

Tarification en assurance IARD



MASTER 1 ECONOMETRIE ET STATISTIQUES PARCOURS
INGENIERIE FINANCIERE

JANVIER 2021

JEREMY BRUNNER

PIERRE DENEUX

Table des matières

1. Statistiques descriptives.....	2
2. Modélisation.....	3
2.1 La régression de poisson	4
2.2 La régression binomiale négative.....	6
2.3 Modèles Zéro-Inflatés (ZIP, ZINB).....	7
3. Modélisation pour la sévérité	8
3.1 Le modèle log normal	8
3.2. Le modèle Gamma	9
4. Sinistres graves.....	9
5. Prime Pure	10
Références :.....	11
Annexe :.....	12
1. Statistiques descriptives :.....	12
2. Régression sur la fréquence :	14
3. Régression pour la sévérité	16

1. Statistiques descriptives

Les données se décomposent en deux fichiers : le fichier des « Effectifs » concernant les informations sur la population assurée et le fichier des « sinistres » portant sur le nombre et le montant des prestations dont ont bénéficié ces mêmes assurés. Les fichiers initiaux concernent les effectifs sur l'exercice « Year 0 ».

On a notre disposition un échantillon de 100 000 observations d'un portefeuille d'une compagnie d'assurance. Pour chaque assuré de notre échantillon, nous disposons de 4 groupes de variables : les caractéristiques du conducteur (sexe, âge, ancienneté du permis, son taux Bonus-Malus, ...), les caractéristiques du véhicule (ancienneté, puissance, valeur, type, ...), le type de contrat (type de formule choisie par l'assuré) et la sinistralité (nombre de sinistres, montant de sinistres).

Le tableau ci-dessous présente la répartition de l'échantillon selon le nombre de sinistres déclarés durant l'année 0 :

Nombre de sinistre	Fréquence	Pourcentage	Pourcentage cumulé
0	88188	88.188	88.188
1	10540	10.540	98.728
2	1144	1.144	99.872
3	109	0.109	99.981
4	15	0.015	99.996
5	3	0.003	99.999
6	1	0.001	100

On peut observer qu'une grande part du notre portefeuille n'ont pas eu de sinistres. La part des assurés ayant eu un sinistre s'élève seulement à 11.812%. Parmi les individus ayant eu un sinistre 10.54% ont eu 1 sinistre, 1,144% ont eu 2 sinistres, 0.109% ont eu 3 sinistres, 0.015% ont eu 4 sinistres, 0.003% ont eu 5 sinistres et 0.001% ont eu 6 sinistres. On peut observer cette distribution sur l'histogramme tracé, que l'on peut retrouver en annexe (figure 1.1). L'allure de cette distribution met en évidence la répartition exprimée ci-dessus.

	Conducteur principal	Conducteur secondaire
Femme	39.766	20.235
Homme	60.234	12.865

On remarque, 60% des conducteurs principaux sont des hommes et 40% des femmes. Ils ont la possibilité de déclarer un deuxième conducteur sur la même voiture, 20% des conducteurs secondaires sont des femmes et 12% des hommes et 67% des individus n'ont pas de conducteurs secondaires.

Sur la figure 1.2-4 graphique en annexe des effectifs en fonction de l'âge du véhicule on peut observer une relation inverse entre l'âge et le nombre de sinistre. En effet lorsque l'âge des véhicules augmente on observe une diminution du nombre de sinistre déclaré.

Afin de mieux cerner la population étudiée, nous avons effectué une série de statistiques. Nous pourrions ainsi mieux appréhender les variables explicatives de la sinistralité sur ce portefeuille.

Type de formule	Mini	Médian 1	Médian 2	Maxi
% des assurés	8.510	9.320	17.316	64.854

La société d'assurance à quatre types de garanties pour l'assurance d'un véhicule. « Mini », l'assurance minimale obligatoire 8,5% de nos assurés ont choisi cette formule. 9,32% des assurés ont opté pour la formule « Median1 », 17,31% des clients ont choisi la formule « Median2 », enfin, 64,85% des clients ont choisi la formule « Maxi ». Ici on remarque que 65% des clients préfèrent se couvrir au maximum on peut donc les caractériser de risquophobes.

2. Modélisation

L'objet de cette section est de définir une méthode de tarification des contrats d'assurance. La méthode la plus couramment utilisée aujourd'hui est l'approche « Fréquence * Coût moyen ».

La réalisation d'un tarif en assurance IARD (auto, MRH, construction, etc.) s'appuie classiquement sur l'analyse de la prime pure dans le cadre d'un modèle fréquence x coût dans lequel l'effet des variables explicatives sur le niveau du risque est modélisé par des modèles de régression de type GLM.

Le cadre usuel de tarification : en pratique la tarification IARD est souvent effectuée dans le cadre très général des modèles fréquence-coût :

$$S = \sum_{i=1}^N C_i + I_G \times G$$

Avec N le nombre de sinistres (souvent supposé suivre une loi de Poisson), C le coût unitaire d'un sinistre (en général gamma ou log-normal), IG l'indicatrice de survenance d'un sinistre grave et G le coût d'un sinistre grave. On observera par la suite qu'on négligera la précision par les sinistres graves en raison de faible représentation

Le cadre usuel de tarification :

Sous réserve de l'indépendance de la fréquence et des coûts, la prime pure à l'intérieur d'une classe de risque est de la forme :

$$E[S|X] = E[N - I_G|X] \times E[C|X] + P(I_G = 1|X) \times E[G|X]$$

On se ramène ainsi à modéliser l'espérance conditionnelle du nombre de sinistres et l'espérance conditionnelle du coût unitaire. Il s'agit donc d'estimer des espérances conditionnelles, ce qui est le cadre général des modèles de régression, et plus particulièrement des modèles de régression non linéaires (GLM).

On étudie dans cette partie la façon dont les modèles linéaires généralisés sont appliqués à notre étude. Notamment on se penche sur les paramétrages nécessaires à la prise en compte des variables explicatives, ainsi que les modélisations qui sont retenues pour nos deux variables expliquées que sont la fréquence moyenne des sinistres et leur coût moyen annuel.

Dans un premier temps, on se propose de modéliser la fréquence des accidents déclarés par l'assureur. Par la suite on présentera les quatre modèles choisis, les variables retenues et leurs limites d'application.

Nous proposons dans ce qui va suivre d'estimer les lois de probabilités permettant de modéliser le phénomène de la survenance d'accidents et de déterminer la tarification en assurance automobile tenant compte du nombre de sinistres passés.

2.1 La régression de poisson

La régression poissonnienne est une loi des événements rares elle modélise bien le nombre de sinistres d'une police individuelle.

La loi de Poisson de paramètre λ , est une distribution qui permet de modéliser la survenance d'accidents durant une période donnée. La probabilité d'avoir « y » accidents est égale à :

$$P_{\lambda}(y) = P_{\lambda}(Y = y) = e^{-\lambda} \frac{\lambda^y}{y!}, y=0,1,2,\dots$$

L'estimation des probabilités p_y de la loi de Poisson revient à estimer λ .

Le paramètre (λ) de la loi de Poisson est estimé ici par l'estimateur du maximum de vraisemblance. L'estimation du maximum de vraisemblance a pour but de trouver les valeurs possibles d'un paramètre afin d'adapter au mieux une densité de probabilité $f(x_i; \theta)$ à un échantillon de données.

$$\hat{\lambda} = \frac{\sum_{i=1}^n x_i}{n}$$

On obtient ainsi l'estimateur du maximum de vraisemblance $\hat{\lambda}$ défini ci-dessus, qui n'est autre que la moyenne empirique.

$E(N) = 0.1324 < V(N) = 0.147$, ce qui montre une surdispersion.

On note une surdispersion dans la régression de Poisson. La variance de Y/X est supérieure à sa moyenne, violant la propriété sous-jacente à l'hypothèse Poisson.

A travers la procédure stepwise sous R, nous proposons le meilleur modèle qui minimise le critère AIC pour une valeur de $AIC=82108.19$ tout en gardant uniquement les variables explicatives qui sont les plus significatives. Les variables sélectionnées pour expliquer la fréquence via un modèle de poisson sont : `nb_sinistres ~ drv_sex1 + pol_coverage + vh_fuel + pol_duration + vh_age + vh_speed + drv_age_lic1 + vh_cyl`.

Dans la figure 1.7, on interprète le graphique de la manière suivante : les points rouges représentent la loi théorique et les histogrammes les fréquences observées, qui sont collés par le sommet à la loi théorique. Tout écart de la base d'un histogramme avec l'axe des abscisses indique donc un mauvais ajustement des observations par la loi théorique. On remarque ainsi que l'ajustement par la loi de

Poisson est peu satisfaisant. En effet les écarts observés sont conséquents et cela quel que soit le nombre d'occurrences.

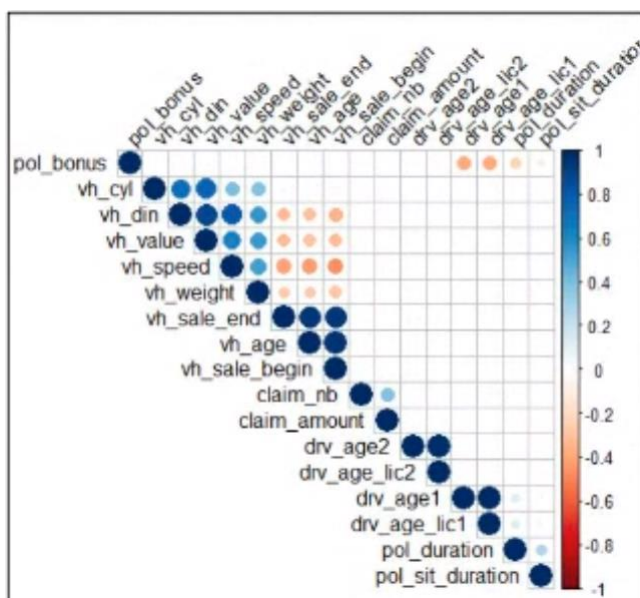
La variance des réponses est, selon la théorie, c'est à dire selon la loi de Poisson, égale à la moyenne des réponses. On dit qu'il y a surdispersion lorsque la variance réelle est supérieure à cette variance théorique. Cela est problématique car dans cette situation, l'erreur standard des paramètres des modèles de régression de Poisson sera sous-estimée. Ceci peut conduire à une p-value excessivement faible, et donc aboutir à une conclusion erronée sur la significativité de la liaison entre les comptages observés et la ou les variables explicatives.

Pour tester la présence d'une éventuelle surdispersion, on peut noter que la surdispersion correspond à une hétérogénéité résiduelle. On utilise simplement `dispersiontest()` de `library(MASS)`. On effectue donc un test de significativité, on testera :

La sortie R nous permet de déterminer le ratio residual deviance / ddl est égal à 57 6330 / 99985 soit 5.764. Ce ratio est très largement supérieur à 1 et permet de mettre en évidence la présence d'une surdispersion. Il est donc nécessaire d'utiliser une autre structure d'erreur dans le modèle de régression. Les principales cause des surdispersion sont :

- Une corrélation variables explicatives,
- L'absence d'une variable explicative importante,
- Une sur-représentation des valeurs zéro par rapport à ce qui est attendu selon la distribution de Poisson de paramètre Lambda.

On va donc étudier la corrélation entre les variables.



On peut voir sur ce graphique qu'il existe une forte corrélation entre les variables : cylindre du véhicule, dynamique du véhicule, valeur du véhicule, et vitesse du véhicule. Ce qui semble logique étant donné qu'ils décrivent la même caractéristique qui est la puissance de la voiture. Plus une voiture est rapide, puissante et à un bon moteur, plus sa valeur est élevée. On observe également un lien entre l'âge du conducteur et l'ancienneté du permis. Ce qui semble évident, en général, les gens passent leurs permis un âge jeune ce qui explique la corrélation entre ces deux variables. Il y a une relation négative entre

le coefficient bonus-malus et les variables « âge du conducteur » et « ancienneté du conducteur dans le portefeuille ». Plus les conducteurs sont âgés moins ils sont risqués et inversement. On peut donc en déduire que la cause de la surdispersion peut être la corrélation entre les variables.

Une autre possibilité est de faire une régression binomiale négative (qui permettra de prendre en compte de la surdispersion). Elle se fait à l'aide de la fonction `glm.nb()` de `library(MASS)`.

La loi Binomiale-Négative est en effet une bonne alternative à la loi de Poisson, en particulier en cas de sur-dispersion des données. En effet, l'utilisation du modèle de Poisson revient à supposer l'égalité entre le nombre moyen de sinistres et la variabilité de ce nombre. Bien souvent, et c'est le cas sur notre jeu d'observation, cette observation n'est pas satisfaite.

2.2 La régression binomiale négative

Le paramètre supplémentaire de la loi Binomiale-Négative nous permet ainsi, par rapport à la loi de Poisson, d'ajuster la variance indépendamment de la moyenne.

On teste aussi l'ajustement à la loi Binomiale-Négative. On dit qu'une variable aléatoire N suit une loi Binomiale-Négative de paramètre

α et q , où $\alpha > 0$ et $0 < q < 1$

Et on a :

$$\Pr[N = k] = \binom{\alpha + k - 1}{k} q^\alpha (1 - q)^k, k \in \mathbb{N}$$

On interprète cette loi à notre situation, on considère les sinistres d'un individu comme une suite d'échecs jusqu'à ce qu'il obtienne un succès, c'est-à-dire un sinistre de coût nul. Le paramètre α est donc égal à 1.

En reprenant la méthode expliquée ci-dessus, on peut déterminer l'estimateur du maximum de vraisemblance pour q . On obtient ainsi :

$$\hat{q} = \frac{\hat{\alpha}}{\hat{\alpha} + k}$$

On estime donc q par la probabilité de succès observée sur l'échantillon. L'étude de cette loi est importante, on constate en effet que la loi de Poisson n'est parfois pas adaptée aux observations de fréquences en assurance, du fait de la présence d'une certaine hétérogénéité dans les observations. A contrario, les mélanges Poissoniens tels que la loi Binomiale-Négative, semblent mieux convenir.

Si on compare maintenant les estimations du modèle de Poisson avec celui de Quasi-Poisson, on remarque qu'au niveau des estimations il n'y a pas de différence :

La fonction `goodfit` de R nous permet d'obtenir la figure 1.8, On remarque tout de suite sur le graphique que l'ajustement est bien meilleur pour la loi Binomiale-Négative que pour la loi de Poisson, les écarts avec l'axe des abscisses sont en effet très faibles.

Par contre, on observe partout une augmentation des valeurs des erreurs standards en passant du Poisson au binomial négatif. Ceci n'est pas du tout surprenant parce qu'en effectuant une régression de Poisson et en présence de sur-dispersion, on va sous-estimer les erreurs standards.

On voit bien que le rapport de la déviance résiduelle et son degré de liberté vaut presque 1 ce qui nous suggère que la loi binomiale négative a bien géré le problème de sur-dispersion.

2.3 Modèles Zéro-Inflatés (ZIP, ZINB)

La sureprésentation de la valeur 0 dans notre portefeuille (88% des assurés) peut amener à l'utilisation du zéro inflatés. Afin de faire face au nombre important de valeurs nulles et à l'hétérogénéité de notre portefeuille, des modèles à inflation de zéros ont été proposés : le modèle de Poisson à inflation de zéros (ZIP) et le modèle binomial négatif à inflation de zéros (ZINB). Ce modèle consiste à combiner deux lois pour la modélisation : (1) loi binomiale pour la survenance ou non de $Y = 0$; (2) loi de Poisson pour le comptage des événements, y compris possiblement la valeur 0. Les modèles ZIP et ZINB comportent donc deux parties : celle relative au modèle de comptage (modèle classique Poisson ou Binomial Négatif) et celle relative à l'inflation de zéros (Logit) qui explique la probabilité de ne pas avoir un sinistre.

La régression ZIP procède d'une modélisation avec la combinaison de 2 lois de distribution : Binomiale et Poisson

Distribution de Y

$$P(Y = y|X, Z) = \begin{cases} \pi + (1 - \pi)e^{-\mu} & \text{si } y = 0 \\ (1 - \pi) \frac{\mu^y e^{-\mu}}{y!} & \text{si } y > 0 \end{cases}$$

Éléments à
modéliser

- π est la probabilité que l'on ait structurellement 0, par conséquent il y a $(1 - \pi)$ de chances d'être dans la situation modélisable par Poisson
- μ est le paramètre de la loi de Poisson

Critère / modèle	Poisson	BN avec thêta estimé	BN avec thêta fixé	ZINB	ZIP
Log vraisemblance		-40822.917		-4.08 e + 0.4	-4.083 e +04
AIC	82108.19	81678	81690	81676.8	81701.5
Déviance	57620	48930	47324	0	0

En ce qui concerne le choix du meilleur des modèles nous allons utiliser le critère d'information d'Akaike, la déviance et la log vraisemblance (test de vuong). Le test de proximité de Vuong est un test basé sur le rapport de vraisemblance pour la sélection de modèle utilisant le critère d'information de Kullback Leibler. Cette statistique fait des déclarations probabilistes sur deux modèles. Ils peuvent être imbriqués, non imbriqués ou se chevauchent. La statistique teste l'hypothèse nulle selon laquelle les deux modèles sont également proches du véritable processus de génération de données, contre l'alternative selon laquelle un modèle est plus proche.

On va également utiliser le critère d'Akaike (AIC), on observe que celui utilisant une loi binomiale négative est beaucoup plus petit (81 690 contre 82108.19 du modèle de Poisson). On peut donc dire que le modèle à loi binomiale négative est probablement le plus adéquat des deux.

La déviance se définit alors comme l'écart de vraisemblance entre le modèle contraint et le modèle non contraint. On compare les deux pour voir si le modèle contraint explique correctement le modèle non contraint. Plus la déviance d'un modèle est proche de 0, plus ce modèle se rapproche du modèle parfait, et donc plus adéquat il est. De même, plus la déviance s'écarte de 0, ins le modèle est bon.

On peut également analyser les résidus, en effet cela permet de diagnostiquer le modèle : identification des points atypiques, recherche des régularités. Le graphe des résidus (ou des résidus réduits) ne doit pas présenter de structure (variance constante sur la verticale et symétrie par rapport aux abscisses). En observant le graphique 2.3 on peut voir que les modèles respectent ces conditions.

On en déduit donc que modèle est le plus adapté est le modèle binomial négatif avec θ fixé.

3. Modélisation pour la sévérité

De la même façon que pour les fréquences de sinistres, on cherche à modéliser le coût moyen de chaque type d'acte par des lois usuelles. Classiquement, on utilise la loi Gamma ou la loi Log-Normale. De même que dans le paragraphe précédent, on estime donc les paramètres de ces lois grâce à la méthode de l'estimateur du maximum de vraisemblance, et l'on compare ensuite la loi théorique à la loi empirique par une représentation graphique.

3.1 Le modèle log normal

On estime dans un premier temps les paramètres de la loi Log-Normale. Une variable est dite Log-Normale si son logarithme suit une loi Normale. Elle présente l'avantage d'être positive, donc adaptée à la modélisation de coût, et permet d'ajuster des phénomènes asymétriques.

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

La méthode utilisée pour déterminer l'estimateur du maximum de vraisemblance, on obtient les estimateurs suivants des paramètres :

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \log(x_i)$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (\log(x_i) - \hat{\mu})^2$$

Dans ce modèle nous observons que le type de carburant utilisé pour une voiture à une forte influence sur le montant des sinistres. Mais cette influence diffère en fonction du carburant utilisés. En effet, on peut voir l'hybride influence négativement le montant du sinistre

contrairement à l'utilisation du gasoil qui a un effet positif sur le montant du sinistre. Cependant nous notons une faible représentation de voiture hybride (0.78%) dans notre portefeuille, ce qui peut représenter un biais dans notre analyse. De plus l'âge du conducteur agit positivement sur le montant du sinistre, ce qui semble intuitif, plus les individus sont âgés plus ils possèdent des véhicules chers ce qui influence le montant du sinistre.

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	6.6953338	0.0275311	243.191	<2e-16	***
vh_fuelGasoline	0.0266606	0.0221229	1.205	0.2282	
vh_fuelHybrid	-0.6070819	0.3085480	-1.968	0.0491	*
drv_age1	0.0014664	0.0007659	1.915	0.0556	.

3.2. Le modèle Gamma

On observe des impacts semblables pour l'âge que le modèle log-normal. On peut exposer ces résultats grâce à la régression :

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	7.3490502	0.0588792	124.816	<2e-16	***
vh_fuelGasoline	0.0966887	0.0473129	2.044	0.041	*
vh_fuelHybrid	-0.8276307	0.6598731	-1.254	0.210	
drv_age1	0.0004753	0.0016380	0.290	0.772	

En étudiant la moyenne de la valeur du véhicule par type d'alimentation, on remarque que le prix moyen des voitures diesel est beaucoup plus élevé que les voitures essence (ce qui est représentatif du parc automobile français). Dans le cas de la loi lognormal et Gamma, les valeurs des coûts prédits sont plus élevées pour les véhicules essence que les véhicules diesel. Cela peut s'interpréter de la manière suivante : les véhicules essences payent pour le risque des véhicules diesels, en ayant une prime plus élevée.

Comment fait précédemment nous allons utiliser les critères d'AIC, la déviance et log-vraisemblance, présentés dans le tableau ci-dessous pour choisir le meilleur modèle :

Critère / modèle	Log Normal	Gamma
Log vraisemblance	-18836	-99258.01
AIC	4178.48	198526
Déviance	16812.18	15453.42

4. Sinistres graves

Il est courant, dans les études de tarification et notamment en tarification automobile, de distinguer les sinistres graves du reste des sinistres pour la modélisation. Cela permet ainsi d'éviter que quelques gros sinistres influencent trop le calcul des coefficients et indicateurs renvoyés par le modèle.

Dans le tableau ci-dessous, les montants des sinistres sont classés par ordre décroissant. On remarque le type de véhicule prédominant est « Tourisme ». L'élément clé, afin de savoir s'il est nécessaire de modéliser les sinistres graves séparément, est d'analyser la proportion

cumulative du montant des sinistres. Ainsi le premier sinistre représente 1.20% de la totalité des pertes de l'assureur. Les trois premiers sinistres représentent 3.84% de la charge total. Il n'est donc pas nécessaire de modéliser les sinistres grave à part, étant donnés de leur poids acceptable dans le portefeuille.

On en déduit donc qu'on n'a pas besoin de prendre en compte ces sinistres graves dans notre analyse.

Montant du sinistre	Nombre de sinistre	Type de véhicule	Proportion cumulative du montant des sinistres
232 104	1	tourisme	1.20
211 112	1	tourisme	1.29
185 065	1	tourisme	3.24
160 086	2	Tourisme	4.06
129 665	1	Tourisme	4.73
50 185	1	tourisme	4.99

5. Prime Pure

Afin d'obtenir notre tarif, il suffit de multiplier nos prédictions de fréquence de sinistres (issue du modèle binomial négatif) et montant de sinistres (issue du modèle log-normal). Le tableau ci-dessous décrit la prime pure moyenne de 227.9 euros, avec un minimum à 101.8 euros et un maximum à 279.9 euros.

	Min	Moyenne	Médiane	Q1	Q3	Max
Tarif	101.8	227.9	229.9	22.4	238.4	279.9

Références :

Mathieu Vautrin, [Mémoire ISUP 2009 Matthieu VAUTRIN \(ressources-actuarielles.net\)](https://ressources-actuarielles.net/ressources/2009/03/01/memoire-isup-2009-matthieu-vautrin/)

Tarification IARD, introduction aux techniques avancées, Frédéric Planchet, Mars 2017 [Présentation PowerPoint \(ressources-actuarielles.net\)](https://ressources-actuarielles.net/ressources/2017/03/01/presentation-iard/)

Arthur Charpentier, Statistique de l'assurance [Statistique de l'assurance \(archives-ouvertes.fr\)](https://archives-ouvertes.fr/theses/2013/03/01/statistique-de-l-assurance/)

Actuariat de l'assurance non vie, Arthur charpentier [slides_ensae_1.pdf \(hypotheses.org\)](https://hypotheses.org/2013/03/01/slides-ensae-1.pdf)

RIAD Meriem, maitre assistante a l'université de tipaza, doctorante à l'ENSSEA, [2515.pdf \(enssea.net\)](https://enssea.net/ressources/2015/03/01/2515.pdf)

Actuariat Introduction, Idris KHARRBOUBI, [Intro-actuariat.pdf \(lpsm.paris\)](https://lpsm.paris/ressources/2015/03/01/intro-actuariat.pdf)

Sandrine BABIN, tarification en assurance emprunteur, CRÉATION DE TABLES DE MORTALITÉ D'EXPÉRIENCE APRÈS SEGMENTATION D'UN PORTEFEUILLE DE PRÊTS PERSONNELS PAR SCORING [link.php \(institutdesactuaires.com\)](https://institutdesactuaires.com/ressources/2015/03/01/link.php)

Julia simaku, [Tarification ANV.pdf](https://anv.fr/ressources/2015/03/01/tarification-anv.pdf)

[Tutoriel : GLM sur données de comptage \(régression de Poisson\) avec R - DellaData](https://delladata.com/ressources/2015/03/01/tutoriel-glm-sur-donnees-de-comptage-regression-de-poisson-avec-r/)

Cannels tamara, analyse du jeu de données de compagne de crises epileptiques traité initialement par Thall et Vail, [Analyse du jeu de données de comptage de crises épileptiques traité initialement par Thall et Vail \(1990\) \(cnrs.fr\)](https://cnrs.fr/ressources/2015/03/01/analyse-du-jeu-de-donnees-de-comptage-de-crisis-epileptiques-traite-initialement-par-thall-et-vail-1990/)

Zero inflated poisson regression, Ricco Rakotomala, Université lumière Lyon 2 [Zero-Inflated Poisson Regression - Régression de Poisson - Modèle de comptage \(univ-lyon2.fr\)](https://univ-lyon2.fr/ressources/2015/03/01/zero-inflated-poisson-regression-regression-de-poisson-modele-de-comptage/)

Mémoire santé aviva, l'assurance santé, [link.php \(institutdesactuaires.com\)](https://institutdesactuaires.com/ressources/2015/03/01/link.php)

rappels de cours et exemples sous R, [R-cours 7 \(univ-mrs.fr\)](https://univ-mrs.fr/ressources/2015/03/01/r-cours-7/)

Annexe :

1. Statistiques descriptives :

Figure 1.1 :



Figure 1.2 :

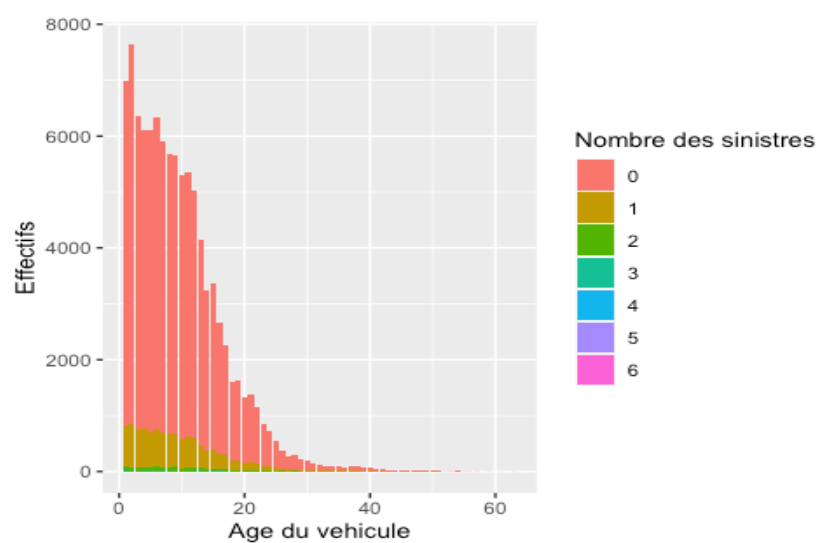


Figure 1.3 :

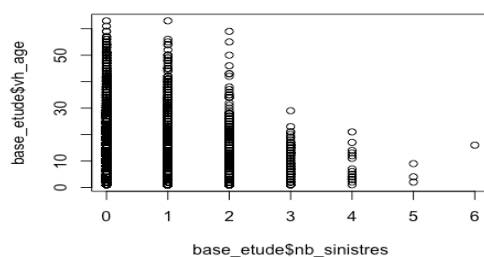


Figure 1.4 :

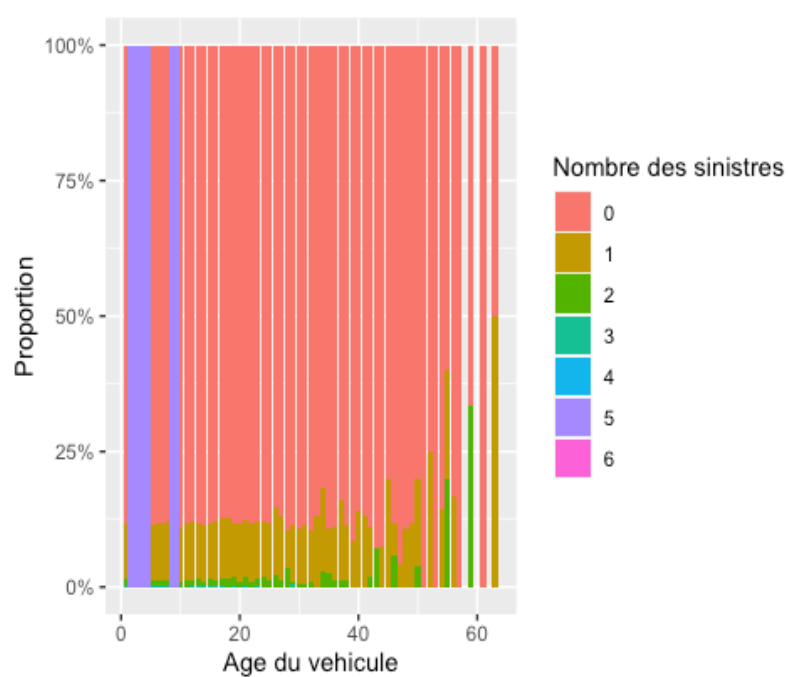


Figure 1.5 :

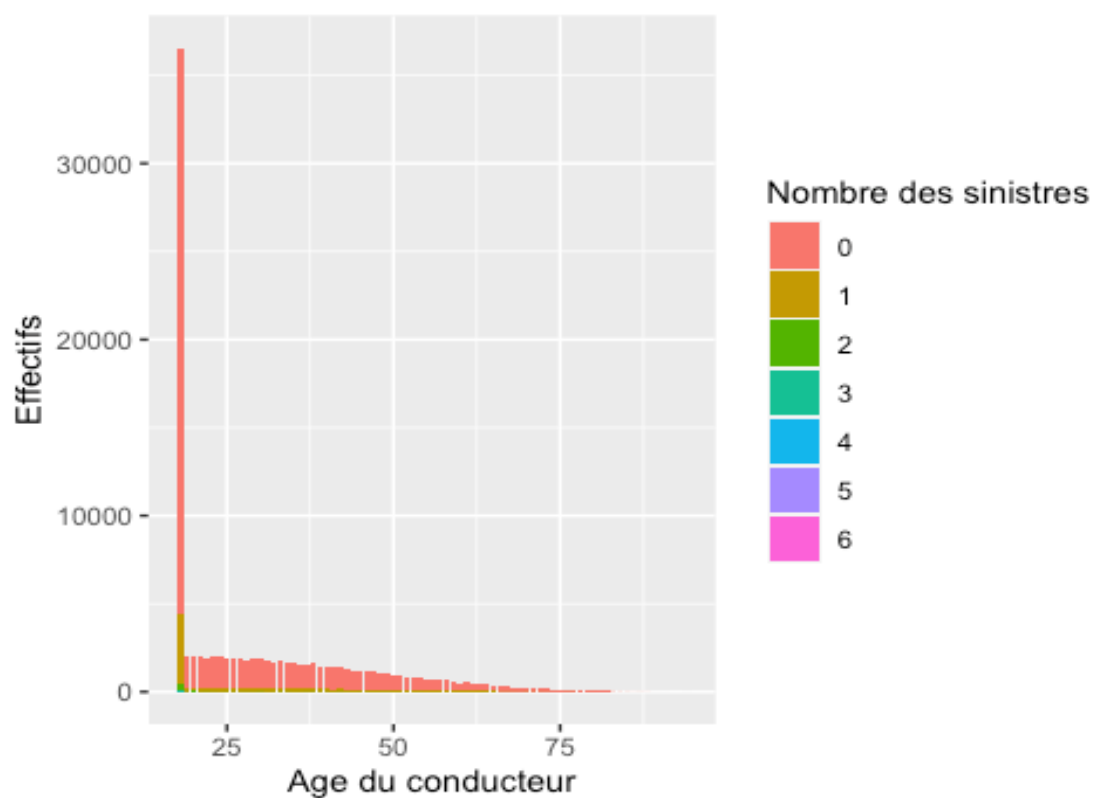
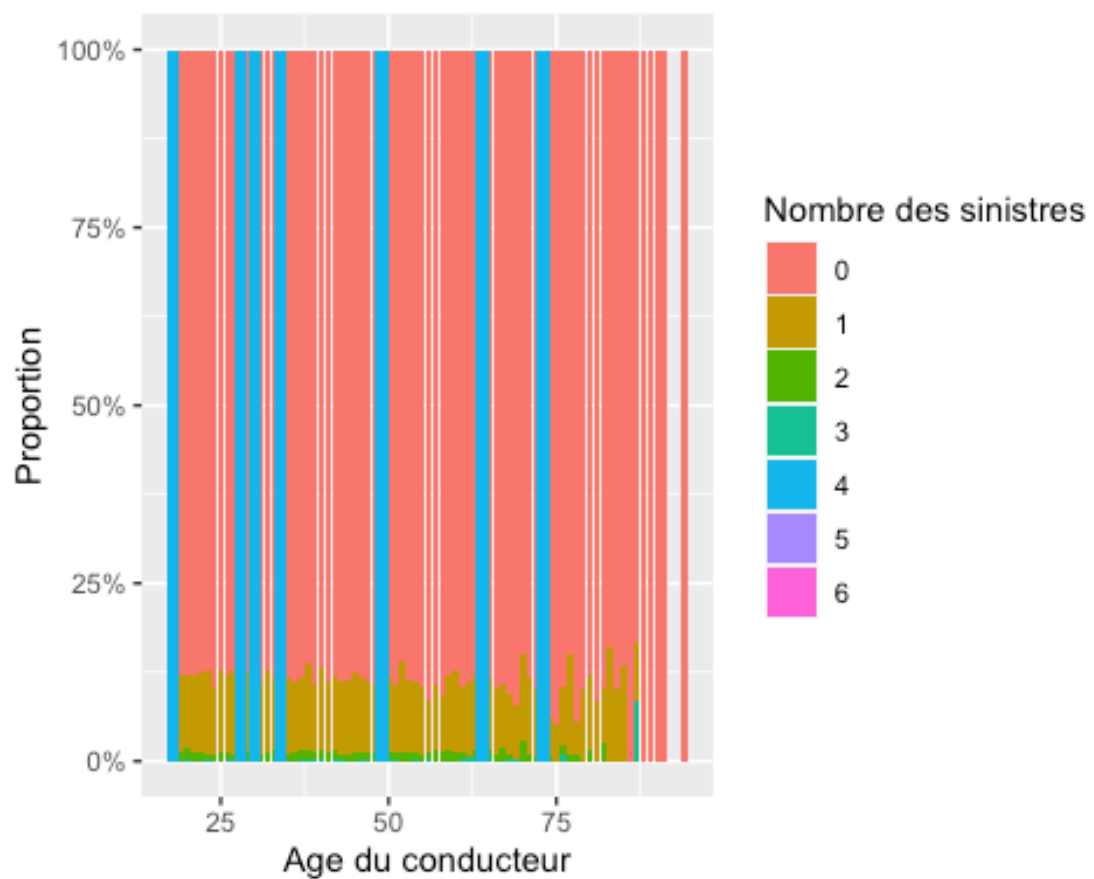


Figure 1.6 :



2. Régression sur la fréquence :

Figure 2.1 :

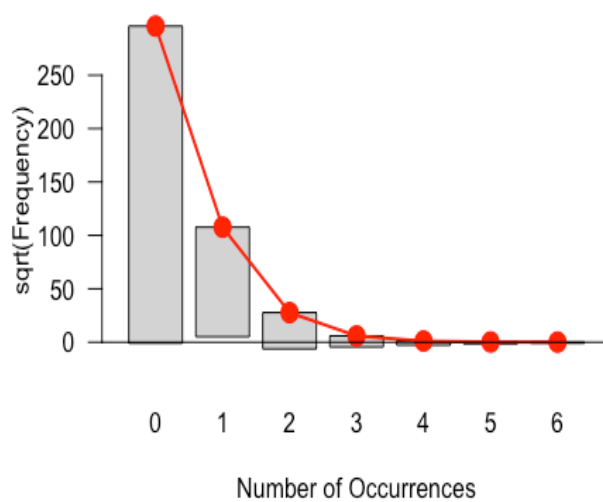


Figure 2.2 :

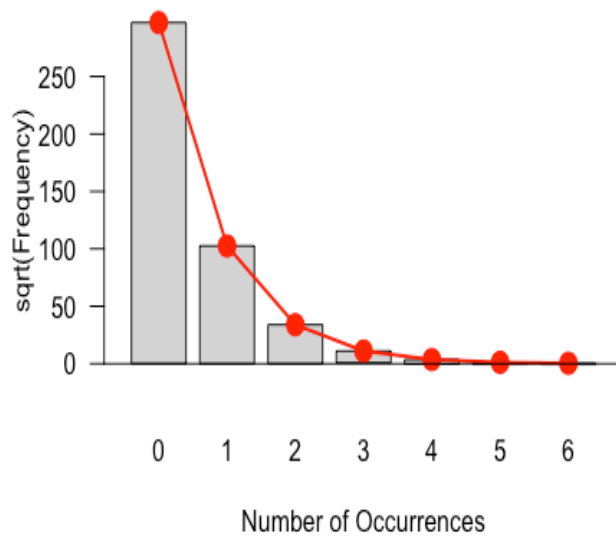
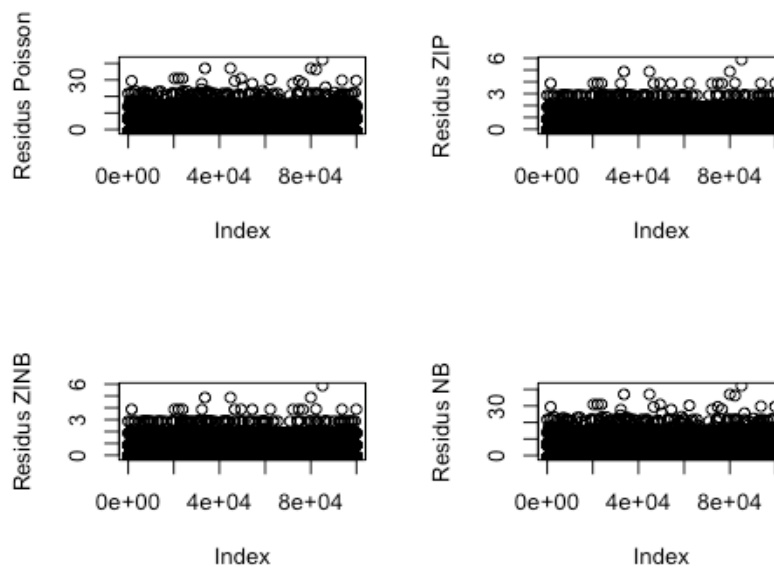


Figure 2.3 : résidus de régression poisson, ZIP, NB, ZINB



3. Régression pour la sévérité

Figure 3.1 : résidus de régression de log normal

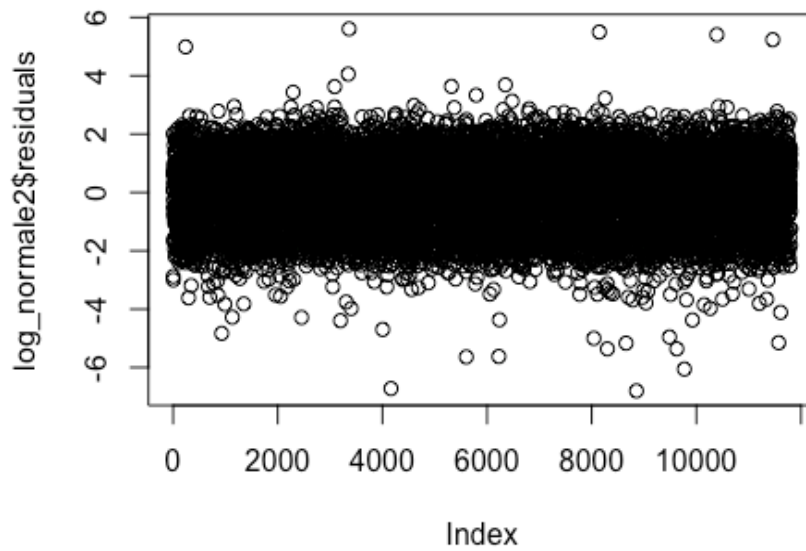


Figure 3.2 : résidus de régressions de gamma

