $\alpha = F^{Usuelle}(u)$ , les densités  $f^{Globale}$  et  $f^{Usuelle}$  sont égales sur  $[0, u]$  et les fonctions de répartition associées le sont également sur ce même intervalle. L'avantage de ce choix de  $\alpha$  est donc qu'il permet de conserver exactement la même fonction de répartition qu'avec une unique loi usuelle sur  $[0, u]$ , c'est-à-dire que la partie correspondant à la sinistralité attritionnelle reste la même dans le modèle à deux lois que dans le modèle à une seule loi. Le coefficient de pondération  $\alpha = F^{Usuelle}(u)$  ici choisi correspond donc à la probabilité qu'un sinistre soit plus petit que le seuil  $u$  en utilisant la loi de probabilité modélisant la sinistralité attritionnelle. D'autres choix pour le coefficient de pondération  $\alpha$  auraient également été possibles.

Nous aboutissons donc à :

$$f^{Globale} = f^{Usuelle} \cdot \mathbb{I}_{[0,u]} + (1 - F^{Usuelle}(u)) f^{GPD}$$

Nous avons bien une densité car  $\int f^{Globale} = 1$ . De plus, sur  $[0, u]$ , la fonction de répartition associée à cette densité est égale à celle de la loi usuelle :  $F^{Globale} = F^{Usuelle}$  sur  $[0, u]$ . En outre,  $F^{Globale}$  est bien continue en temps qu'intégrale entre 0 et  $x$  de sa densité, fonction mesurable, donc il n'y a pas de problème de saut lors de la tarification de tranches.

Par intégration, on obtient la fonction de répartition  $F^{Globale}$  associée à la densité  $f^{Globale}$  :

$$F^{Globale}(x) = F^{Usuelle}(\min(x, u)) + (1 - F^{Usuelle}(u))F^{GPD}(x)$$

$F^{Globale}$  s'exprime également sous la forme :

$$F^{Globale}(x) = \begin{cases} F^{Usuelle}(x) & \text{si } x < u \\ F^{Usuelle}(u) + (1 - F^{Usuelle}(u))F^{GPD}(x) & \text{si } x \geq u \end{cases}$$

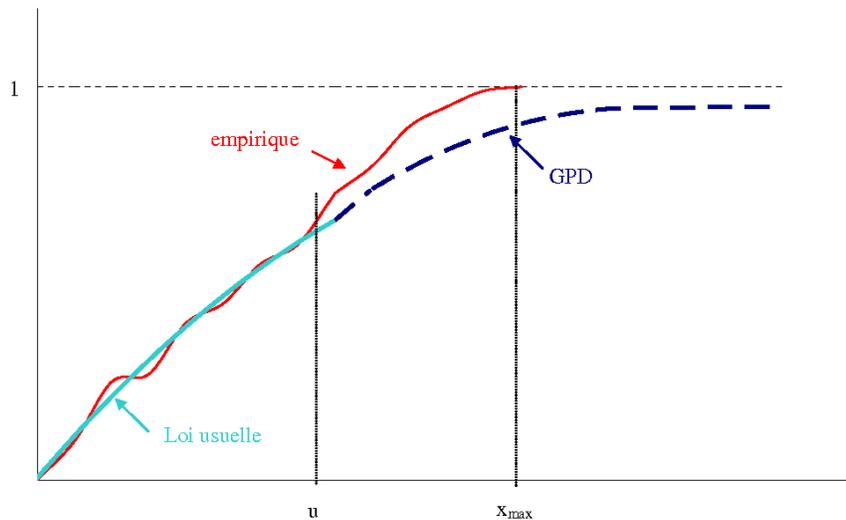


FIG. 11 – Illustration de la modélisation de la sinistralité par deux distributions. La fonction de répartition globale est composée d'une loi usuelle sur  $[0, u]$  puis d'une GPD sur  $[u, +\infty[$ .

*Remarque :* Une alternative à ce modèle peut consister à prendre  $F^{Globale}$  de classe  $C^1$ , pour encore plus de régularité au niveau de la tarification. Dans ce cas, il faut imposer  $f^{Globale}$  de classe  $C^0$  (c'est-à-dire continue). Au point  $u$ , nous aurons donc égalité entre  $f^{Usuelle}$  et  $f^{GPD}$ . Nous avons alors une équation reliant les deux paramètres de la loi usuelle et les paramètres de queue  $\xi$  et d'échelle  $\sigma$  de la loi de Pareto généralisée. Nous pouvons alors par exemple déterminer les deux paramètres de la loi usuelle, déterminer l'un des paramètres de la GPD (par exemple le paramètre de queue par l'estimateur de Hill), le dernier paramètre de la GPD est donc déterminé par l'équation supplémentaire.

### 4.3 Application du modèle à deux lois sur données simulées

Nous avons précédemment défini un modèle à deux lois permettant de modéliser la sinistralité dans son ensemble. La sinistralité attritionnelle est modélisée par une loi usuelle comme une loi lognormale ou une loi de Weibull tandis que la sinistralité extrême est modélisée par une loi de Pareto généralisée. Ce modèle comporte cinq paramètres : le seuil  $u$ , les deux paramètres de la loi usuelle et les deux paramètres restant pour la loi de Pareto

généralisée, qui sont le paramètre de queue  $\xi$  et le paramètre d'échelle  $\sigma$ .

Nous souhaiterions savoir si les paramètres sont bien identifiables, notamment le seuil  $u$ . Pour cela, nous allons considérer un échantillon contenant des réalisations de la distribution définie en 4.2 (modèle à deux lois). Nous avons simulé 1 000 réalisations de cette distribution à deux lois, la loi choisie pour la sinistralité attritionnelle étant une loi lognormale. La densité associée à cette distribution a donc pour expression :

$$f^{Globale} = f^{LN} \cdot \mathbb{I}_{[0,u]} + (1 - F^{LN}(u))f^{GPD}$$

Les cinq paramètres dans cet exemple sont les suivants : le seuil  $u$  a été fixé à 100 000, les deux paramètres de la loi lognormales sont  $\mu^{LN} = 10$  et  $\sigma^{LN} = 1,5$  et les deux paramètres de la loi Pareto généralisée sont  $\xi = 0,7$  et  $\sigma = 80\,000$ .

Nous allons essayer de mettre en place un protocole d'étude sur ces données simulées, en utilisant les considérations de cette partie 4 et les outils de la partie 3. Ce protocole d'étude sera repris avec plus de détails dans la partie 5 sur deux jeux de données réelles.

## Ajustement GEV

Tout d'abord, afin de savoir si nous sommes dans le cas  $\xi = 0$  (cas Gumbel) ou  $\xi > 0$  (cas Fréchet), nous allons tracer les graphiques pour la loi de Gumbel et de Fréchet définis théoriquement dans la section 2.3.2. Rappelons qu'un alignement des points correspond pour le premier graphe (figure 12) à l'acceptation de l'hypothèse  $\xi = 0$  (cas Gumbel) et pour le second (figure 13) de l'hypothèse  $\xi > 0$

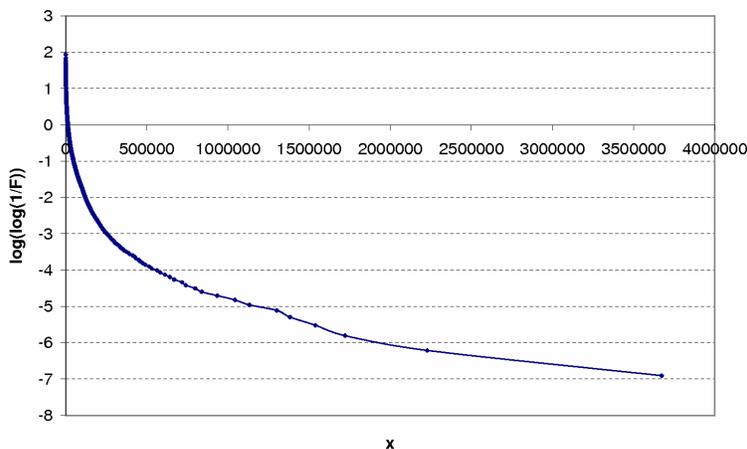


FIG. 12 – Test de  $\xi = 0$ .

Nous constatons aisément que la première courbe est strictement convexe, d'où le rejet de l'hypothèse  $\xi = 0$ . De plus, la dernière partie de la seconde courbe, qui correspond à la queue de distribution, est assimilable à une droite, d'où l'acceptation de l'hypothèse  $\xi > 0$ .

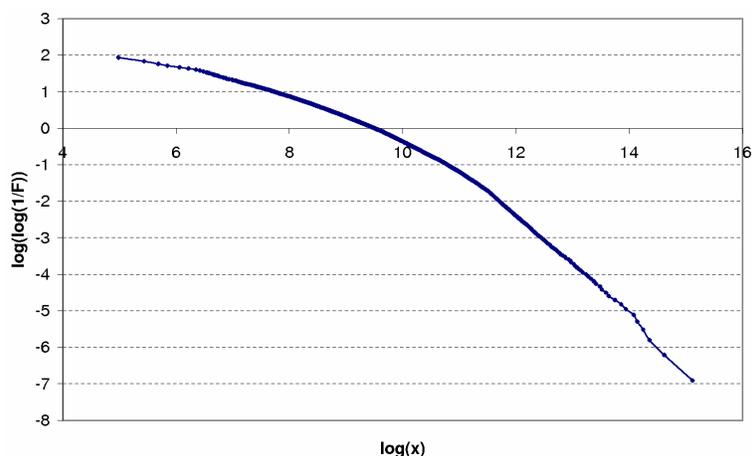


FIG. 13 – Test de  $\xi > 0$ .

### Détermination du seuil

Pour déterminer le seuil d'entrée dans la zone extrême, nous allons appliquer le protocole défini au début de la section précédente. Nous allons donc tout d'abord tracer le Gerstengarbe Plot (figure 14), puis nous regarderons la cohérence de ces résultats avec le Hill plot (figure 16) et le graphe de la mean excess function (figure 17).

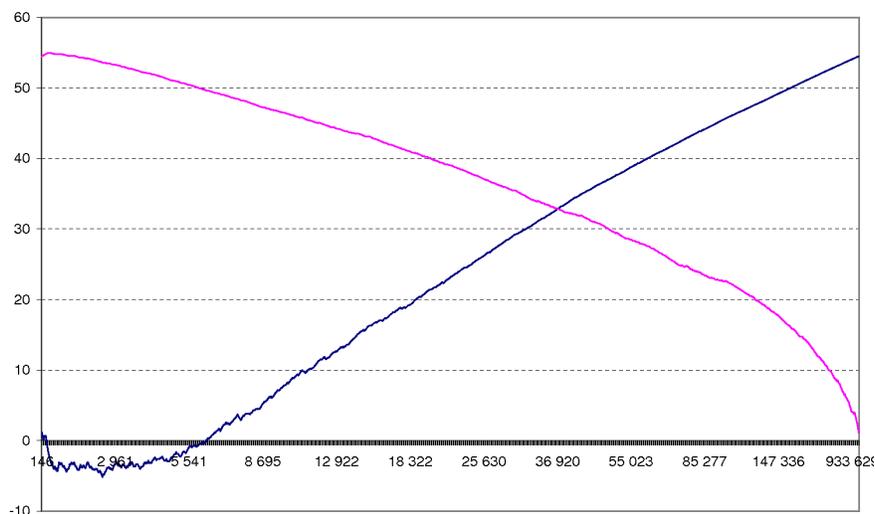


FIG. 14 – Gerstengarbe plot en fonction du seuil sur données simulées.

Le résultat obtenu par la méthode de Gerstengarbe n'est pas satisfaisant. En effet, il suggère un seuil d'environ 36 900, alors que la vraie valeur est de 100 000 !

De plus, cette méthode n'est pas du tout robuste. En effet, nous l'avons à nouveau appliquée à notre échantillon simulé en enlevant les 300 premières observations. Le seuil obtenu s'approchait alors des 60 000.

Ainsi, la méthode Gerstengarbe ne semble pas convenir à notre problématique. Peut-être

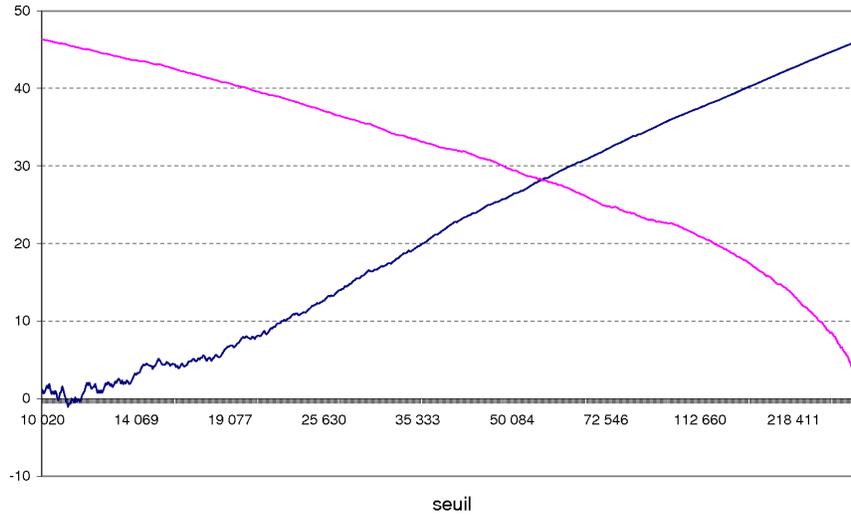


FIG. 15 – Gerstengarbe plot en fonction du seuil sur données simulées tronquées.

est-ce dû au fait qu'elle ne repose pas sur la théorie des valeurs extrêmes ? Dans tous les cas, nous avons décidé de ne pas l'utiliser et continuons notre protocole de détermination du seuil avec le Hill Plot.

Comme nous sommes dans le cas  $\xi > 0$ , nous pouvons utiliser l'estimateur de Hill qui est bien défini<sup>5</sup>. L'étude des zones de stabilité de cet estimateur en fonction du seuil considéré permet de déterminer un seuil adapté au problème. Le graphique suivant représente l'estimateur de Hill en fonction du seuil et du nombre d'excès considérés.

<sup>5</sup>Nous rappelons que l'estimateur de Hill n'est défini que pour des lois ayant un paramètre de queue  $\xi > 0$ . L'estimateur de Hill n'est donc pas valide dans le cas  $\xi = 0$ . D'où l'intérêt d'effectuer un ajustement GEV au préalable, afin de vérifier que nous sommes bien dans le cas  $\xi > 0$

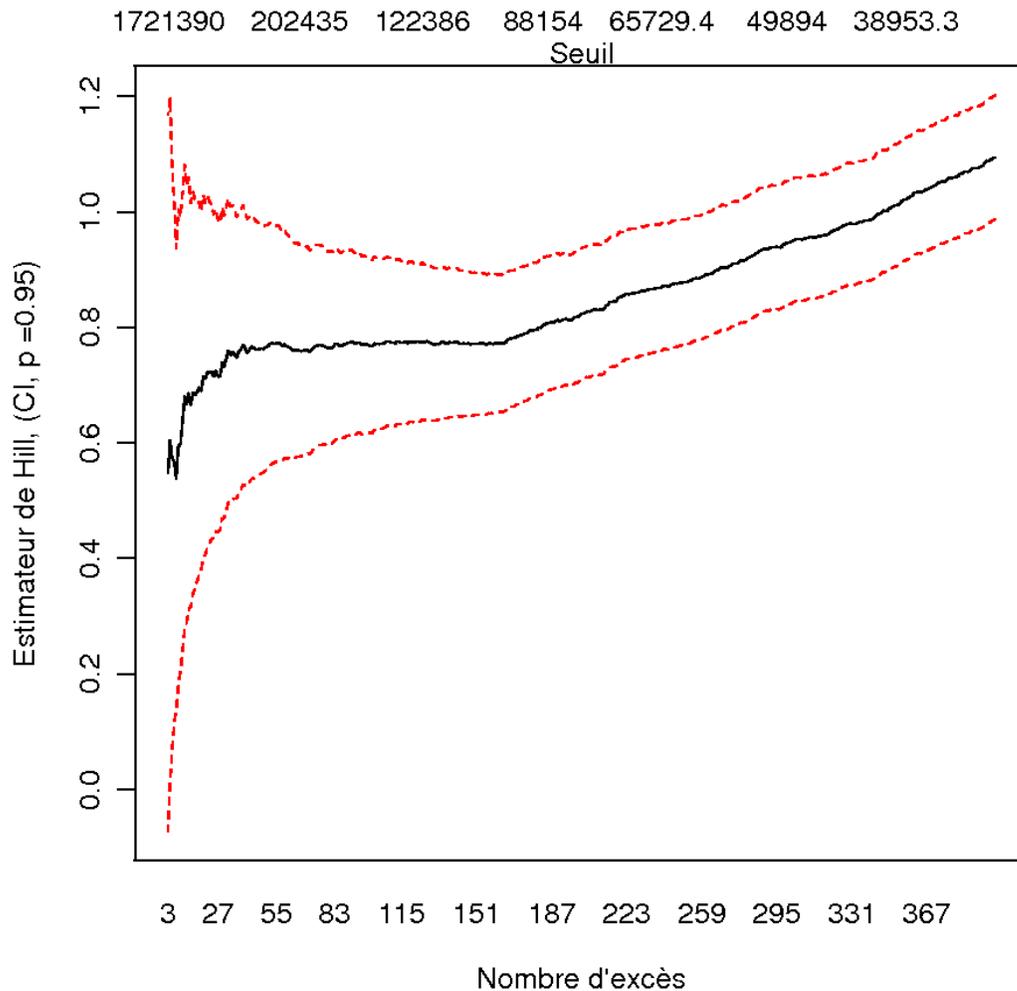


FIG. 16 – Estimateur de Hill en fonction du seuil et du nombre d’excès sur données simulées.

La situation sur données simulées est très nette et en accord avec la théorie. Au tout début de la courbe, quand le seuil est élevé, l’estimateur de Hill est très volatil parce que le nombre d’excès considérés n’est pas suffisant. C’est ce que l’on constate sur le graphique quand on considère moins de 50 excès. Par contre, sur la partie suivante, l’estimateur de Hill est très stable et constant. A partir d’environ 160 excès, le comportement de la courbe change radicalement : la courbe s’identifie à une fonction affine strictement croissante.

Le point de cassure, situé à environ 160 excès, correspond à un seuil de 100 000, soit le seuil théorique utilisé pour les simulations. La partie gauche de la courbe, avant 160 excès, correspond à la partie de la distribution qui utilise une loi de Pareto généralisée tandis que la partie droite correspond aux réalisations de la loi lognormale.

Sur cet échantillon simulé, nous constatons donc clairement que la partie de l’échantillon provenant de réalisations de la loi de Pareto généralisée conduit à un estimateur de Hill constant, hormis pour les premiers excès où l’estimateur de Hill est trop volatil, alors que la partie de l’échantillon provenant de réalisations de la loi lognormale conduit à une portion de

courbe strictement croissante.

Le Hill plot semble donc un outil pertinent pour déterminer le seuil d'entrée dans la zone des sinistres extrêmes.

Afin de vérifier la cohérence du choix du seuil basé sur l'estimateur de Hill, nous pouvons utiliser un autre outil, la mean excess function. Sa représentation graphique est donnée sur la figure suivante.

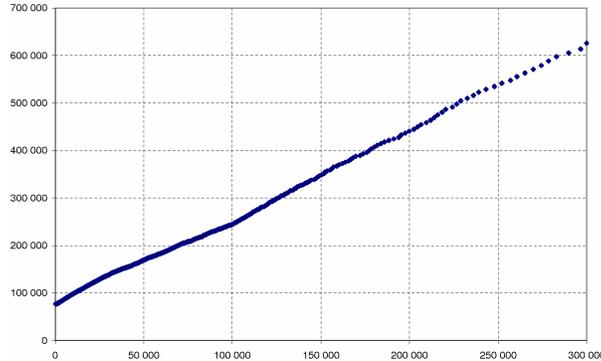


FIG. 17 – Mean excès function sur données simulées.

Nous constatons une légère cassure à 100 000. D'autre part, avant 100 000, la courbe a une légère tendance convexe, alors qu'après les points sont alignés, ce qui est caractéristique de la loi de Pareto généralisée. La mean excès function, même si elle permet de conclure, apporte sur cet exemple une réponse moins nette que l'estimateur de Hill.

De manière générale, il ne faudrait pas croire à la lecture des graphiques de l'estimateur de Hill sur données simulées que la détermination du seuil soit aisée. Nous verrons sur données réelles que c'est loin d'être le cas.

### Estimation des paramètres pour la sinistralité attritionnelle et pour la sinistralité extrême

Maintenant que le seuil a été identifié, il est possible d'estimer les quatre autres paramètres du modèle. La méthode utilisée est le maximum de vraisemblance. Les résultats sont synthétisés dans le tableau suivant.

Paramètre	Valeur réelle	Valeur estimée
$\mu^{LN}$	10	10,0094
$\sigma^{LN}$	1,5	1,5199
$\xi$	0,7	0,7620
$\sigma$	80 000	78 573

Nous constatons que les valeurs estimées sont proches des valeurs réelles. Sur données simulées, il est donc possible d'identifier les paramètres du modèle.

Pour le seuil, comme une méthode graphique est utilisée pour la détermination, il est difficile de quantifier la qualité de l'estimation. Cependant, l'estimateur de Hill permet une identification très nette sur ces données simulées.

## Modélisation globale

Les fonctions de répartition empirique et modélisée relatives à la sinistralité dans son ensemble sont représentées sur la figure suivante.

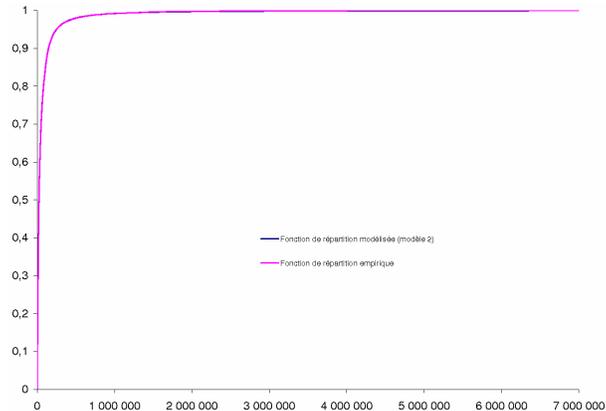


FIG. 18 – Fonction de répartition empirique et modélisée.

Sans surprise, l'adéquation est visuellement excellente. Le modèle à deux lois présenté ici semble donc robuste et nous pourrions dans la partie 5 l'appliquer sur données réelles.

## 4.4 Approches alternatives

Après avoir présenté deux modèles, un modèle classique à une seule loi et un modèle à deux lois, nous allons présenter quelques pistes d'extensions possibles. Ces modèles seront présentés uniquement de manière théorique et ne feront pas l'objet ici d'applications. Ils peuvent constituer des voies futures de recherche.

### 4.4.1 Cas plus élaboré à trois lois

L'approche à deux distributions pour modéliser les deux types de sinistralité est une bonne solution, notamment pour améliorer l'ajustement au niveau des queues de distribution, c'est à dire des gros sinistres. Cependant, le passage d'une distribution à l'autre au niveau du seuil  $u$  peut paraître un peu "brutal". D'où l'idée d'un cas plus élaboré à trois lois : deux lois pour modéliser chaque type de sinistralité et une loi hybride effectuant le raccord entre les deux, permettant ainsi un ajustement plus en douceur. L'utilisation d'une loi hybride peut permettre de mieux prendre en compte la sinistralité entre le seuil et le plus gros sinistre.

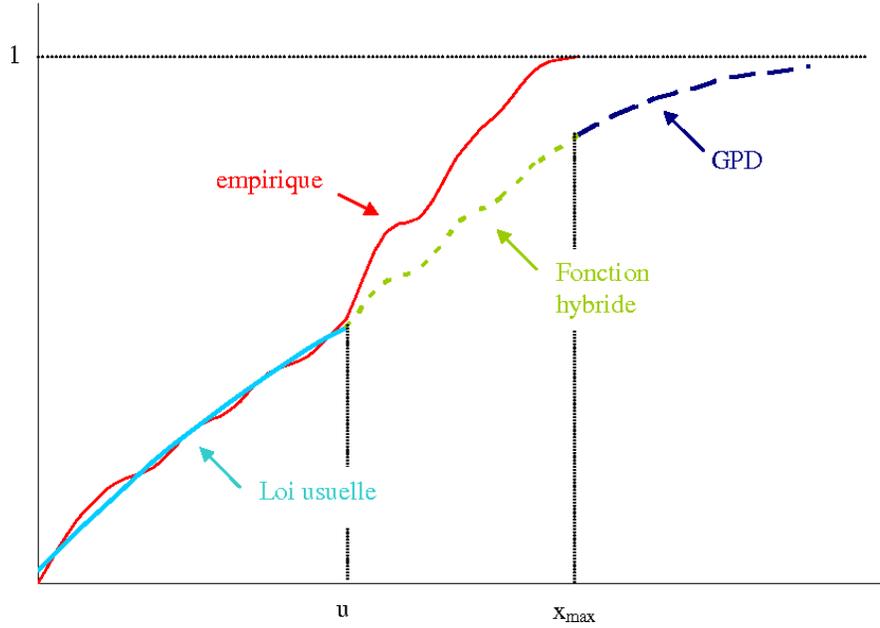


FIG. 19 – Illustration de la modélisation de la sinistralité par trois distributions. La fonction de répartition globale est composée d'une loi usuelle sur  $[0, u]$ , d'une loi hybride sur  $[u, X_{max}]$  puis d'une GPD sur  $[X_{max}, +\infty[$ .

Une forme de la fonction hybride a été développée à la SCOR, menant à la fonction de répartition  $F^{hybrid}$  suivante :

- Si  $0 \leq x \leq u$ ,  $F^{hybrid}(x) = F^{Usuelle}(x)$
- Si  $u < x \leq x_{max}$ ,  $F^{hybrid}(x) = F^{Usuelle}(u) + [F^{GPD}(x_{max}) \times (F^{Usuelle}(x) - F^{Usuelle}(u))]^{p(x)}$ ,  
avec

$$p(x) = \frac{\ln(F^{Usuelle}(x) \times (F^{GPD}(x_{max}) - F^{Usuelle}(u)))}{\ln(F^{Usuelle}(x) \times F^{GPD}(x_{max}) \times (1 - F^{Usuelle}(u)))}$$

- Si  $x > x_{max}$ ,  $F^{hybrid}(x) = F^{GPD}(x)$

où les notations sont les mêmes que pour la partie précédente.  $F^{Usuelle}$  représente la fonction de répartition de la loi usuelle ajustée à la sinistralité attritionnelle (une lognormale par exemple) et  $F^{GPD}$  représente la fonction de répartition de la loi GPD qui ajuste la sinistralité extrême.  $u$  est le seuil d'entrée dans la zone extrême et  $x_{max}$  le plus gros sinistre observé.

Cette fonction hybride vérifie bien les propriétés des fonctions de répartition. En effet, il a été démontré qu'elle vérifie les cinq contraintes suivantes :

- $F^{hybrid}(u) = F^{Usuelle}(u)$
- $F^{hybrid}(x_{max}) = F^{GPD}(x_{max})$
- $F^{hybrid}$  est croissante entre 0 et 1
- $\forall x, u \leq x \leq x_{max}, F^{GPD}(x_{max}) \times F^{Usuelle}(x) \leq F^{hybrid}(x) \leq F^{Usuelle}(x)$
- $F^{hybrid}(x) \rightarrow F^{GPD}(x_{max}) \times F^{Usuelle}(x)$ , quand  $F^{Usuelle}(x) \rightarrow 0$

Par souci de confidentialité, les démonstrations ne figurent pas dans le rapport.

#### 4.4.2 Modèle de minimisation contraint

Dans les méthodes présentées ci-dessus, le seuil est déterminé de manière graphique. Il peut être intéressant de proposer une méthode alternative. Nous formulons ici notre proposition.

Il s'agit d'un problème de minimisation contraint. La distance entre la distribution empirique et la distribution globale ajustée est minimisée afin de déterminer le seuil  $u$ . L'élément clé est de choisir une bonne distance. Il peut être judicieux de choisir une norme qui dérive d'un produit scalaire, comme la norme 2. Considérons le produit scalaire :

$$\langle f, g \rangle_2 = \int fg$$

Le carré de la norme associée à ce produit scalaire est :

$$\|f\|_2^2 = \langle f, f \rangle_2 = \int f^2$$

Le problème de cette norme dans le cas qui nous intéresse est qu'elle va minimiser l'écart entre les fonctions de répartition théoriques et empiriques (respectivement notées  $F^{Globale}$  et  $F^{Empirique}$ , mais en considérant les points de  $[0, +\infty[$  de la même manière. Or  $F^{Empirique}$  est constante et égale à 1 après le plus gros sinistre. La fonction de répartition  $F^{Globale}$  qui va résulter du programme de maximisation aura donc tendance à être la plus proche possible de  $F^{Empirique}$  sur  $[X_{max}, +\infty[$ , qui vaut exactement 1.  $F^{Globale}$  sera donc très proche de 1 sur  $[X_{max}, +\infty[$ , c'est-à-dire que la queue sera exagérément fine par rapport à ce qu'elle devrait.

Il nous est donc apparu qu'une distance pondérée serait plus adaptée. On considère le produit scalaire :

$$\langle f, g \rangle_{2,\omega} = \int fg\omega$$

Et le carré de la norme associée :

$$\|f\|^2 = \langle f, f \rangle_{2,\omega} = \int f^2\omega$$

où  $\omega$  est une fonction poids telle que  $\omega(x) = \exp(-p.x)$  ou encore  $\omega(x) = 1/(x^2 + 1)$ .

Par exemple, en prenant  $\omega(x) = \exp(-x)$  :

$$\|f\|_{2,\omega}^2 = \langle f, f \rangle_{2,\omega} = \int_0^{+\infty} f^2(x) \exp(-x) dx$$

Ainsi, plus les  $x$  augmentent, moins grand est l'impact de la valeur de  $f(x)$ .

Pour la densité de  $F^{Globale}$ , on peut reprendre la densité du cas 2 :

$$f^{Globale} = f^{Usuelle} \cdot \mathbb{I}_{[0,u]} + (1 - F^{Usuelle}(u)) f^{GPD}$$

A présent,  $f^{Globale}$  (et donc  $F^{Globale}$ ) dépend de cinq paramètres : le seuil  $u$ , mais aussi les deux paramètres de la loi usuelle et les paramètres de queue  $\xi$  et d'échelle  $\sigma$  de la loi de Pareto généralisée. Sous cette forme,  $F^{Globale}$  est déjà continue, mais on peut imposer en plus que  $F^{Globale}$  soit  $C^1$ .

Le programme de minimisation contraint s'écrit alors :

$$\min_{u, a^{Usuelle}, b^{Usuelle}, \xi^{GPD}, \sigma^{GPD}} \|F^{Globale} - F^{Empirique}\|_{2, \omega}^2$$

$$sc : F^{Globale} \text{ est } C^1$$

En explicitant la forme de la norme et en détaillant la condition de régularité, sachant que le point où la densité  $f^{Globale}$  est éventuellement non continue est le seuil  $u$ , nous obtenons :

$$\min_{u, a^{Usuelle}, b^{Usuelle}, \xi^{GPD}, \sigma^{GPD}} \int_0^{+\infty} (F^{Globale}(x) - F^{Empirique}(x))^2 \exp(-x) dx$$

$$sc : F^{Usuelle'}(u) = F^{GPD'}(u)$$

La loi usuelle peut être, comme précédemment, soit une loi lognormale, dans ce cas les deux paramètres relatifs à la sinistralité attritionnelle sont  $a = \mu^{LN}$  et  $b = \sigma^{LN}$ , soit une loi de Weibull, dans ce cas les paramètres sont  $a = \alpha$  et  $b = \beta$ .

Ce problème de minimisation contraint ne semble pas pouvoir être résolu analytiquement. Nous sommes donc contraints d'utiliser une méthode numérique si nous voulons aboutir à un résultat pratique.

Néanmoins, il peut être intéressant de savoir s'il est possible de donner une autre forme à ce problème de minimisation sous certaines hypothèses. Faisons l'hypothèse forte que les quatre paramètres autres que le seuil sont fixés, soit par une première application du modèle 2, soit par des considérations externes sur les valeurs que doivent prendre les paramètres pour les sinistralités attritionnelle et extrême. Cherchons également une fonction de répartition globale qui est continue mais pas nécessairement  $C^1$ , ce qui permet de supprimer la contrainte. Le problème de minimisation non-contraint s'écrit alors :

$$\min_u \int_0^{+\infty} (F^{Globale}(x) - F^{Empirique}(x))^2 \exp(-x) dx$$

Le seul paramètre dont dépend la minimisation est alors uniquement le seuil. Comme il n'y a pas de contrainte, une condition nécessaire d'optimalité est alors que la dérivée partielle par rapport à  $u$  de l'expression à minimiser soit nulle, soit :

$$\frac{\partial}{\partial u} \int_0^{+\infty} (F^{Globale}(x) - F^{Empirique}(x))^2 \exp(-x) dx = 0$$

En notant  $H(u)$  le membre de gauche de cette égalité, nous avons donc une équation de la forme  $H(u) = 0$ . Une condition nécessaire d'optimalité du seuil  $u$  est donc  $H(u) = 0$ . Il semble impossible de trouver une solution analytique à une telle équation. Dans un cas pratique, la

détermination du seuil  $u$  doit se baser soit sur la résolution numérique de l'équation  $H(u) = 0$ , qui est une condition nécessaire d'optimalité, soit directement en résolvant numériquement le problème de minimisation associé.

## 5 Applications

Dans cette partie, nous allons mettre en application les méthodes décrites précédemment sur des données réelles. Les modèles à une loi (définition en 4.1) puis deux lois (définition en 4.2) seront utilisés et comparés. Le premier jeu de données concerne des sinistres de traités non proportionnels en responsabilité civile automobile et sera traité sous l'angle de la problématique du département Dynamic Financial Analysis (DFA). Le second jeu de données contient des sinistres de traités proportionnels en réassurance incendie et sera envisagé sous une problématique de tarification.

La partie pratique de ce mémoire a donné lieu au développement d'un outil Excel qui regroupe tous les outils nécessaires à la mise en œuvre des modèles 1 et 2 présentés dans la partie 4. Le fonctionnement de ce programme Excel est détaillé en annexe 1.

### 5.1 DFA

Le département DFA a pour objectif notamment d'étudier la solvabilité de la société à partir de l'analyse des flux financiers futurs.

La base de données dont nous disposons contient l'historique des coûts des sinistres responsabilité civile automobile à la charge de la SCOR, pour des traités de réassurance en excédent de sinistre. Pour un sinistre  $i$  de montant  $X_i$ , nous disposons donc du coût à la charge de la SCOR  $Y_i$  défini par :

$$Y_i = \begin{cases} 0 & , X < d_i \\ (X_i - d_i) \times tx_{SCOR,i} & , d_i < X < d_i + p_i \\ p_i \times tx_{SCOR,i} & , d_i + p_i < X_i \end{cases}$$

où  $d_i$  et  $p_i$  sont respectivement la priorité et la portée du traité en excédent de sinistre relatif au sinistre  $i$  et  $tx_{SCOR,i}$  est la part du sinistre réassurée par SCOR.

Prenons un exemple et considérons un sinistre d'un montant de 3 millions d'euros réassuré par un traité en (5 million) $XS$ (1 millions), à hauteur de 40% par SCOR, 30% par Swiss Re et 30% par Munich Re. Dans ce cas, le coût à la charge de la SCOR sera  $Y_i = (3 \text{ millions} - 1 \text{ million}) \times 40\% = 800 \text{ 000}\text{£}$ .

On se trouve alors face à un premier problème pour modéliser les coûts des sinistres à la charge de SCOR. On ne peut théoriquement pas les modéliser tous de la même manière, puisque la franchise et la part réassurée par la SCOR seront différentes selon les traités. On suppose alors que les conditions sont plus ou moins uniformes et qu'il existe une certaine homogénéité dans la structure du portefeuille de la SCOR. Ainsi, en supposant que la structure est la même dans le temps, nous pourrions prévoir les flux futurs à l'aide des flux historiques.

La base de données dont nous disposons comporte deux triangles de liquidation recensant les montants des sinistres à la charge de la SCOR de 1996 à 2005 pour la branche responsabilité civile automobile en Grande-Bretagne. Dans cette branche, les sinistres peuvent avoir un

développement long, plusieurs années peuvent se passer entre la survenance et le règlement complet du sinistre.

Le premier triangle fournit les montants "*Incurred*", c'est à dire l'évaluation totale de chaque sinistre. Ces montants évoluent au cours du temps, car l'évaluation initiale du sinistre n'est pratiquement jamais égale à ce qui sera finalement payé au total. En effet, pour les sinistres à développement long, la première évaluation du sinistre est difficile puisque le montant dû à la victime lui sera payé en plusieurs fois, sur une longue période et pourra être modifié par l'intervention des instances judiciaires. Le second triangle fournit les montants payés cumulés (*SP*). De ces deux triangles, nous pouvons en déduire le triangle des montants restant à payer (*SAP*), puisque nous avons la relation :  $Incurred = SP + SAP$ , ainsi que le triangle des montants payés décumulés.

Nous disposons en plus des indices sinistres de 1996 à 2020 allant servir à la revalorisation.

### 5.1.1 Notations

Dans la suite, nous utiliserons les notations suivantes :

- $I_i$  : l'indice sinistre relatif à l'année  $i$
- $S_{i,j}$  : le montant payé cumulé à la fin de l'année  $j$  pour un sinistre survenu l'année  $i$
- $Z_{i,j}$  : le montant payé l'année  $j$  pour un sinistre survenu l'année  $i$
- $SAP_{i,j}$  : le montant restant à payer l'année  $j$  pour un sinistre survenu l'année  $i$
- $Inc_{i,j}$  : le montant Incurred l'année  $j$  pour un sinistre survenu l'année  $i$

Prenons l'exemple d'un sinistre survenu en 2000. Le tableau suivant montre son évolution annuelle. On remarque que l'évaluation faite du coût total du sinistre à la fin de l'année de survenance (108 960£) était insuffisante puisque la compagnie aura finalement payé 128 236£ pour ce sinistre.

Année de développement $j$	0	1	2	3	4	5
Paiement Annuel ( $Z_{i,j}$ )	0	0	101 540	20 900	5 095	701
Paiement Cumulé ( $S_{i,j}$ )	0	0	101 540	122 440	127 535	128 236
Restant ( $SAP_{i,j}$ )	108 960	108 960	7 420	5 095	0	0
Coût Total ( $Inc_{i,j}$ )	108 960	108 960	108 960	127 535	127 535	128 236

### 5.1.2 Revalorisation des sinistres

Pour prévoir la sinistralité future en Responsabilité Civile Automobile (RCA), nous allons nous appuyer sur l'historique du portefeuille RCA. Ainsi, chaque année passée est considérée comme un scénario possible pour l'année d'étude. Ceci suppose que toutes les conditions d'une année passée soient identiques aux conditions de l'année étudiée. En particulier, la revalorisation des sinistres est une étape indispensable dans le traitement des données.

Un sinistre passé a un certain coût. Le but de la revalorisation des sinistres est de déterminer le coût de ce même sinistre s'il était survenu l'année étudiée, on parle également de "mise as if" des sinistres. Plusieurs méthodes sont envisageables pour revaloriser les sinistres : des méthodes prospectives, des méthodes rétrospectives ou encore un mélange entre les deux. Pour les sinistres à développement long, nous avons privilégié une approche prospective reposant sur l'estimation des valeurs futures des indices sinistres qui vont servir à

la revalorisation de nos données. Ainsi, l'inflation est représentée par un indice sinistre ( $I$ ) connu jusqu'en 2006 et dont l'évolution future a été estimée jusqu'en 2020.

La mise as if s'effectue sur les montants payés non cumulés et sur les montants restant à payer. Supposons que l'on veuille revaloriser les sinistres pour l'année de survenance  $N$ . Il faut alors multiplier chaque montant par le facteur  $FI_N = I_{AC+j}/I_{N+j}$  où  $j$  est l'année de développement du sinistre et  $AC$  l'année de cotation.

Année de Survenance	Année de développement			
	0	1	2	3
2003	a	b	c	d
2004	e	f	g	
2005	h	i		
2006	j			



	Année de développement			
	0	1	2	3
Scénarii pour 2007	$a \frac{I_{07}}{I_{03}}$	$b \frac{I_{08}}{I_{04}}$	$c \frac{I_{09}}{I_{05}}$	$d \frac{I_{10}}{I_{06}}$
	$e \frac{I_{07}}{I_{04}}$	$f \frac{I_{08}}{I_{05}}$	$g \frac{I_{09}}{I_{06}}$	
	$h \frac{I_{07}}{I_{05}}$	$i \frac{I_{08}}{I_{06}}$		
	$j \frac{I_{07}}{I_{06}}$			

FIG. 20 – Revalorisation des sinistres.

Appliquons la mise as if à l'exemple précédent du sinistre survenu en 2000. On obtient les montants revalorisés à l'année 2007 suivants, sur la base de l'indice SCOR utilisé en RCA en Grande-Bretagne :

Année de développement $j$	0	1	2	3	4	5
Païement Annuel ( $Z_{i,j}$ )	0	0	213 422	44 780	10 816	1480
Païement Cumulé ( $S_{i,j}$ )	0	0	213 422	258 202	269 018	270 498
Restant ( $SAP_{i,j}$ )	225 313	225 871	15 595	10 916	0	0
Coût Total ( $Inc_{i,j}$ )	225 313	225 871	229 017	269 118	269 018	270 498

Cette méthode de revalorisation permet d'estimer les valeurs qu'aurait pris chaque sinistre s'il était survenu en 2007. Les indices sinistres nous ont été fournis par la SCOR, mais chaque réassureur détermine l'indice à utiliser selon la branche considérée. En RCA, il est fréquent d'utiliser un indice corporel ou bien un cumul de deux indices : un premier représentant l'inflation économique, comme le coût du travail horaire, et un second représentant l'inflation juridique (Super Imposed Inflation).

### 5.1.3 Projection à l'ultime de la charge des sinistres

La deuxième étape va consister à estimer le coût ultime des sinistres. On procède pour cela à ce que l'on appelle la "liquidation" des triangles. La méthode la plus utilisée est la méthode chain ladder. Cette méthode s'appuie sur le modèle de développement pour les facteurs. Le modèle de développement pour les facteurs suppose l'existence de paramètres  $\varphi_1, \dots, \varphi_n \in (1, \infty)$  tels que

$$E[S_{i,j}] = E[S_{i,j-1}] \times \varphi_j$$

pour tout  $i \in \{0, 1, \dots, n\}$  et  $j \in \{1, \dots, n\}$ . Les  $\varphi_j$  sont les facteurs de développement. Ils sont inconnus et doivent être estimés. Ainsi, pour chaque année de développement  $j \in \{1, \dots, n\}$ , on définit le facteur de Chain Ladder  $\hat{F}_j^{CL}$  comme estimateur du paramètre  $\varphi_j$  :

$$\hat{F}_j^{CL} = \frac{\sum_{k=0}^{n-j} S_{k,j}}{\sum_{k=0}^{n-j} S_{k,j-1}}$$

Et pour chaque année d'origine  $i$ , on en déduit l'estimateur de Chain Ladder de l'espérance  $E[S_{i,j}]$  :

$$\hat{S}_{i,j}^{CL} = S_{i,n-i} \prod_{l=n-i+1}^k \hat{F}_l^{CL}$$

Ainsi, on calcule des coefficients de passage d'une année sur l'autre, ce qui nous permet de compléter le triangle. Cette méthode se base donc sur l'hypothèse que la cadence historique de l'évolution des sinistres se répète. Elle a l'avantage d'être simple et d'être fondée sur l'observation plutôt que sur la connaissance a priori. Elle présente toutefois des défauts non négligeables. Tout d'abord, cette méthode n'est pas additive, ce qui n'est pas gênant dans notre cas puisque nous l'appliquerons au triangle des "incurred" uniquement. De plus, pour les années récentes, l'incertitude peut être très importante puisque le coefficient multiplicatif est alors le produit de  $n-1$  estimations de coefficients de proportionnalité. Enfin, cette méthode ne permet pas d'obtenir de mesure de précision sur les estimations.

Conscients de ces inconvénients, nous avons appliqué la méthode Chain Ladder aux montants incurred afin d'obtenir le coût ultime des sinistres à la charge de SCOR.

### 5.1.4 Etude de la série des ultimes

#### Statistiques descriptives

Afin de nous familiariser avec les données, nous avons commencé par effectuer quelques statistiques descriptives sur les sinistres RCA.

Nombre d'observations	1851
Sinistre minimum	18£
Sinistre moyen	157 472£
Sinistre maximum	2 044 115£
Ecart-Type	185 646£
Sinistre médian	101 222£
Quantile à 90%	363 436£
Quantile à 95%	481 491£
Quantile à 99%	858 008£

On constate une grande étendue des valeurs. En effet, le minimum est de 18£, alors que le sinistre maximum dépasse les 2 millions de £, soit près de 13 fois la moyenne. L'écart-type est important. On retrouve la particularité des données de réassurance : beaucoup de petits sinistres et quelques sinistres très importants. En effet, la médiane est nettement inférieure à la moyenne qui est augmentée par les gros sinistres. On remarque également que le quantile à 90% est faible par rapport au sinistre maximum, puisqu'il ne représente que 18% de ce dernier.

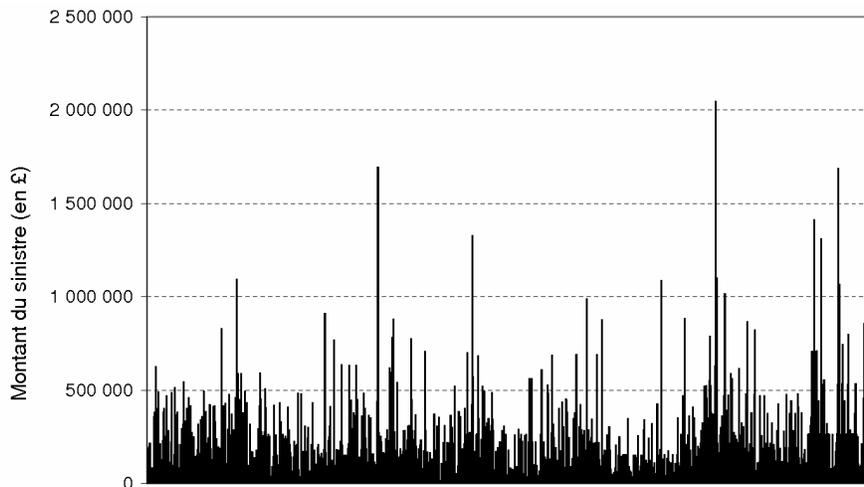


FIG. 21 – Sinistres RCA.

Le graphe de la série des sinistres (figure 21) nous permet d'identifier les plus gros montants et leur nombre d'occurrence. Il nous permet également de voir s'il y a un regroupement des gros sinistres ou une tendance dans leur fréquence d'apparition, ce qui remettrait en cause notre hypothèse d'indépendance. Cela ne semble pas être le cas pour nos données RCA.

## QQ Plot

Afin d'avoir une première idée de la distribution du montant des sinistres, nous avons tracé les graphiques Quantiles Quantiles étudiés dans la partie 3, pour quatre distributions usuelles : Exponentielle, Lognormale, Pareto et Weibull.

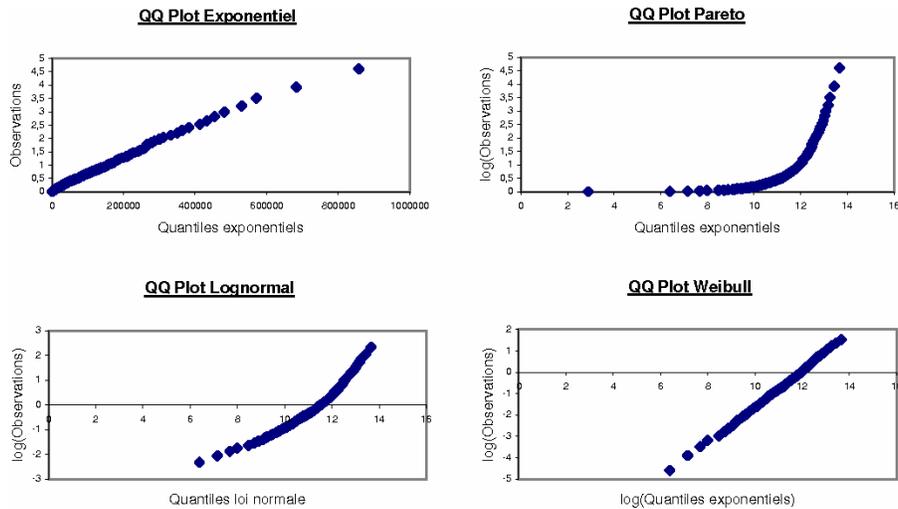


FIG. 22 – QQ Plot de lois usuelles pour les sinistres RCA.

Le graphique obtenu pour le QQ Plot Weibull ressemble à une droite. Il est possible qu'il y ait une relation affine entre les quantiles de l'échantillon et ceux de la loi de Weibull et l'adéquation de la loi de Weibull à nos données est donc probable. Cependant, l'adéquation à une loi exponentielle paraît également possible à la vue du graphique. En revanche, l'étude des QQ Plot élimine les lois lognormale et de Pareto comme candidats à l'adéquation à notre échantillon.

## Mean excess function

Nous avons ensuite tracé le graphique de la mean excess function.

Nous avons volontairement omis les quatre derniers points du graphique. Ces points, étant des moyennes d'au mieux quatre observations, auraient pu être très irréguliers, et, comme nous l'avons vu sur les simulations de la partie 3, la mean excess function souffre d'une instabilité sur les dernières observations.

Théoriquement, une fonction de moyenne en excès constante correspond à une distribution exponentielle, une fonction affine correspond à une Pareto. En dehors de ces deux cas particuliers, nous pouvons dire que si les points du graphique de la mean excess function présentent une tendance à la baisse, c'est le signe d'une distribution à queue fine, et inversement, les données issues de distributions à queue épaisse présentent une tendance à la hausse.

Le graphique obtenu (figure 23) correspond plutôt à ce dernier cas. Ceci nous amène à penser que nous sommes en présence d'une distribution du montant des sinistres à queue épaisse, ce qui est fréquent avec les données de réassurance. Les dernières observations sont étalées et difficilement interprétables, ce qui est courant avec les observations extrêmes d'un échantillon

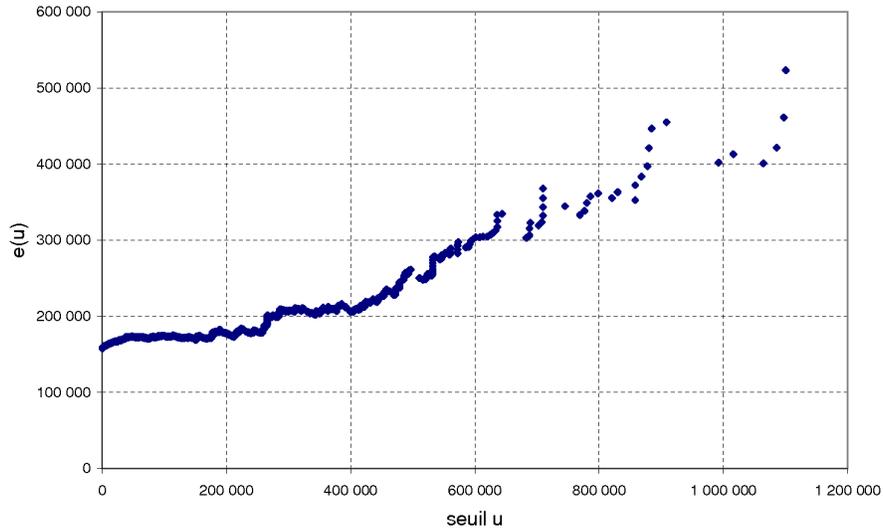


FIG. 23 – Fonction de moyenne en excès pour les sinistres RCA.

réel. Toutefois, la détermination du seuil à partir duquel les points sont alignés permet de connaître le seuil à partir duquel la modélisation par une loi de Pareto généralisée est justifiée.

### 5.1.5 Modélisation à l'aide d'une seule distribution

La première étape de notre modélisation consistait à étudier si une seule loi permettait d'ajuster correctement la sinistralité, aussi bien sur la partie centrale de la distribution qu'en queue de distribution. Nous avons pour cela appliqué les tests d'adéquation de Kolmogorov-Smirnov et Anderson Darling à plusieurs lois usuelles. Le seuil de significativité des tests est toujours 5%. Nous avons obtenu les résultats suivants :

Loi usuelle	Test de Kolmogorov-Smirnov	Test d'Anderson-Darling
Lognormale	Rejet	Rejet
Gamma	Rejet	Rejet
Weibull	Acceptation	Acceptation
Pareto	Rejet	Rejet
Pareto généralisée	Rejet	Rejet

Ainsi, les tests d'adéquation confirment l'impression visuelle donnée par les QQ Plots, puisque l'adéquation à la loi de Weibull de paramètres  $\alpha = 0,87$  et  $\beta = 147025$  est acceptée par ces tests. De plus les lois de Weibull avec un paramètre  $\alpha < 1$  entrent dans la classe des lois à queues épaisses puisque leur fonction génératrice des moments n'est pas finie :  $M_X(t) = E(\exp(Xt)) = +\infty$  pour  $t > 0$ , ce qui confirme l'intuition laissée par le graphe de la mean excess function.

On constate (figure 24) que l'ajustement à l'aide d'une seule distribution de Weibull est bonne. Cependant, comme nous l'avons montré dans la partie 4, il est rare qu'une seule distribution permette de bien reproduire à la fois la sinistralité attritionnelle et la sinistralité

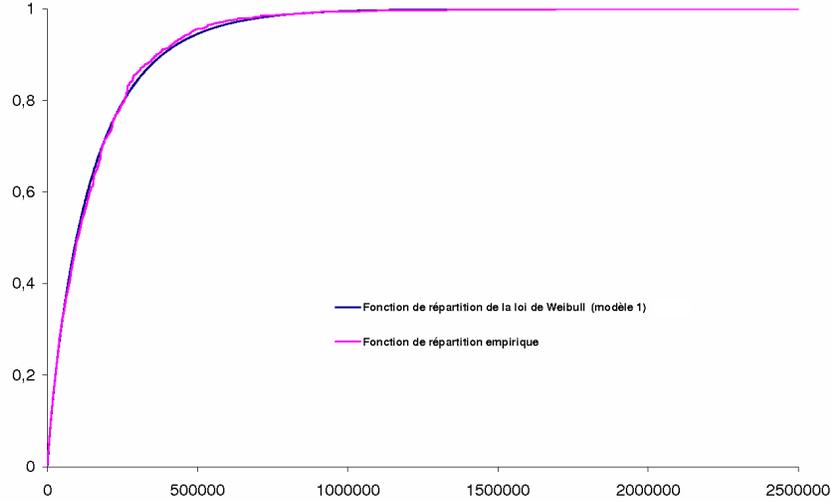


FIG. 24 – Fonctions de répartition empirique et de la loi de Weibull.

extrême. De plus, l'adéquation en queue de distribution ne se voit pas forcément bien graphiquement à cause du faible nombre d'observations et il n'existe pas de tests statistiques efficaces permettant de la vérifier. C'est pourquoi nous allons cependant chercher à appliquer la modélisation à l'aide de deux distributions, afin de pouvoir comparer les résultats obtenus.

### 5.1.6 Modélisation à deux lois

Dans le but de mieux modéliser la queue de distribution, nous allons mettre en oeuvre le modèle 2 défini dans la partie 4 sur ces mêmes données.

#### Ajustement GEV

Tout d'abord, afin d'acquérir certaines informations qualitatives relatives à la queue de distribution, nous avons tracé des graphiques qui permettent de distinguer les cas  $\xi = 0$  et  $\xi > 0$ . Le premier graphique permet d'avoir une représentation visuelle de la pertinence de l'hypothèse  $\xi = 0$ . Si cette hypothèse est vérifiée, alors les points du graphique sont alignés. Le second graphique est le pendant du premier pour l'hypothèse  $\xi > 0$ . Si cette hypothèse est vérifiée, alors les points du graphique sont alignés.

Les deux figures sont présentées ci-après.

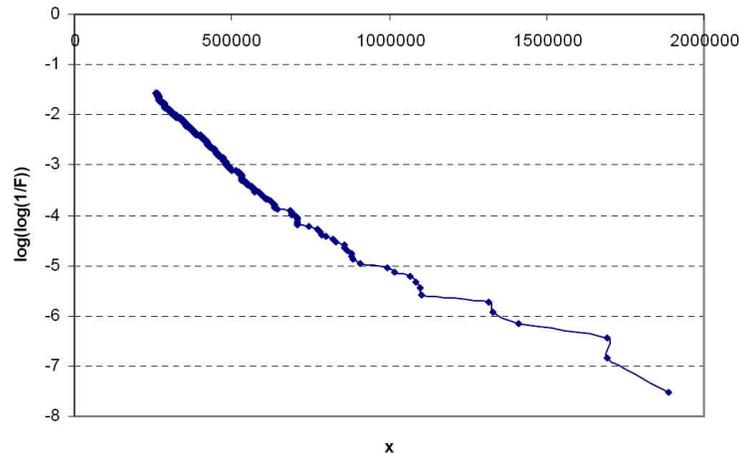


FIG. 25 – Test de  $\xi = 0$ .

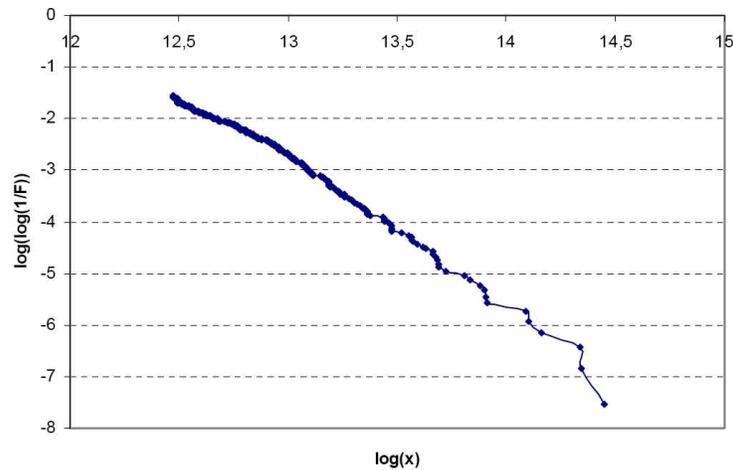


FIG. 26 – Test de  $\xi > 0$ .

Ces graphiques ne permettent pas de différencier nettement le cas  $\xi = 0$  du cas  $\xi > 0$ . Cependant, la première courbe est légèrement convexe, alors que la seconde est plus aisément indentifiable à une droite. L'hypothèse  $\xi > 0$  semble donc légitime.

### Détermination du seuil

Nous pouvons à présent mettre en oeuvre l'estimateur de Hill, dont l'utilisation est licite puisque  $\xi > 0$ , afin de déterminer graphiquement le seuil départageant les sinistres attritionnels des sinistres extrêmes.

Nous avons tracé le Hill plot (figure 27) en fonction du seuil  $u$  et du nombre d'excès  $k$ , en mettant en évidence l'intervalle de confiance à 95% autour de cet estimateur.

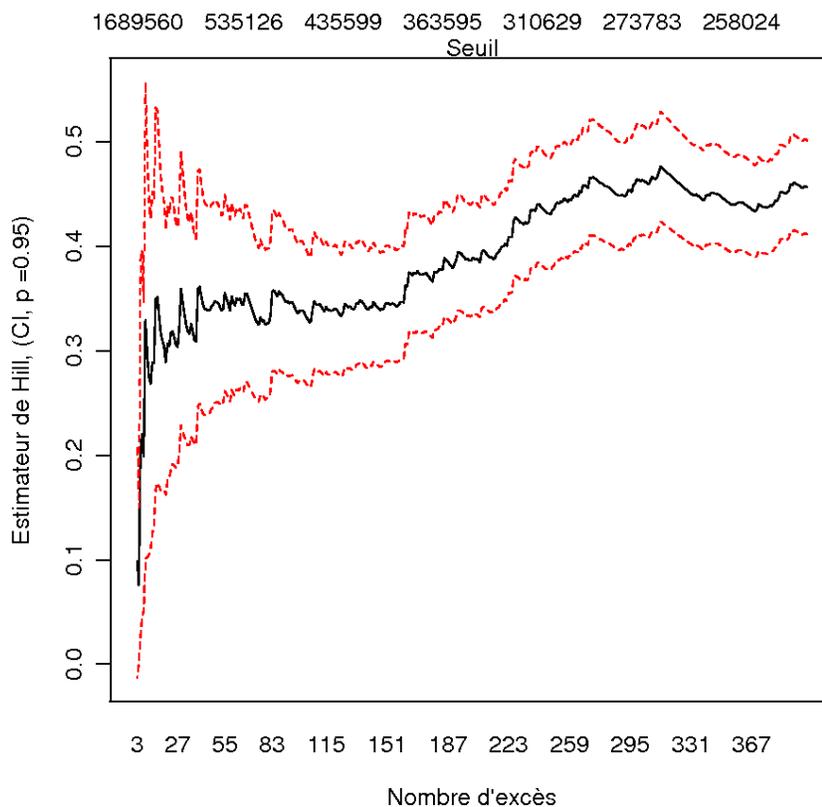


FIG. 27 – Estimateur de Hill en fonction du nombre d'excès considérés et du seuil  $u$ .

Nous voyons sur la figure que la zone de stabilité semble se situer autour de 160 excès, ce qui correspond à un seuil de 400 000 £. Un sinistre inférieur à 400 000 £ sera donc considéré comme attritionnel et un sinistre supérieur comme extrême.

Ce seuil de 400 000 £ semble cohérent avec la fonction de moyenne en excès. En effet, sur la représentation graphique de cette dernière (figure 23), on constate visuellement un léger décrochage à 400 000, après quoi les points semblent être identifiables à une fonction affine.

Maintenant que l'on connaît le seuil, on peut modéliser la sinistralité attritionnelle, à l'aide d'une loi usuelle tronquée, ainsi que la sinistralité extrême, à l'aide d'une loi de Pareto généralisée.

### Ajustement de la partie centrale

La partie centrale, qui représente la sinistralité attritionnelle, est modélisée par une loi usuelle tronquée. Comme dans le cas de la modélisation par une seule loi, les tests d'adéquation de Kolmogorov-Smirnov et d'Anderson-Darling nous conduisent à ne retenir que la loi de Weibull.

L'estimation des deux paramètres de la loi de Weibull par la méthode du maximum de vraisemblance conduit à  $\alpha = 0,87$  et  $\beta = 147025$ .

## Ajustement de la queue de distribution

La queue de la distribution, qui représente la sinistralité extrême, est modélisée par une loi de Pareto généralisée. L'estimation des deux paramètres de la loi de Pareto généralisée par la méthode du maximum de vraisemblance conduit à  $\xi = 0,31$  et  $\sigma = 144494$ . La valeur de  $\xi$  trouvée par maximum de vraisemblance est cohérente avec l'estimateur de Hill. En effet, pour un seuil  $u = 400\,000\mathcal{L}$ , l'estimateur de Hill correspondant est  $\hat{\xi}^{Hill} = 0,34$  avec un intervalle de confiance à 95% de  $[0,29; 0,40]$ .

Les fonctions de répartition empirique et modélisée relatives à la sinistralité extrême sont représentées sur la figure suivante.

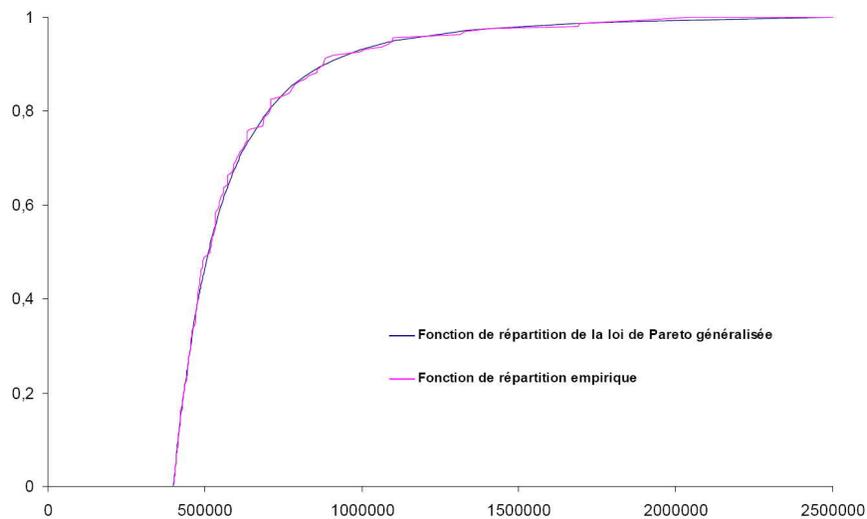


FIG. 28 – Sinistralités extrêmes empirique et modélisée.

L'adéquation est visuellement très bonne. La distribution de Pareto généralisée modélise très bien la sinistralité extrême. Le seuil choisi semble donc pertinent.

## Modélisation globale

Maintenant que nous avons modélisé les sinistralités attritionnelle et extrême, ainsi qu'estimé tous les paramètres, nous pouvons modéliser la sinistralité dans son ensemble en utilisant le modèle à deux lois défini dans la partie 4. Les fonctions de répartition globales empirique et modélisée sont représentées sur la figure 29 suivante.

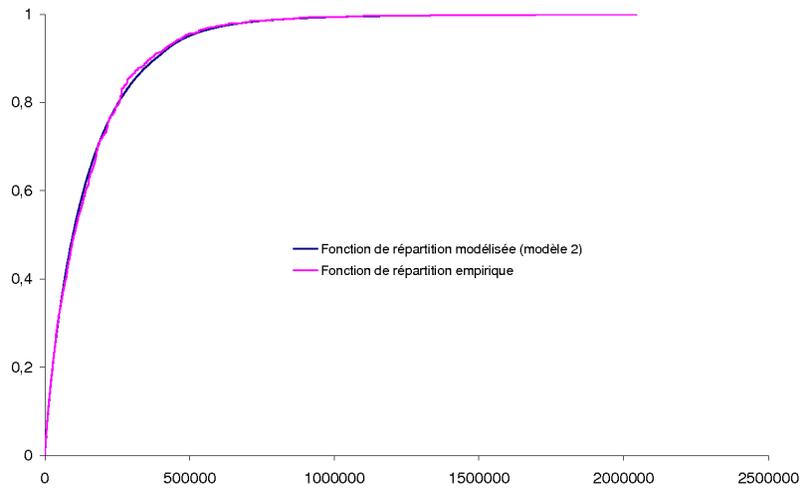


FIG. 29 – Sinistralités globales empirique et modélisée (modèle à deux lois).

Il est également intéressant de faire un zoom sur la queue de distribution et de comparer la modélisation avec deux lois (loi de Weibull puis loi de Pareto généralisée) à la modélisation avec une seule loi de Weibull.

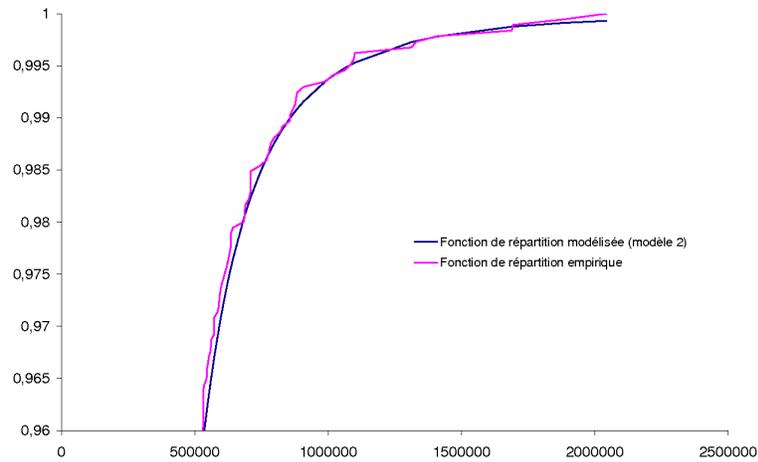


FIG. 30 – Sinistralités globales empirique et modélisée (modèle à deux lois), détails de la queue de distribution.

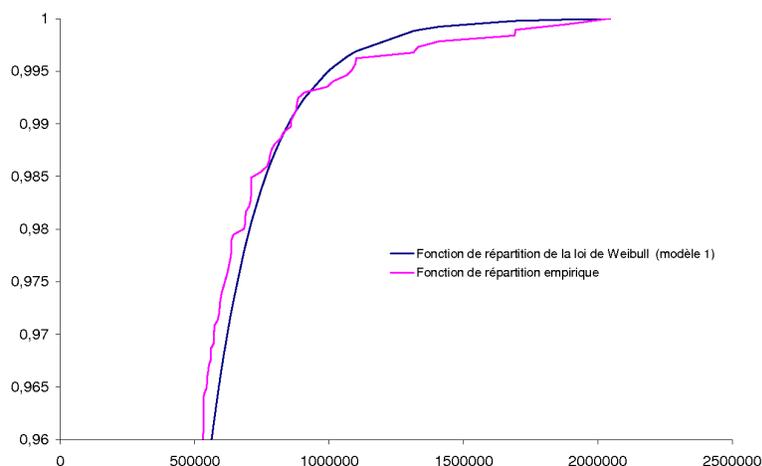


FIG. 31 – Sinistralités globales empirique et modélisée (modèle à une loi de Weibull), détails de la queue de distribution.

On constate que le modèle à deux lois est bien meilleur : d'une part il se rapproche plus de la répartition empirique, d'autre part la fin de la queue est plus épaisse (fonction de répartition plus faible), ce qui diminue le risque de sous-estimer les pertes en cas de sinistres extrêmes.

### 5.1.7 Estimation de quantiles extrêmes

Dans la partie précédente, nous avons proposé deux modélisations possibles pour le montant des sinistres à la charge de SCOR. La première modélisation ajuste une seule loi de Weibull à l'ensemble des données, tandis que la seconde modélisation se veut plus précise et ajuste une loi de Weibull sur la sinistralité attritionnelle et une loi GPD sur la sinistralité extrême. Nous avons vu graphiquement que la deuxième modélisation semblait meilleure. Nous allons à présent essayer de quantifier ces différences. Pour cela, nous allons étudier l'impact de ces deux modélisations sur le niveau des quantiles extrêmes. Nous allons en particulier nous intéresser au quantile à 99,5%, sur lequel est fondé le calcul du capital requis dans la réglementation prudentielle.

Pour rappel, le quantile d'ordre  $p \in ]0, 1[$  associé à la variable aléatoire  $X$  de fonction de répartition  $F_X$  se définit par :

$$x_p = F^{-1}(p) = \sup\{x \in \mathbb{R}, F(x) \geq p\}$$

On parle aussi de Value-at-risk de niveau  $p$ . Nous allons comparer les quantiles obtenus en les estimant de manière empirique, puis à l'aide du modèle à une loi (modèle 1) et enfin à l'aide du modèle à deux lois (modèle 2).

## Estimation empirique

L'estimation empirique du quantile à 99,5% nous conduit à  $x_{99,5}^{emp} = 1\,080\,559\mathcal{L}$ .

## Estimation selon le modèle 1

Dans le modèle à une loi, nous avons ajusté une loi de Weibull à notre échantillon. Les paramètres ont été estimés par maximum de vraisemblance :  $\hat{\alpha} = 0,87$  et  $\hat{\beta} = 147025$ . Nous estimons le quantile à 99,5% selon le modèle 1 ( $\hat{x}_{99,5}^{(1)}$ ) par le quantile à 99,5% d'une loi de Weibull de paramètres  $\hat{\alpha}$  et  $\hat{\beta}$ . Nous obtenons  $\hat{x}_{99,5}^{(1)} = 997\,190\mathcal{L}$ .

## Estimation selon le modèle 2

Dans le modèle à deux lois, nous avons modélisé la sinistralité attritionnelle (la partie centrale de la distribution) par une loi de Weibull et la sinistralité extrême (la queue de distribution) par une loi GPD. Nous allons utiliser cette modélisation globale pour estimer le quantile à 99,5% selon le modèle 2 ( $\hat{x}_{99,5}^{(2)}$ ).

Nous rappelons que nous notons  $u$  le seuil au delà duquel nous ajustons une loi GPD, (nous avons trouvé  $u = 400\,000\mathcal{L}$ ). Le modèle 2 nous a conduit à ajuster à nos données une fonction de répartition globale de la forme :

$$F^{Globale}(x) = \begin{cases} F^W(x) & \text{si } x < u \\ F^W(u) + (1 - F^W(u))F^{GPD}(x) & \text{si } x \geq u \end{cases}$$

où  $F^W$  est la fonction de répartition de la loi de Weibull qui modélise la sinistralité attritionnelle.

Cela signifie en particulier que pour  $x \geq u$ , nous pouvons utiliser  $\hat{F}(x) = F^W(u) + (1 - F^W(u))F^{GPD}(x)$  comme approximation de la queue de distribution de  $X$ . Or, il est possible de montrer (cf McNeil [19]) que  $\hat{F}(x)$  est également une distribution de pareto généralisée, avec le même paramètre de queue  $\xi$ , mais avec un paramètre d'échelle  $\tilde{\sigma} = \sigma(1 - F^W(u))^\xi$  et un paramètre  $\tilde{\mu} = \mu - \tilde{\sigma}((1 - F^W(u))^{-\xi} - 1)/\xi$ .

De plus, pour une loi  $GPD(\xi, \sigma, \mu)$ , le quantile d'ordre  $p$  est :

$$x_p = \mu + \frac{\sigma}{\xi}((1-p)^{-\xi} - 1)$$

Finalement, pour le modèle à deux lois, nous obtenons  $\hat{x}_{99,5}^{(2)} = 1\,078\,538\mathcal{L}$

## Comparaison des modèles

Le tableau suivant récapitule les valeurs obtenues pour divers quantiles extrêmes selon chaque modèle.

ordre $p$ du quantile	$\hat{x}_p^{emp}$	$\hat{x}_p^{(1)}$	$\hat{x}_p^{(2)}$
95%	481 491£	518 108£	495 977£
99%	858 008£	848 793£	858 183£
99,5%	1 080 559£	997 190£	1 078 538£
99,9%	1 722 704£	1 351 933£	1 813 346£

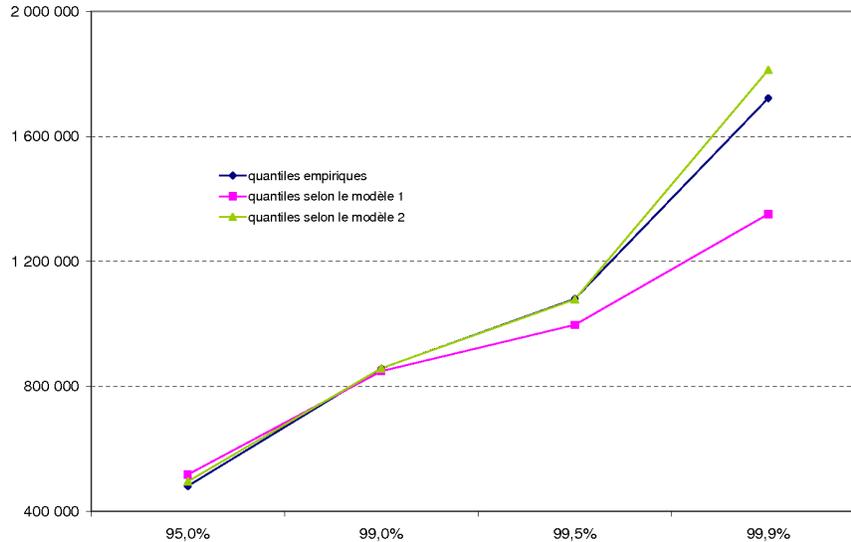


FIG. 32 – Comparaison des quantiles extrêmes pour les différents modèles.

A la vue de ces chiffres, nous pouvons constater l'intérêt de la modélisation à deux lois. L'ajustement à une seule loi a tendance à sous-estimer le risque, car il ne permet pas de reproduire une queue de distribution suffisamment épaisse. A l'inverse, la modélisation à deux lois permet de calculer des quantiles très proches des valeurs observées. Elle a même tendance à sur-estimer le risque pour le quantile à 99,9% ce qui lui confère un caractère prudentiel. L'objectif n'est pas, bien entendu, de sur-estimer le risque de manière démesurée, sous peine de devoir le tarifier trop cher et de voir partir les clients à la concurrence. Cependant l'aspect prudentiel d'une modélisation est toujours intéressant, surtout étant donné le nouveau cadre réglementaire en matière de gestion des risques.

Ainsi, le principal intérêt de cette modélisation pour la DFA réside dans le calcul des quantiles extrêmes. Cependant, il est également important d'avoir une bonne estimation de toute la distribution des montants de sinistres, afin d'avoir une meilleure précision en moyenne, notamment dans une optique de business plan.

## 5.2 Pricing

Nous passons à présent à une application tarification de la méthode présentée dans ce mémoire. Dans l'approche tarification, ce ne sont plus les quantiles de la distribution des sinistres qui vont nous intéresser, mais la valeur de la prime pure de traités en excédent de sinistre. Ceci peut être utile pour améliorer la tarification de la tranche de capacité non travaillante en cas de traité proportionnel.

La base de données dont nous disposons contient l'historique des coûts des sinistres incendie en France, pour des traités de réassurance proportionnelle.

Nous avons vu que la tarification des traités proportionnels s'effectuait en deux parties : la tarification de la sinistralité attritionnelle dans un premier temps puis la tarification de la sinistralité extrême qui s'apparente à une tarification XS. Notre méthode vise surtout à améliorer la tarification en queue de distribution, afin d'évaluer au mieux les gros sinistres et d'éviter toute sur- ou sous-tarification.

### 5.2.1 Revalorisation des sinistres

Comme pour l'application DFA précédente, nous allons nous appuyer sur l'historique du portefeuille Incendie France pour tarifier les traités proportionnels Incendie France. Ainsi, chaque année passée est considérée comme un scénario possible pour l'année d'étude et il est nécessaire que toutes les conditions d'une année passée soient identiques aux conditions de l'année étudiée. Nous commençons donc par effectuer une mise as if rétrospective des montants de sinistre de la base afin de faire disparaître de nos données tout effet dû à l'inflation.

### 5.2.2 Obtention des ultimes

Contrairement à l'application DFA qui concernait une branche à développement long, nous sommes ici en présence de sinistres de la branche incendie qui est une branche à développement court. En effet, nous avons pu observer sur nos données qu'au bout de trois ans au maximum, les sinistres étaient clos (le montant des sinistres à payer est nul et le montant incurred est égal au montant payé cumulé). Nous avons décidé de supprimer les 3 dernières années d'historique sinistre et nous avons donc conservé les sinistres survenus entre 1983 et 2003. Ceci nous permet de garder les vrais valeurs ultimes des sinistres au lieu de les estimer par la méthode chain ladder, ce qui entraînerait un risque de perdre l'indépendance entre nos données, en plus de l'erreur d'estimation. On peut ainsi considérer que nous sommes en présence de données indépendantes et identiquement distribuées (notamment grâce à la revalorisation).

### 5.2.3 Etude de la série des ultimes

#### Statistiques descriptives

Pour traiter l'application pricing, nous avons suivi la même démarche que pour l'application DFA. Ainsi, nous avons commencé par effectuer quelques statistiques descriptives sur les sinistres Incendie France.

Nombre d'observations	378
Sinistre minimum	63€
Sinistre moyen	547 495€
Sinistre maximum	16 596 085€
Ecart-Type	1 149 631€
Sinistre médian	160 944€
Quantile à 90%	1 408 456€
Quantile à 95%	2 489 314€
Quantile à 99%	3 846 749€

La base de données dont nous disposons contient 378 observations. On constate la présence d'un très gros sinistre de plus de 16 millions d'euros. Hormis celui-ci, les sinistres sont tous inférieurs à 5 millions d'euros, le sinistre médian n'est que de 160 944€ et le quantile à 90% s'élève à 1 408 456€.

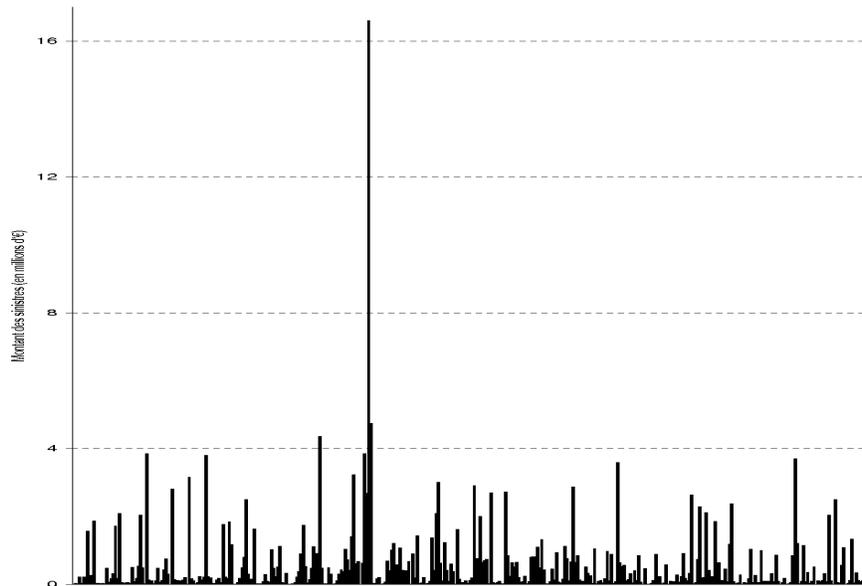


FIG. 33 – Sinistres incendie France entre 1983 et 2003.

La figure 33 permet d'identifier les sinistres importants. Il ne semble pas y avoir ici de tendance particulière dans l'occurrence de ces sinistres. Notre hypothèse d'indépendance n'est donc pas remise en cause.

### QQ Plot

Afin d'avoir une première idée de la distribution du montant des sinistres, nous avons tracé les graphiques Quantiles Quantiles étudiés dans la partie 3, pour quatre distributions usuelles : Exponentielle, Lognormale, Pareto et Weibull.

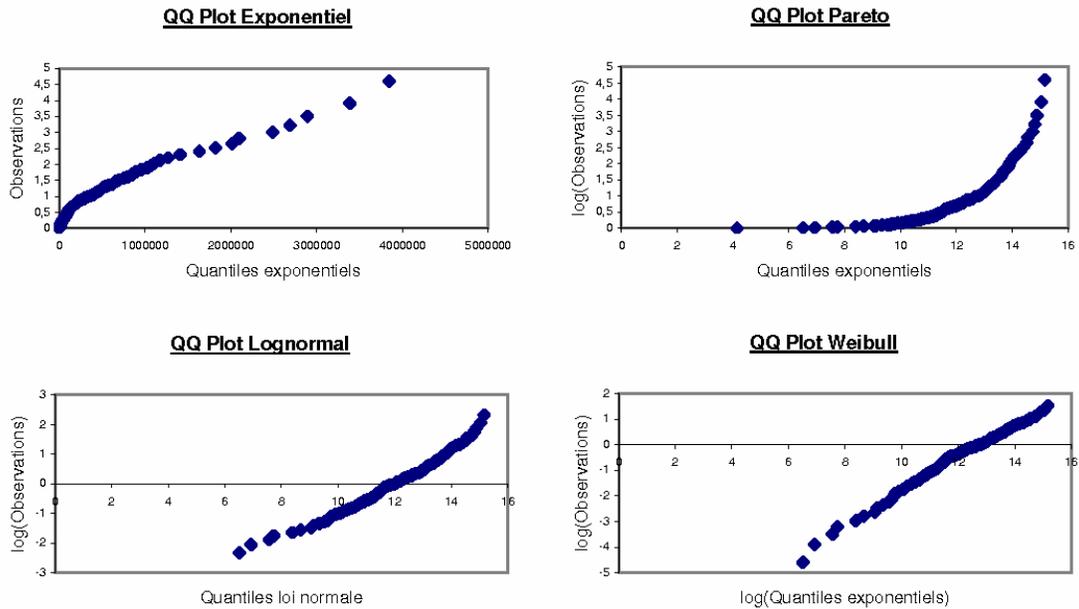


FIG. 34 – QQ Plot de lois usuelles pour les sinistres Incendie France.

Le graphe obtenu pour le QQ Plot Weibull ressemble à une droite. Il est possible qu'il y ait une relation affine entre les quantiles de l'échantillon et ceux de la loi de Weibull ; l'adéquation de la loi de Weibull à nos données est donc probable. L'adéquation à une loi lognormale paraît également possible au vu du graphique, bien qu'un peu moins bonne. En revanche, l'étude des QQ Plot élimine les lois exponentielle et de Pareto comme adéquation à notre échantillon.

### Mean excess function

Nous avons ensuite tracé le graphe de la mean excess function (figure 35).

Il est difficile de tirer une interprétation de ce graphique. La mean excess function est un outil qui fonctionne assez bien sur les données simulées. En revanche, les résultats obtenus sur des données réelles sont parfois difficilement interprétables. On peut cependant observer un léger décrochage de la fonction moyenne en excès aux alentours de 1 million d'euros. La fonction prend alors la forme d'une droite avant de se stabiliser et de repartir à la hausse aux alentours de 2.5 millions d'euros. La fonction n'étant pas horizontale, on peut postuler une queue épaisse pour la distribution du montant des sinistres incendie.

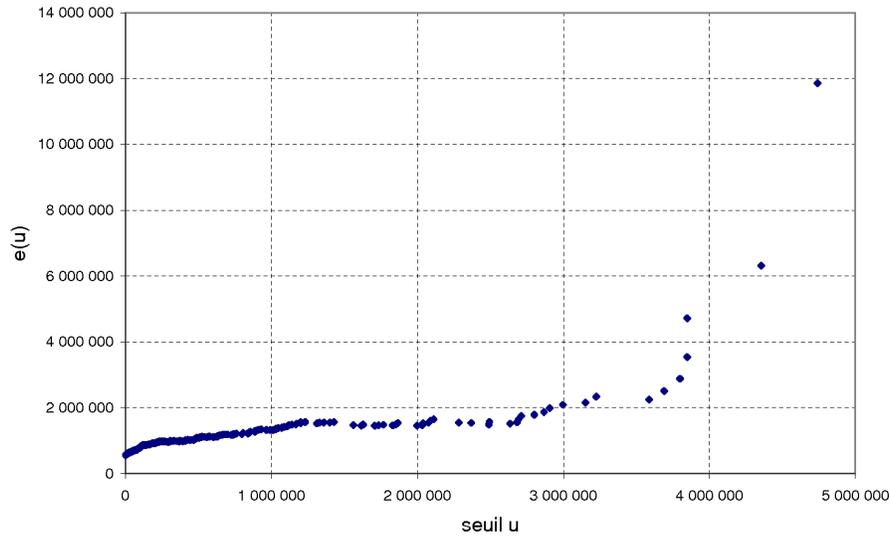


FIG. 35 – Fonction de moyenne en excès pour les sinistres incendie France.

#### 5.2.4 Modélisation à l'aide d'une seule distribution

La première étape de notre modélisation consistait à étudier si une seule loi permettait d'ajuster correctement la sinistralité, aussi bien sur la partie centrale de la distribution qu'en queue de distribution. Nous avons pour cela appliqué les tests d'adéquation de Kolmogorov-Smirnov et Anderson Darling à plusieurs lois usuelles. Le seuil de significativité est toujours 5%. Nous avons obtenu les résultats suivants :

Loi usuelle	Test de Kolmogorov-Smirnov	Test d'Anderson-Darling
Lognormale	Acceptation ( $D=0,053$ )	Acceptation ( $A=2,00$ )
Gamma	Rejet	Rejet
Weibull	Acceptation ( $D=0,067$ )	Acceptation ( $A=1,12$ )
Pareto	Rejet	Rejet
Pareto généralisée	Rejet	Rejet

Ainsi, les tests d'adéquation confirment l'impression visuelle donnée par les QQ Plots, puisque l'adéquation à la loi de Weibull de paramètres  $\alpha = 0,6$  et  $\beta = 355276$  est acceptée, de même que l'adéquation à une loi lognormale de paramètres  $\mu = 11,8$  et  $\sigma = 1,9$ . De plus, on constate que le test de Komogorov-Smirnov accepte en priorité l'adéquation à une loi lognormale (la distance  $D$  est inférieure pour la loi lognormale), alors que le test d'Anderson Darling accepte en priorité l'adéquation à la loi de Weibull (la distance  $A$  est inférieure pour la loi de Weibull). Ceci est certainement dû au fait que la loi de Weibull ajuste mieux nos données au niveau de la queue de distribution.

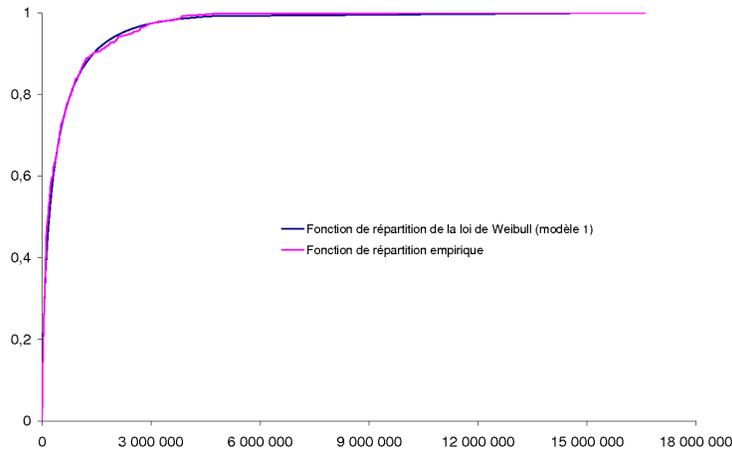


FIG. 36 – Fonction de répartition empirique et fonction de répartition de la loi de Weibull.

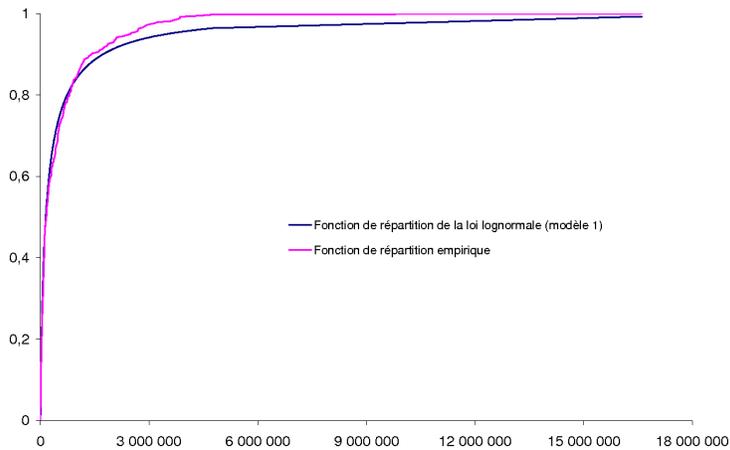


FIG. 37 – Fonction de répartition empirique et fonction de répartition de la loi lognormale.

De manière générale, la loi de Weibull semble mieux ajuster la distribution de nos observations et encore plus particulièrement au niveau de la queue. Il est donc logique que le test d'Anderson Darling, qui accorde plus de poids aux queues de distribution, préfère la loi de Weibull. Nous allons maintenant tenter d'améliorer la modélisation en considérant deux lois.

### 5.2.5 Modélisation à deux lois

#### Ajustement GEV

Nous nous intéressons tout d'abord à la queue de distribution et traçons les deux graphiques qui permettent de distinguer les cas  $\xi = 0$  et  $\xi > 0$ . Si les points sont alignés sur le premier graphique, alors l'hypothèse  $\xi = 0$  est validée. Au contraire, si les points sont alignés sur le deuxième graphique, alors c'est l'hypothèse  $\xi > 0$  qui est validée.

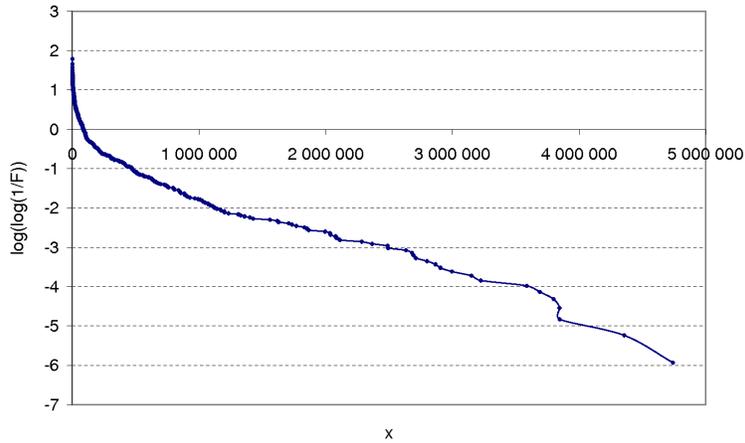


FIG. 38 – Test de  $\xi = 0$ .

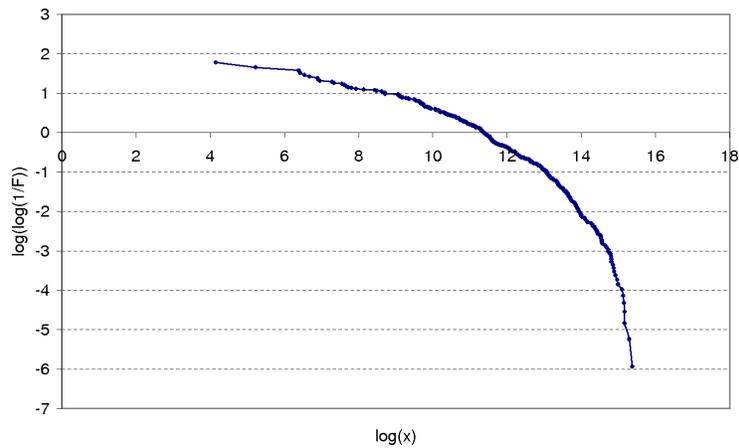


FIG. 39 – Test de  $\xi > 0$ .

Il est important de savoir dans lequel des deux cas ( $\xi = 0$  ou  $\xi > 0$ ) on se trouve car cela conditionne la suite de la modélisation. En effet, pour déterminer le seuil de sinistralité extrême, nous nous basons sur l'estimateur de Hill. Mais ce dernier n'a de sens que si nous sommes dans le cas  $\xi > 0$ . Cependant, ces graphiques ne permettent pas de différencier nettement le cas  $\xi = 0$  du cas  $\xi > 0$  puisqu'aucune des deux courbes obtenues n'est identifiable à une droite.

Nous continuerons la modélisation en supposant que nous sommes dans le cas  $\xi > 0$  et donc que l'estimateur de Hill existe. Cette hypothèse est tout à fait possible étant donné que nous sommes en présence de distributions à queues épaisses (le paramètre  $\alpha$  de la loi de Weibull qui modélise l'attritionnel est strictement inférieur à 1). Or le cas  $\xi = 0$  correspond théoriquement à des distributions de type exponentielle, c'est-à-dire à queues fines.

### Détermination du seuil

Nous mettons à présent en oeuvre l'estimateur de Hill afin de déterminer graphiquement

le seuil départageant les sinistres attritionnels des sinistres extrêmes.  
 Le graphique suivant représente l'estimateur de Hill en fonction du seuil et du nombre d'excès considérés.

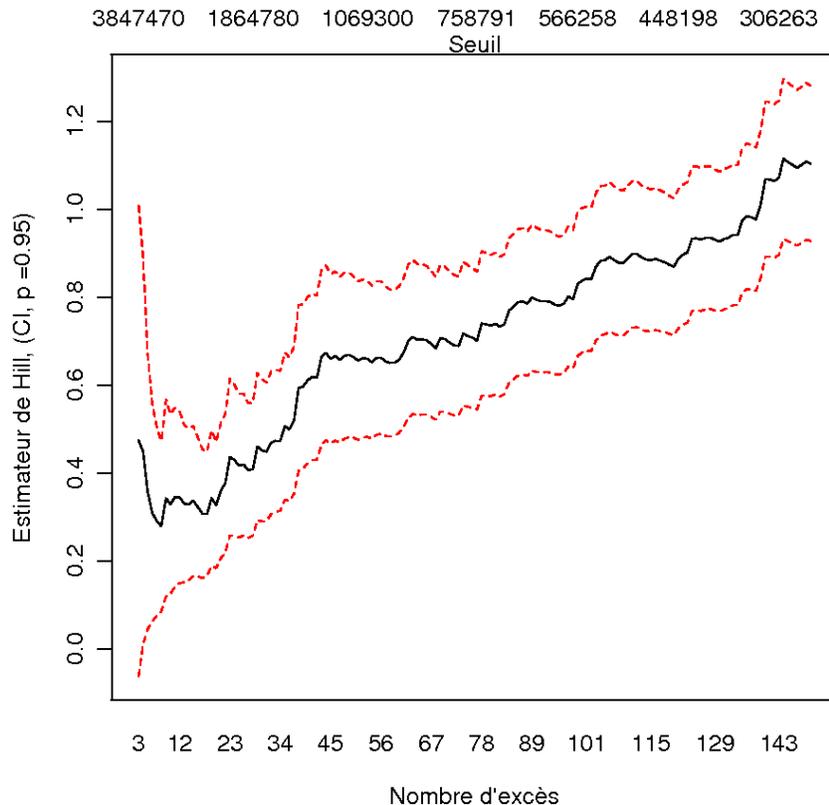


FIG. 40 – Estimateur de Hill en fonction du seuil et du nombre d'excès considérés.

Nous remarquons une zone de stabilité entre 40 et 60 excès. Au delà de 60 excès, l'estimateur n'est plus du tout stable. Nous considérerons donc que l'adéquation à une GPD débute au niveau du 60<sup>ème</sup> excès, soit un seuil à 1 000 000€.

Ainsi, un sinistre inférieur à 1 000 000€ sera considéré comme attritionnel, alors qu'un sinistre supérieur à 1 000 000€ sera considéré comme extrême.

### Ajustement de la partie centrale

Nous avons vu dans la partie "modélisation à une loi" que deux lois étaient acceptées pour modéliser le montant des sinistres : la loi lognormale et la loi de Weibull. Nous nous demandons alors laquelle des deux lois est-il le plus judicieux de garder pour modéliser la partie attritionnelle de notre sinistralité. Sur les figures 36 et 37, nous constatons que sur la partie centrale de la distribution, la courbe de la loi lognormale se situe légèrement au dessus de la courbe empirique, alors que la courbe de la loi de Weibull se situe sur la courbe empirique ou légèrement dessous. La courbe de Weibull aura donc tendance à tarifier un peu plus cher le risque (cf. figure 44). Mais, par principe de prudence, il vaut mieux prendre en

compte plus de risques que pas suffisamment (tout en restant conscient qu'une sur-tarification entraînerait la perte de clients). C'est pourquoi nous avons décidé de conserver l'adéquation à la loi de Weibull pour l'attritionnel.

L'estimation des deux paramètres de la loi de Weibull par la méthode du maximum de vraisemblance conduit à  $\alpha = 0,60$  et  $\beta = 355276$ .

### Ajustement de la queue de distribution

Nous modélisons la sinistralité extrême à l'aide d'une loi de pareto généralisée. L'estimation des deux paramètres de la GPD par la méthode du maximum de vraisemblance conduit à  $\xi = 0,28$  et  $\sigma = 909442$ .

Les fonctions de répartition empirique et modélisée relatives à la sinistralité extrême sont représentées sur la figure suivante.

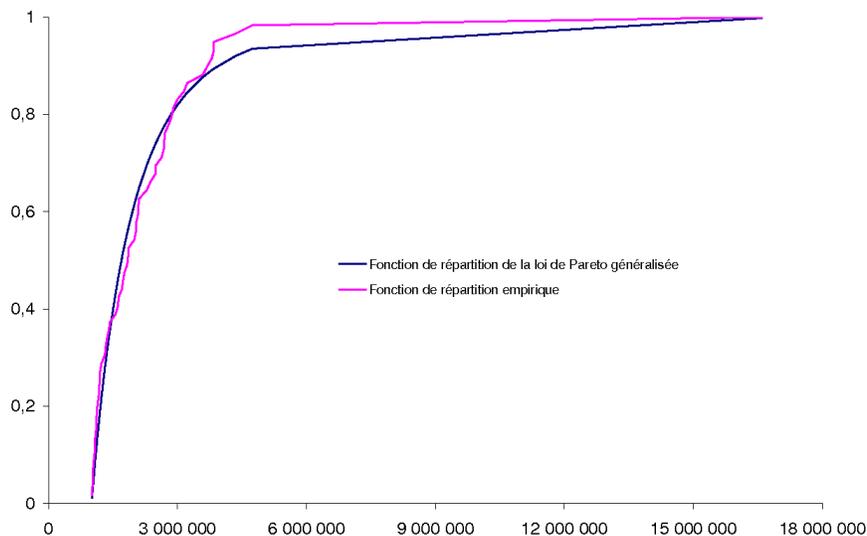


FIG. 41 – Sinistralité extrême empirique et modélisée.

### Modélisation globale

Nous reconstituons à présent le modèle global à deux lois défini dans la partie 4. La figure suivante représente la fonction de répartition empirique et la fonction de répartition modélisée.

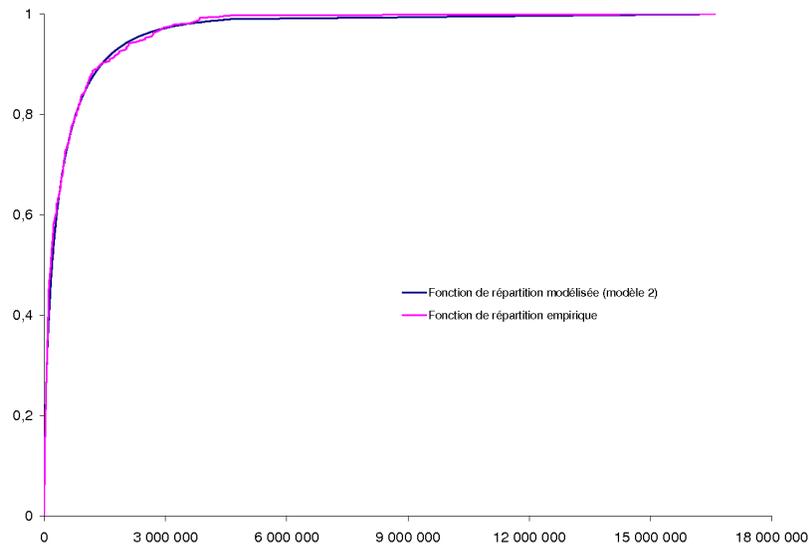


FIG. 42 – Sinistralité empirique et modélisée (modèle à deux lois).

Visuellement, l'adéquation semble bonne. Nous effectuons un zoom sur la queue de distribution afin de comparer le modèle à une loi et le modèle à deux lois.

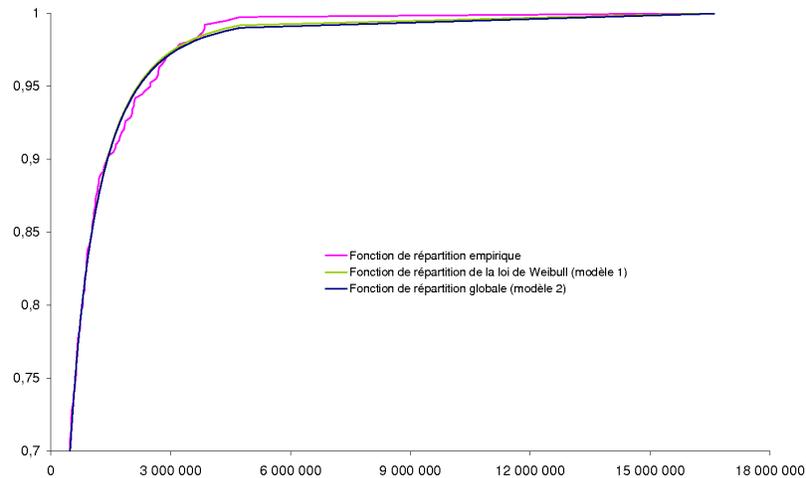


FIG. 43 – Comparaison des modèles 1 et 2 : zoom sur la queue de distribution.

Le résultat obtenu est nettement moins bon que pour l'application DFA. Les deux courbes des deux modèles sont quasiment superposées, l'apport de la modélisation à deux lois semble visuellement ici faible. Cependant, nous allons voir dans la section suivante que les queues de distribution des deux modèles ont des comportements différents, ce qui ne se voit pas ici, mais aura des implications fortes en termes de tarification.

### 5.2.6 Tarification de traités en excédent de sinistre

Maintenant que nous avons modélisé la distribution des sinistres dans son intégralité, il est possible de coter n'importe quel traité en excédent de sinistre, dit "traité XS". Bien que les sinistres proviennent d'une base proportionnelle, il peut être imaginé une structure de réassurance XS sur cette base qui s'apparenterait à une réassurance XS de "pool". Considérons un traité  $pX Sf$  où  $f$  désigne la franchise du sinistre et  $p$  sa portée. Notons  $X$  le montant du sinistre et  $Y$  le montant à la charge du réassureur. Rappelons que le coût à la charge du réassureur  $Y$  est nul si le coût du sinistre est inférieur à la franchise  $f$ ,  $X - f$  si le sinistre est compris entre  $f$  et  $f + p$  et  $p$  si le sinistre excède  $f + p$ . Mathématiquement, on a :

$$Y = \begin{cases} 0 & , X < f \\ X - f & , f < X < f + p \\ p & , f + p < X \end{cases}$$

Plaçons nous dans le cadre du modèle fréquence-sévérité, présenté en annexe 2. Notons  $N$  le nombre de sinistres,  $Y_i$  les différents sinistres et  $S$  la sinistralité, qui a donc pour expression  $S = \sum_{i=1}^N Y_i$ . Afin de nous concentrer sur l'impact des modèles 1 et 2 en terme de prime pure, nous allons utiliser le modèle fréquence-sévérité, mais en faisant l'hypothèse simplificatrice que l'espérance du nombre de sinistre  $\mathbb{E}(N)$  est égale à 1. Dans ce cadre, la prime pure est simplement  $\mathbb{E}(S) = \mathbb{E}(Y) \cdot \mathbb{E}(N) = \mathbb{E}(Y)$ .

Notons  $f_X$  la densité de  $X$ ,  $F_X$  sa fonction de répartition et  $\bar{F}_X = 1 - F_X$  sa fonction de survie,  $\mathbb{E}(Y)$  s'exprime alors sous la forme :

$$\begin{aligned} \mathbb{E}(Y) &= \int_f^{f+p} (x - f) f_X(x) dx + \int_{f+p}^{+\infty} p \cdot f_X(x) dx \\ &= \int_f^{f+p} x \cdot f_X(x) dx - f \int_f^{f+p} f_X(x) dx + p \int_{f+p}^{+\infty} f_X(x) dx \\ &= \int_f^{f+p} x \cdot f_X(x) dx - f \cdot (F_X(f+p) - F_X(f)) + p \cdot (1 - F_X(f+p)) \\ &= [x \cdot (F_X(x) - 1)]_f^{f+p} - \int_f^{f+p} (F_X(x) - 1) dx - f \cdot (F_X(f+p) - F_X(f)) + p \cdot (1 - F_X(f+p)) \\ &= \int_f^{f+p} \bar{F}_X(x) dx \end{aligned}$$

La valeur de la prime pure du contrat XS a donc une interprétation géométrique simple : elle est égale à l'aire située au dessus de la courbe de la fonction de répartition, plus précisément l'aire délimitée par la courbe de la fonction de répartition, les droites verticales d'équation  $y = f$  et  $y = f + p$ , et la droite horizontale d'équation  $x = 1$ .

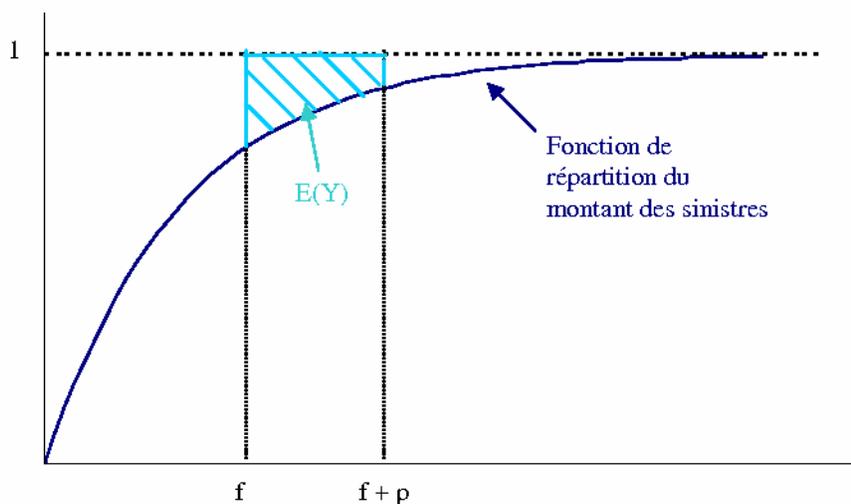


FIG. 44 – Illustration de la prime pure d'un traité  $pXSf$ .

Il est maintenant possible de calculer la prime pure. Le premier contrat XS considéré cote une tranche située de part et d'autre du seuil, c'est-à-dire qu'elle fait intervenir la sinistralité attritionnelle et la sinistralité extrême, tandis que les contrats XS suivants font intervenir la queue de distribution. Les valeurs des primes pures obtenues en utilisant les modèles 1 et 2 sont synthétisées dans le tableau suivant.

Contrat	Modèle 1	Modèle 2	Ecart relatif
1 000 000 XS 500 000	166 074	166 359	0,17%
1 000 000 XS 2 000 000	39 783	40 972	2,99%
1 000 000 XS 3 000 000	18 925	20 529	8,48%
1 000 000 XS 4 000 000	9 784	11 526	17,80%
1 000 000 XS 5 000 000	5 356	7 019	31,05%
1 000 000 XS 10 000 000	432	1184	174,07%
1 000 000 XS 15 000 000	55	364	561,82%
1 000 000 XS 20 000 000	9	150	1566,67%

Les contrats  $XS$  choisis ont tous une portée de 1 000 000€. Nous constatons qu'à portée constante, les primes pures diminuent avec la franchise, ce qui est normal car les tranches ont de moins en moins de chance d'être traversées. En comparant les modèles 1 et 2 sur le premier contrat, qui fait intervenir à la fois la sinistralité attritionnelle et la sinistralité extrême puisque le seuil  $u$  est de 1 000 000€, on constate que les deux primes sont très proches (écart relatif de 0,17%), le modèle 2 ayant tendance à être légèrement plus prudent. Les deux contrats suivants concernent des tranches totalement travaillantes et les résultats sont encore assez proches, le modèle 2 étant cette fois clairement plus prudent. Pour les contrats suivants, qui concernent des tranches peu travaillantes à non travaillantes, les résultats des deux modèles s'éloignent de plus en plus. Le modèle 1 fait intervenir une loi de Weibull, plus fine que la loi de Pareto généralisée, et semble mal adapté à la tarification des tranches

hautes.

Ainsi, quand les tranches sont très peu travaillantes, comme les deux avant-dernières (seul le sinistre maximum traverse), ou totalement non-travaillantes, comme la dernière, le premier modèle peut présenter un risque important de sous-tarification. En revanche, le modèle 2 semble fournir des résultats plus cohérents, la loi de Pareto généralisée étant beaucoup mieux adaptée à la modélisation des queues de distribution.

## Conclusion

La sinistralité issue de l'activité de réassurance, bien que de nature complexe, a donc pu être modélisée dans son ensemble. Les outils mathématiques utilisés ont été principalement la théorie des tests et la théorie des valeurs extrêmes. Des outils graphiques de détermination du seuil ont été présentés, le plus efficace étant selon nous l'estimateur de Hill. Nous préconisons donc son utilisation pour la détermination du seuil, en suivant ce protocole : en partant des observations les plus extrêmes, rechercher une zone de stabilité de l'estimateur de Hill, puis prendre le seuil à la fin de cette zone de stabilité, juste avant un changement de comportement. La détermination n'est cependant pas toujours des plus aisées. L'utilisation de méthodes graphiques présente toutefois l'avantage de visualiser la situation, la modélisation des queues de distribution nécessitant souvent beaucoup d'attention. Plusieurs modèles de la sinistralité dans son ensemble ont été définis théoriquement et deux d'entre eux ont été implémentés informatiquement sous la forme d'un outil Excel, ce qui nous a permis d'étudier leur comportement sur données simulées afin de les valider puis de les utiliser sur données réelles.

Sur les deux jeux de données réelles considérés, la sinistralité a pu être modélisée dans son ensemble. Concernant le premier jeu de données, relatives à la responsabilité civile automobile en Grande-Bretagne, le modèle 2 rend particulièrement bien compte de la sinistralité dans son ensemble. Dans le cadre d'une approche allocation de capital, nous avons estimé les quantiles extrêmes en utilisant les différentes approches afin de les comparer. Avec le second jeu de données, relatives à des sinistres incendie en France, l'adéquation au modèle 2 est moins spectaculaire. Cependant, l'étude de la tarification de traités en excédent de sinistre permet de se rendre compte de l'importance de l'utilisation d'une loi adaptée pour la modélisation de la queue de distribution. De manière générale, pour les deux jeux de données étudiés, le deuxième modèle présentait un caractère plus prudentiel que le premier qui avait tendance à sous-estimer le risque et à le tarifer moins cher.

## Références

- [1] ANDERSON C., O'HAGAN A., TANCREDI A. (2006), 'Accounting for threshold uncertainty in extreme value estimation', *Extremes*, Vol. 9, 87-106.
- [2] BEIRLANT J., MATTHYS G., DIERCKX G. (2001), 'Heavy-Tailed Distributions and rating', *ASTIN Bulletin*, Vol. 31, No. 1.
- [3] BEIRLANT J., TEUGELS J., VYNCKIER P. (1996), *Practical Analysis of Extreme Values*, Leuven University Press.
- [4] BLONDEAU J., PARTRAT C. (2003), *La Reassurance : Approche Technique*, Economica.
- [5] CEBRIAN A. C., DENUIT M., LAMBERT P. (2003), 'Generalized Pareto fit to the society of actuaries large claims database', *North American Actuarial Journal* 7, 18-36.
- [6] CHARPENTIER A., *Théorie des extrêmes et couverture des catastrophes*, Polycopié de cours de l'ENSAE.
- [7] COLES S. (2001), *An Introduction to Statistical Modelling of Extreme Values*, Springer London.
- [8] D'AGOSTINO R. ET STEPHENS M. (1986), *Goodness of Fit Technics*.
- [9] DAVID C. M. DICKSON, LEANNA M. TEDESCO, BEN ZEHNWIRTH (1998), 'Predictive Aggregate Claims Distributions', *The Journal of Risk and Insurance*, Vol. 65, No. 4.
- [10] DAVISON A., SMITH R. (1990), 'Models for exceedances over high thresholds', *Journal of Royal Statistical Society*, Vol. 52, No. 3, 393-442.
- [11] DUMAS S. (2000), *Mise en place d'un outil de test statistique appliqué aux données de réassurance*, Mémoire ISUP.
- [12] EMBRECHTS P., KLÜPPELBERG C., MIKOSH T. (1997), *Modelling extremal events for insurance and finance*, Springer Berlin.
- [13] GARRIDO M. (2002), *Modélisation des évènements rares et estimation des quantiles extrêmes, méthodes de sélection de modèles pour les queues de distribution*, Thèse de l'Université Joseph Fourier, Grenoble I.
- [14] GONZALO J., OLMO J. (2004), 'Which Extreme Values Are Really Extreme?', *Journal of Financial Econometrics*, Vol. 2, No. 3, 349-369.
- [15] GUILLOU A., HALL P. (2001), 'A diagnostic for selecting the threshold in Extreme Value Analysis', *Journal of the royal Statistical Society*, Vol. 63, No. 2.
- [16] HALL P. (1990), 'Asymptotic Properties of the bootstrap for Heavy-Tailed Distributions', *The Annals of Probability*, Vol. 18, No. 3.
- [17] HOSKING J., WALLIS J. (1987), 'Parameter and quantile estimation for the Generalized Pareto Distribution', *Technometrics* 29.

- [18] KLUGMAN S., PANJER H. ET WILLMOT G. (1998), *Loss Models : From data to decisions*, Wiley Series in Probability and Statistics.
- [19] MC NEIL A. (1997), 'Estimating the tails of loss severity distributions using extreme value theory', *ASTIN Bulletin*, Vol. 27, 117-137.
- [20] MC NEIL A. (1998), 'Calculating Quantile Risk Measures for Financial Return Series using Extreme Value Theory'.
- [21] MC NEIL A. (1999), 'Extreme Value Theory for Risk Managers'.
- [22] NAESS A., CLAUSEN P.H. (2000), 'The peaks over threshold method and bootstrapping for estimating long return period design values'.
- [23] RONCALLI T., *Théorie des valeurs extrêmes ou modélisation des évènements rares pour la gestion des risques*, Polycopié de cours du DESS 203 de l'Université Paris IX Dauphine.
- [24] SCOR (2006), 'Le Modèle Fréquence - Sévérité : Applications au pricing et fondements théoriques'.
- [25] SCOR, GROUP ACTUARIAL DEPARTMENT (2006), 'Hybrid Severity PDF. Benchmark'.
- [26] STEPHENS M. (1974), 'EDF Statistics for Goodness of Fit and some Comparisons', *Journal of the American Statistical Association*, Vol. 69, No. 347, 730-737.

## Annexe 1 : L'outil Excel

La partie pratique de ce mémoire a donné lieu au développement d'un outil Excel qui regroupe tous les outils nécessaires à la mise en oeuvre des modèles 1 et 2 présentés dans la partie 4 de ce mémoire. Il fournit en particulier les outils utilisés dans la détermination du seuil de sinistralité extrême (Graphe de l'estimateur de Hill, mean excess function et Gerstengarbe plot) et il permet de construire le modèle d'ajustement à une loi (modèle 1) ainsi que le modèle d'ajustement à deux lois (modèle 2).

Nous allons présenter plus précisément les fonctionnalités de l'outil ainsi que la manière dont elles ont été mises en place, et ce feuille par feuille du fichier excel. Nous précisons également que pour fonctionner, le programme a besoin du solveur d'Excel. Il est donc nécessaire de l'installer à partir du menu "outils / Macros complémentaires".

### Les feuilles exploratoires

Dans notre outil, trois feuilles sont consacrées à l'étude exploratoire préliminaire de l'échantillon étudié.

#### Feuille "Données"

C'est dans cette feuille que l'utilisateur est invité à entrer la série des montants qu'il souhaite modéliser. Il s'agit des montants ultimes déjà revalorisés. L'utilisateur entre également le seuil de sinistralité extrême qu'il aura choisi à l'aide des outils présentés plus loin.

En cliquant sur le bouton "GO", il obtiendra diverses statistiques descriptives sur son échantillon. Le bouton "Mise en forme" permet d'initialiser les autres feuilles du fichier, en vue de leur utilisation future.

#### Feuille "Répartition"

Cette feuille permet d'obtenir la fonction de répartition empirique de l'échantillon étudié, ainsi que l'histogramme de la série (comme la figure [21] pour les données Responsabilité Civile Automobile).

#### Feuille "QQPlot"

Cette feuille fournit le graphe des QQ plot pour les quatre lois usuelles qui sont la loi exponentielle, la loi lognormale, la loi de Pareto et la loi de Weibull.

	A	B	C	D	E	F	G	H	I
1	Echantillon								
2	3401,07957						<b>SINISTRALITE TOTALE</b>		
3	19966,73303		Gal				Nombre d'observations	378	
4	13263,63533						Sinistre minimum	63	
5	212169,2909						Sinistre moyen	547 495	
6	30238,84608		Mise en forme				Sinistre maximum	16 596 085	
7	211050,9309						Ecart-type	1 149 631	
8	44120,59016						Sinistre médian	160 944	
9	1560615,404						Quantile à 90%	1 408 458	
10	257069,0037						Quantile à 95%	2 489 314	
11	613,9043011						Quantile à 99%	3 646 749	
12	1864784,904		Seuil (à rentrer par l'utilisateur)	1 000 000					
13	25673,2169								
14	23481,24685						<b>SINISTRALITE ATTRITIONNELLE</b>		
15	8741,320737						Nombre d'observations	319	
16	30100,84793						Sinistre minimum	63	
17	16890,39939						Sinistre moyen	221 476	
18	471245,0178						Sinistre maximum	991 527	
19	67041,22949						Ecart-type	251 602	
20	156659,3422						Sinistre médian	104 071	
21	305680,1017						Quantile à 90%	635 775	
22	1707617,94						Quantile à 95%	799 064	
23	171449,7724						Quantile à 99%	909 259	
24	2076911,07								
25	49194,90536								
26	40030,96528						<b>SINISTRALITE EXTREME</b>		
27	17294,33149						Nombre d'observations	59	

FIG. 45 – La feuille "données" de l'outil excel pour les données Incendie France.

## Les feuilles de détermination du seuil

Après la phase exploratoire, des feuilles sont consacrées à l'étude du seuil de sinistralité extrême. On retrouve ainsi les outils présentés en partie 3 de ce mémoire.

### Feuille "Mean excess"

Cette feuille fournit le graphe de la mean excess fonction empirique de notre échantillon

### Feuille "Gerstengarbe"

Cette feuille trace le graphe de Gerstengarbe et permet d'obtenir la valeur numérique précise du seuil à retenir selon cette méthode. Cependant, nous avertissons l'utilisateur du manque de robustesse de cette méthode.

### Feuille "Ajustement GEV"

Cette feuille trace les deux graphes d'adéquation aux lois de Gumbel et de Fréchet. Il est nécessaire de l'étudier avant d'appliquer la feuille "Estimateur de Hill", afin de s'assurer de la validité de l'estimateur de Hill.

## Feuille "Estimateur de Hill"

Cette feuille présente sans doute l'outil le plus pertinent pour déterminer le seuil d'entrée dans la zone de sinistre extrême. Elle permet d'obtenir le Hill Plot en fonction du nombre d'excès considérés.

## Les feuilles de modélisation

Une fois le seuil déterminé grâce aux outils précédents, nous pouvons passer à la modélisation de la sinistralité. Elle s'effectue en trois étapes. Tout d'abord nous ajustons la sinistralité à l'aide d'une seule loi. Puis nous essayons d'améliorer l'estimation en queue de distribution (au delà du seuil choisi). Enfin, nous construisons la modélisation à deux lois.

### Modélisation à une loi

La modélisation à une loi a été mise en place dans le cas particulier des lois lognormales et de Weibull (feuilles "Fit Weibull" et "Fit LN"). Le programme estime les paramètres à partir de l'échantillon en utilisant la méthode du maximum de vraisemblance. Pour cela, il faut une première estimation des paramètres de la loi considérée : c'est le calcul effectué par la première approche. Les valeurs trouvées sont ensuite recopiées dans les cases paramètres et servent au calcul des densités de probabilité  $f_X(X_i)$ . Enfin, dans la case "somme", nous calculons la valeur de la logvraisemblance :

$$l = \sum_{i=1}^n \ln[f_x(x_i)] - n \ln[F_X(s) - F_X(d)]$$

où  $s$  est le seuil de troncature à gauche et  $d$  le seuil de troncature à droite (facultatifs).

Pour la recherche des estimateurs, nous devons maximiser cette fonction ; le programme va donc lancer le solveur qui va maximiser la valeur de  $l$  en modifiant la valeur des paramètres. Le solveur utilise la méthode de résolution numérique de Newton. Au  $l$  maximum correspondront les estimateurs du maximum de vraisemblance.

Les estimateurs calculés en première approche sont très importants car le solveur d'Excel a besoin de valeurs initiales les plus proches possibles de la solution pour converger.

Pour la première approche de la loi lognormale, nous avons utilisé la transformation de la loi Lognormale en une loi Normale par passage au logarithme ; nous obtenons alors en première approche :

$$\begin{cases} \hat{\mu} = \frac{1}{n} \sum_{i=1}^n \ln(X_i) \\ \hat{\sigma} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\ln X_i)^2 - \hat{\mu}^2} \end{cases}$$

Pour la loi de Weibull, la première approche par la méthode des moments est difficilement exploitable car l'espérance et la variance de la loi de Weibull dépendent de la fonction Gamma

qui n'est pas directement intégrée sous Excel. Par contre, la fonction de répartition d'une loi de Weibull ayant une écriture simple, nous pouvons utiliser les quartiles pour l'initialisation du solveur.

Soit  $Q_{25}$  le premier quartile (25e percentile) et  $Q_{75}$  le troisième quartile (75e percentile); nous obtenons le système suivant :

$$\begin{cases} \frac{1}{4} = 1 - \exp[-(\frac{Q_{25}}{\beta})^\alpha] \\ \frac{3}{4} = 1 - \exp[-(\frac{Q_{75}}{\beta})^\alpha] \end{cases}$$

Ce qui nous conduit à :

$$\begin{cases} \hat{\alpha} = \frac{\ln(\frac{\ln 4}{\ln 4 - \ln 3})}{\ln(\frac{Q_{75}}{Q_{25}})} \\ \hat{\beta} = \frac{Q_{75}}{(\ln 4)^{1/\hat{\alpha}}} \end{cases}$$

Une fois les estimateurs obtenus, nous pouvons calculer la fonction de répartition théorique et tracer le graphe la comparant à la fonction de répartition empirique.

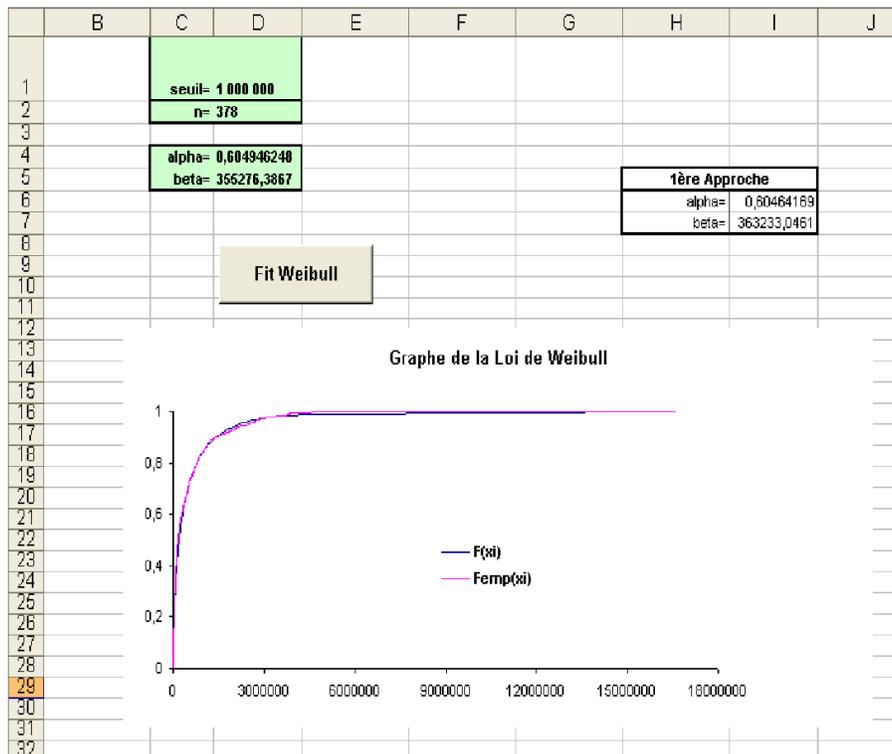


FIG. 46 – La feuille "Fit Weibull" de l'outil excel pour les données Incendie France.

## Feuille "fit GPD"

Pour améliorer l'adéquation en queue de distribution, nous ajustons une loi GPD aux données excédant le seuil entré par l'utilisateur. L'ajustement se fait de la même manière que dans la partie précédente. Nous calculons les estimateurs par la méthode du maximum de vraisemblance. Les estimateurs utilisés en première approche sont les estimateurs des moments (Hosking and Wallis 1987) définis par :

$$\begin{cases} \hat{\xi} = \frac{1}{2} \left(1 - \frac{\bar{X}^2}{s^2}\right) \\ \hat{\sigma} = \frac{1}{2} \bar{X} \left(1 + \frac{\bar{X}^2}{s^2}\right) \end{cases}$$

où  $\bar{X}$  est la moyenne de l'échantillon et  $s^2$  la variance.

## Feuille "fit Global"

Enfin, la feuille "fit Global" permet d'effectuer la modélisation à deux lois en compilant les résultats des ajustements précédents et en les pondérant par le facteur de normalisation  $\alpha$  défini dans la partie 4 du mémoire.

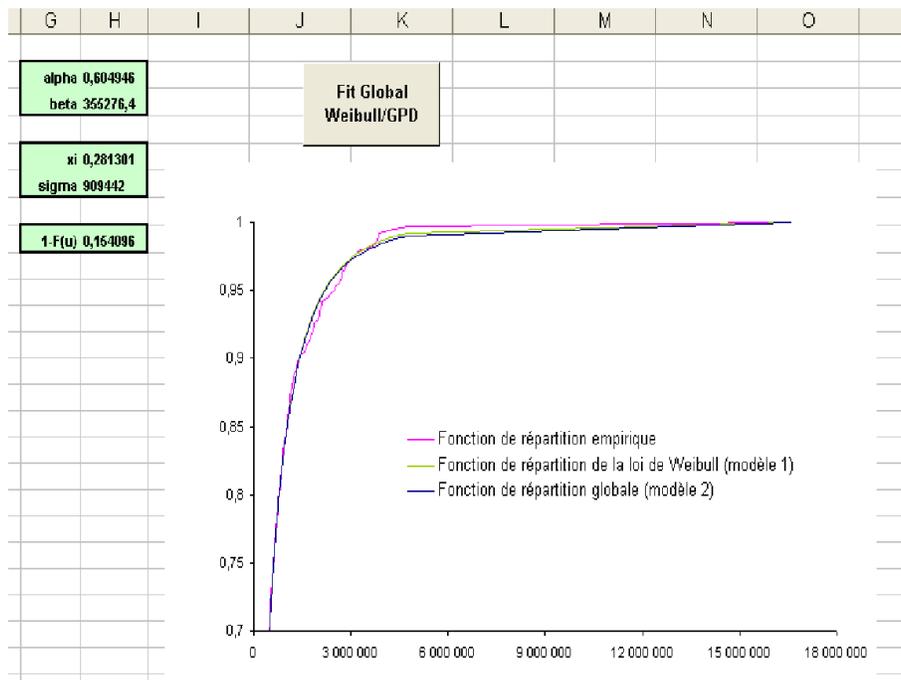


FIG. 47 – La feuille "Fit Global" de l'outil excel pour les données Incendie France.

## Annexe 2 : Rappel sur le modèle de risque collectif

Le modèle de risque collectif considère le montant total des sinistres d'un portefeuille composé de plusieurs polices homogènes. Une police peut donner lieu à plusieurs sinistres. Si les  $X_i$  sont les montants des sinistres individuels et  $N$  est le nombre total de sinistres pour toutes les polices du portefeuille, alors le montant agrégé est donné par :

$$S = X_1 + X_2 + \dots + X_N$$

Les hypothèses du modèle sont les suivantes :

- les montants des sinistres  $X_i$  sont indépendants et identiquement distribués,
- le nombre de sinistres  $N$  est indépendant des  $X_i$ .

De ce modèle découle la formule fréquence / sévérité utile dans le calcul des primes pures. Nous avons effectivement le résultat suivant :

$$E(S) = E(N)E(X)$$

En effet :

$$\begin{aligned}\mathbb{E}(S) &= E[E(S/N)] \\ &= \sum_{n=0}^{\infty} P(N = n)E(S/N = n) \\ &= \sum_{n=0}^{\infty} P(N = n)E(X_1 + \dots + X_n) \\ &= \sum_{n=0}^{\infty} P(N = n)nE(X) \\ &= E(N)E(X)\end{aligned}$$

## Annexe 3 : Remarque sur l'impact de la méthode du Chain Ladder

Dans le cas de traités à développement long, le prix du sinistre à payer par le réassureur évolue au cours du temps, en fonction des réévaluations successives du sinistre. Ces réévaluations successives sont courantes en responsabilité civile. Ce problème ne se pose pas pour les traités à développement court, où le montant du sinistre varie très rarement après quelques années.

Le développement long conduit à utiliser la méthode du Chain Ladder pour obtenir les ultimes. Ainsi est estimée la sinistralité empirique, qui est ensuite modélisée. Cependant, plus l'année de survenance du sinistre est récente, moins nous disposons d'information sur l'historique d'évolution des montants. La méthode du Chain Ladder a pour but de rendre les conditions identiques quelle que soit l'année de survenance du sinistre. Toutefois, comme il y a de moins en moins d'information disponible, la méthode du Chain Ladder introduit de la volatilité. Dans le cadre de l'utilisation de ces ultimes, il est légitime de se demander si cette méthode du Chain Ladder n'implique pas un épaississement de la queue de distribution des ultimes ou une modification du seuil départageant un sinistre attritionnel d'un sinistre extrême. La réponse à cette question est complexe et nous tenterons uniquement d'illustrer la situation par un exemple.

Considérons le jeu de données relatives aux sinistres responsabilité civile automobile en Grande Bretagne que nous avons étudié dans la partie 5 de ce mémoire. Nous sommes bien dans le cas de traités à développement long et nous avons utilisé la méthode du Chain Ladder pour obtenir les sinistres ultimes. Commençons par représenter les moyennes et écarts type par année de survenance du sinistre, puis différents quantiles resserrés au niveau des quantiles extrêmes.

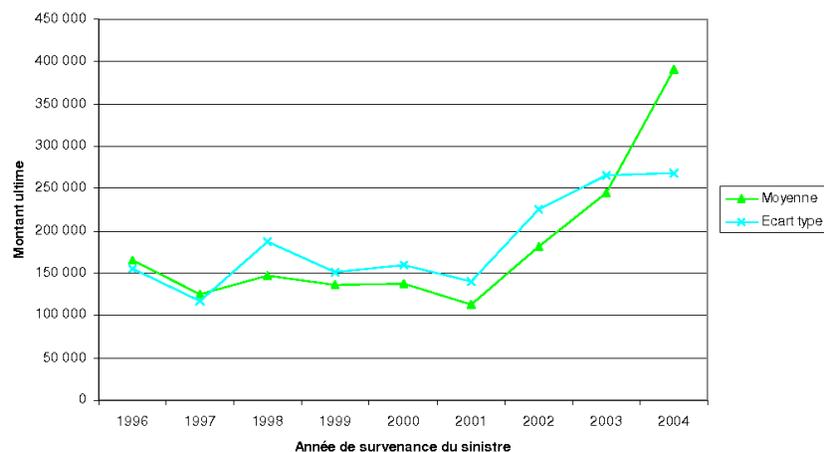


FIG. 48 – Moyennes et écarts type.

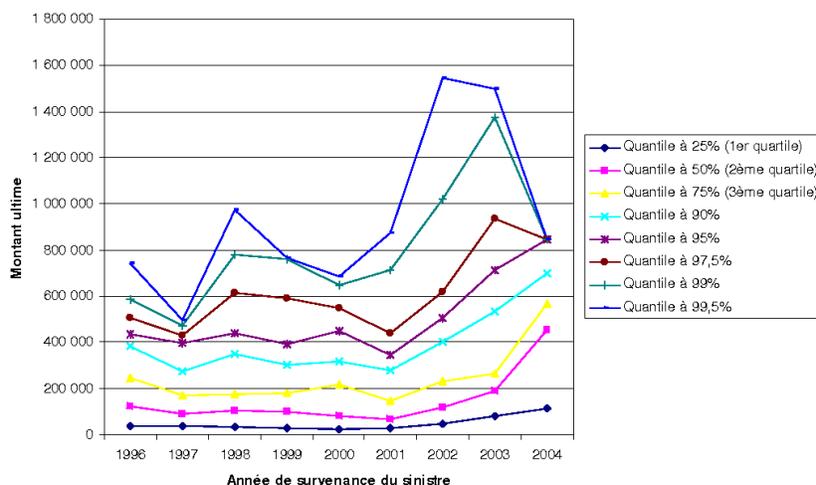


FIG. 49 – Quantiles.

Nous observons bien une tendance à l'augmentation de la volatilité de 2001 à 2004, quantifiée ici par l'écart type. La moyenne semble également augmenter sensiblement sur la même période. En observant la tendance des quantiles extrêmes, nous constatons que la queue de distribution a tendance à s'épaissir, phénomène dû à l'utilisation de la méthode du Chain Ladder.

Traçons maintenant l'estimateur de Hill en ne considérant d'une part que la première moitié de l'échantillon, c'est-à-dire relative aux sinistres les plus anciens, d'autre part la seconde moitié de l'échantillon, relative aux sinistres les plus récents. Nous cherchons à répondre à deux questions : d'une part, y a-t-il un épaississement de la queue indiqué par l'estimateur de Hill, d'autre part, est-ce que la sélection du seuil est impactée ?

En observant ces deux graphiques (figures 50 et 51), nous constatons que l'indice de queue  $\xi$  est plus élevé sur la seconde moitié de l'échantillon. Cependant, en cherchant une zone de stabilité de l'estimateur de Hill juste avant une zone de croissance, nous sommes amenés dans les deux cas à considérer entre 75 et 80 excès, ce qui conduit à des seuils comparables. Sur cet exemple, la conclusion est que la méthode du Chain Ladder tend à épaissir la queue de la distribution empirique, mais la détermination du seuil en utilisant l'estimateur de Hill est très peu impactée.

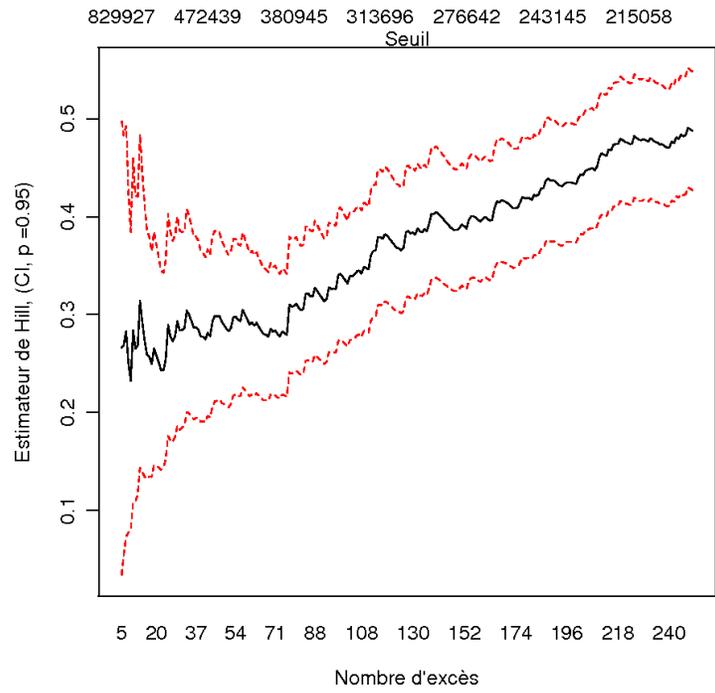


FIG. 50 – Estimateur de Hill relatif à la première moitié de l'échantillon.

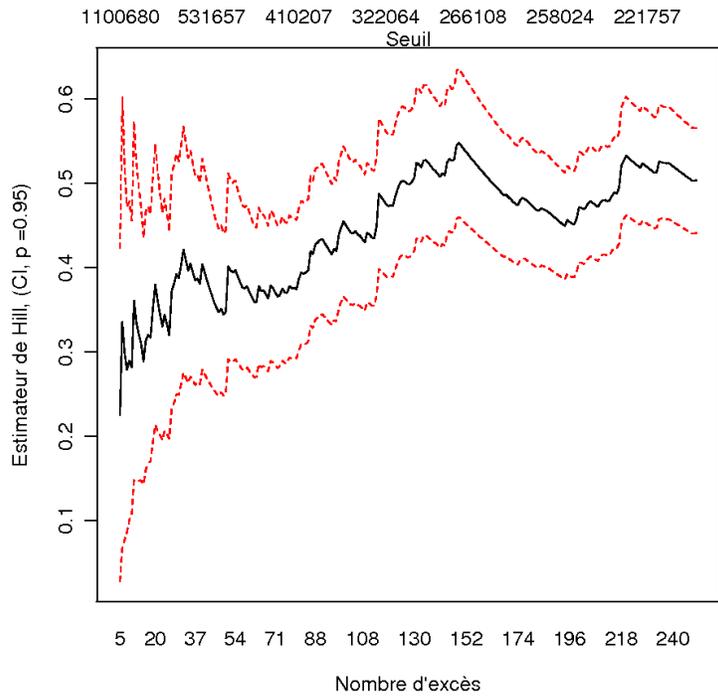


FIG. 51 – Estimateur de Hill relatif à la seconde moitié de l'échantillon.