

# Adaboost for Neural networks applied to an insurance Database

Lucien Ledune

16 april 2019

# **Abstract**

This will be the abstract of the document

# 1 Introduction

Partout dans le monde, la facilité montante de l'accès à une puissance de calcul importante a motivé un grand nombre de changement dans la façon d'approcher les problèmes, et ce dans la plupart des domaines de recherche. Cette avancée technologique a permis d'utiliser des techniques de modélisations et de prédictions qui étaient jusqu'à lors trop longues à mettre en place pour avoir un réel intérêt pratique. Le secteur des assurances n'a pas été épargné par le phénomène et évolue avec son temps. Le machine learning<sup>1</sup>, outil de modélisation très puissant, est maintenant facilement accessible et les compagnies d'assurances souhaitent en tirer le meilleur parti en matière de prédiction et d'aide à la décision.

La modélisation de prédiction n'est pas nouvelle dans ce secteur, mais les algorithmes utilisés changent avec l'essor actuel de la technologie. Ainsi nous pouvons maintenant facilement utiliser des réseaux de neurones afin de répondre aux problèmes pour lesquels des algorithmes moins efficaces étaient utilisés. Le but de ce travail est de montrer que ces nouveaux<sup>2</sup> algorithmes sont efficaces et représentent un espoir d'amélioration pour le futur. TODO : expliquer le mémoire rapidement.

## 2 Dataset

### 2.1 Présentation des données

Afin de réaliser ce travail, nous avons besoin d'une base de donnée adéquate, nous présentons celle-ci dans cette sous-section. La base de données est constituée d'informations sur les clients d'une compagnie d'assurance nommée Wasa, entre 1994 et 1998. Les véhicules assurés sont composés uniquement de motos. Il est difficile d'obtenir des données plus récentes à cause des clauses de confidentialité des assurances. La version originelle de ce jeu de données est utilisée dans une étude cas du livre "Non-life insurance pricing with GLM", écrit par Ohlsson et Johansson, et celle-ci est disponible sur le site web du livre<sup>3</sup>. Le jeu de données est constitué de 64505 observations des 9 variables suivantes :

- OwnersAge : L'âge du conducteur.
- Gender : Le sexe du conducteur.
- Zone : Variable catégorielle représentant la zone dans laquelle le véhicule est conduit..
- Class : Variable catégorielle représentant la classe du véhicule. Les classes sont assignées dans une des 7 catégories selon le ratio :  $EV = \frac{kW \times 100}{kg + 75}$ .
- VehiculeAge : L'âge du véhicule en années.

---

<sup>1</sup>Apprentissage automatique : Méthodes statistique permettant à un ordinateur d'apprendre à résoudre un problème donné à l'aide d'une base de données pertinente.

<sup>2</sup>La plupart de ces algorithmes sont en fait assez peu récents mais étaient difficilement applicables en raison d'une puissance de calcul trop faible, citons par exemple le perceptron, élément de base des réseaux de neurones qui a été inventé dès 1957.

<sup>3</sup><http://staff.math.su.se/esbj/GLMbook/case.html>

- **BonusClass** : Le bonus du conducteur, un nouveau conducteur commence à 1 et sera incrémenté à chaque année complète passée dans la compagnie sans sinistre déclaré, jusqu'à un maximum de 7.
- **Duration** : Le nombre d'année passées dans la compagnie.
- **NumberClaims** : Le nombre de sinistres.
- **ClaimCost** : Le coût des sinistres.

La variable **Zone** est décrite dans la table 1.

Table 1: Descriptions des variables catégorielles		
Variable	Classe	Description
Zone géographique	1	Parties centrales et semi-centrales des trois plus grandes villes de Norvège.
	2	Banlieues et villes moyennes.
	3	Petites villes (à l'exception de celles des catégories 5 et 7).
	4	Villages (à l'exception de ceux des catégories 5 et 7).
	5	Villes du nord de la Suède.
	6	Campagnes du nord de la Suède.
	7	Gotland (Grande île).

## 2.2 Analyse exploratoire

Maintenant que les différentes variables ont été brièvement présentées et leur fonction plus claire, nous allons maintenant passer à l'analyse exploratoire de celles-ci. Le but de cette analyse est de mieux comprendre les données qui serviront à entraîner les différents algorithmes, ainsi que de repérer d'éventuelles anomalies. Durant l'analyse exploratoire d'une base de données, il est important de regarder la distribution des variables, celle-ci nous donne beaucoup d'informations quant aux données.

Sur ces deux premières figures nous observons respectivement les distributions des variables **Gender** et **OwnersAge**. La première chose que nous pouvons observer est la grande disparité entre le nombre d'hommes et de femmes clients de l'assurance. Notre jeu de données est composé d'hommes pour la grande majorité. Pour ce qui est de la variable **OwnersAge**, nous constatons que les valeurs sont réparties entre 16 et 92 ans, avec deux "pics" vers 25 et 45 ans. Les valeurs maximales et minimales de nos variables continues seront importantes pour la suite, car elles sont nécessaires afin d'appliquer une normalisation des données, qui sera discutée plus tard dans ce travail. Ci-dessous les distributions des variables **VehiculeAge** et **Zone** :

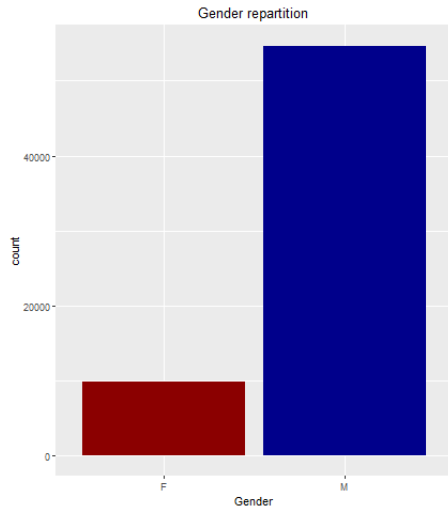


Figure 1: Distribution Gender

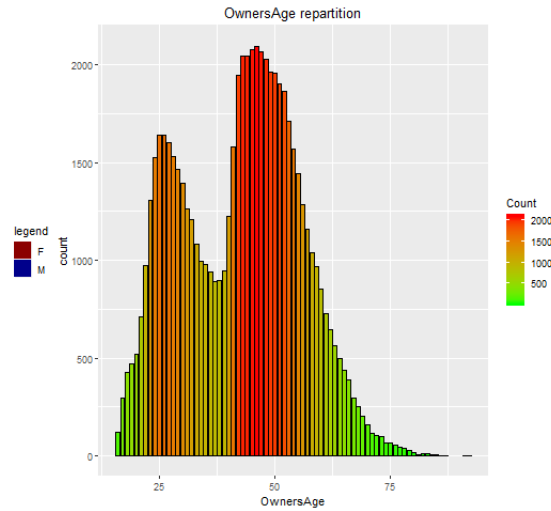


Figure 2: Distribution OwnersAge

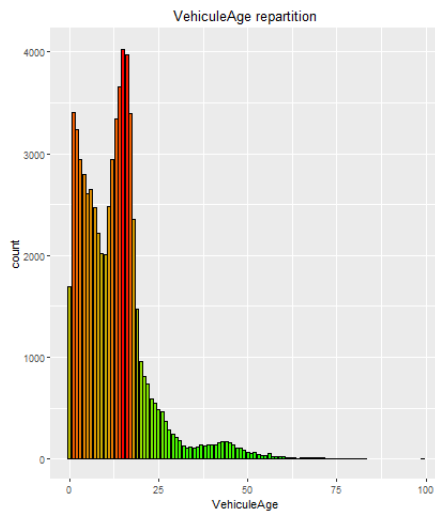


Figure 3: Distribution VehicleAge

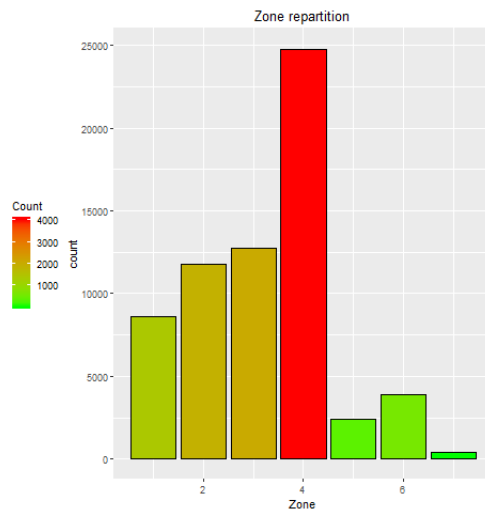


Figure 4: Distribution Zone

La distribution de vehicule est indiquée en années, et nous pouvons donc observer que si la plupart des véhicules assurés ont moins de 20 ans, un grand nombre de ceux-ci sont bien plus vieux, avec comme maximum 99 ans. Il est possible que ce véhicule soit considéré comme outlier et il sera reconsidéré dans la partie preprocessing.

La distribution de la variable Zone est intéressante, elle nous révèle que la plupart des véhicules assurés sont conduits dans des villages, mais aussi que très peu d'entre eux le sont dans le Gotland. Ceci n'est pas surprenant puisque la population de la Suède est d'environ 10 millions d'habitants, pour seulement

60.000 habitants la région du Gotland. Les classes 5 et 6 sont elles aussi minoritaires, cela était aussi à prévoir puisque ces catégories représentent le nord de la Suède alors que la plupart de la population vis dans le sud du pays.

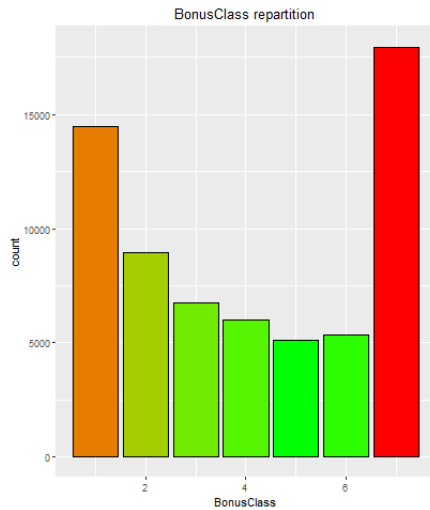


Figure 5: Distribution BonusClass

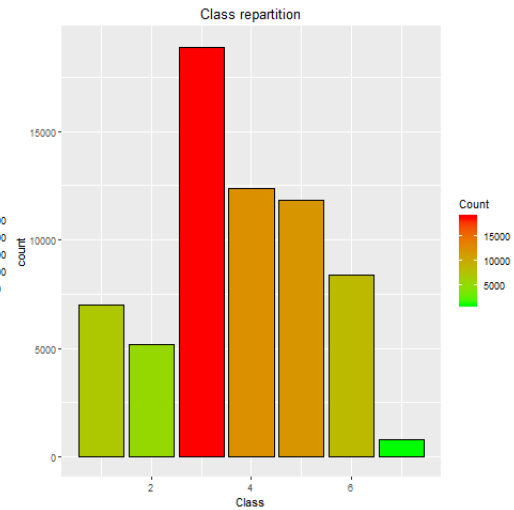


Figure 6: Distribution Class

Les classes de bonus les plus fréquentes sont les classes 1 et 7, respectivement le minimum (entrée dans la compagnie d'assurance) et le maximum (client fidèle depuis sept années au minimum).

Pour la variable Class nous observons qu'assez peu de véhicules appartiennent à la catégorie la plus puissante. En fait, la plupart des véhicules se situent dans les classes 3, 4 et 5, ce qui montre que les véhicules les moins puissants et les plus puissants sont minoritaires. Les graphiques suivants nous montrent la répartition des sinistres par bonus et par classe de véhicule.

Il paraît logique de supposer que plus un client a un bonus élevé, moins celui-ci causera d'accident, puisque le bonus monte uniquement si le client parvient à compléter une année sans causer d'accident. Cependant en observant la figure 7, il apparaît que la majorité des cas d'accidents sont déclarés par des clients appartenant aux classes 3 à 6 de bonus. Ce qui est d'autant plus étonnant lorsque l'on associe ce résultat avec la distribution de la variable Class (figure 5) : les classes 3 à 6 sont celles contenant le moins d'utilisateurs. Les clients appartenant à la classe de bonus 7 semblent cependant causer très peu d'accidents malgré le fait qu'ils soient la classe de bonus majoritaire.

Les résultats de la figure 8 sont moins surprenant : plus un véhicule est puissant, plus le risque d'accident sera important. Les déviations de cette règle par les classes 3 et 7 sont expliquées par la distribution de la population à travers les différentes classes (figure 6), ainsi il y a peu d'accidents pour les véhicules de classe 7 simplement car ceux-ci sont peu nombreux, une conclusion similaire peut être énoncée pour la classe 3.

La répartition des sinistres par sexe indique que les hommes sont plus suscep-

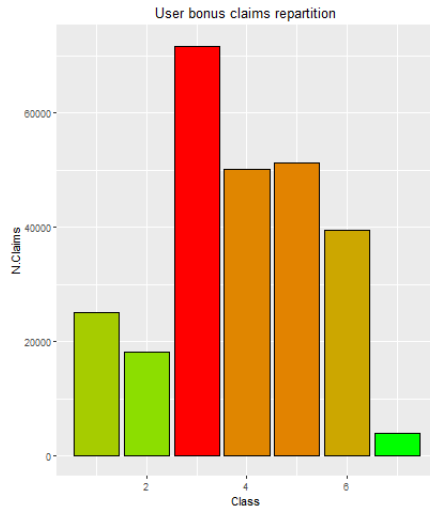


Figure 7: Bonus Claims

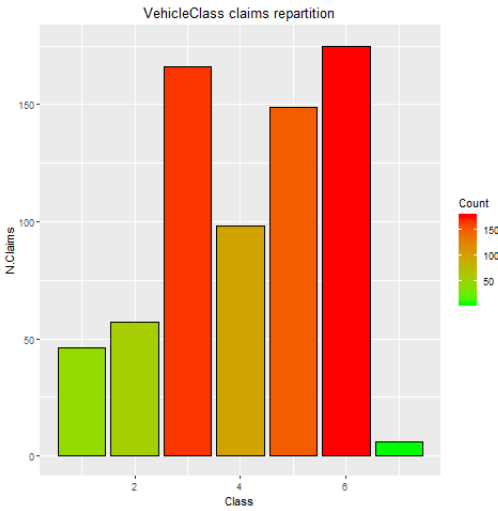


Figure 8: Class Claims

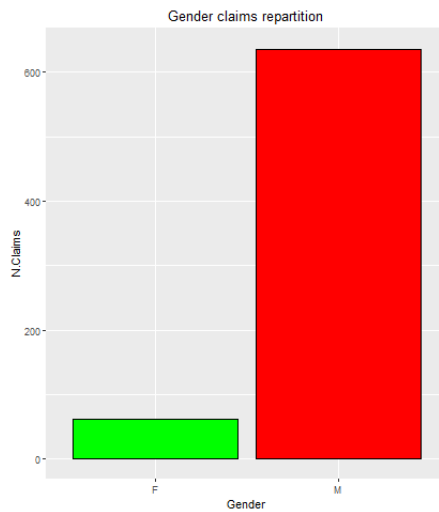


Figure 9: Gender claims

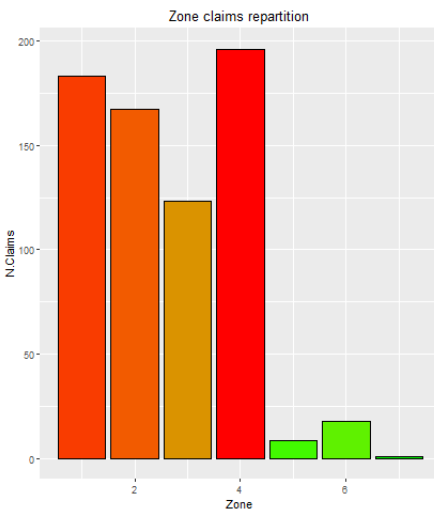


Figure 10: Zone claims

tibles de causer des accidents que les femmes. Il faut prendre en compte que les hommes sont bien plus nombreux que les femmes dans nos données, mais la conclusion ne change pas puisque la proportion de femmes est de 18%, alors que celles-ci ne causent que 8.7% des sinistres. Sur la figure 10, nous observons que les classes 1 et 3 sont celles qui déclarent le plus de sinistres. La première classe étant plutôt minoritaire (figure 4) nous observons que les personnes habitant dans les parties centrales et semi-centrales de Norvège semblent causer bien plus d'accidents que les autres catégories, la même conclusion peut être faite pour la deuxième classe (Banlieues et villes moyennes). Pour ce qui est de nord de la Suède, nous constatons l'inverse puisque ceux-ci semblent causer assez peu d'accidents.

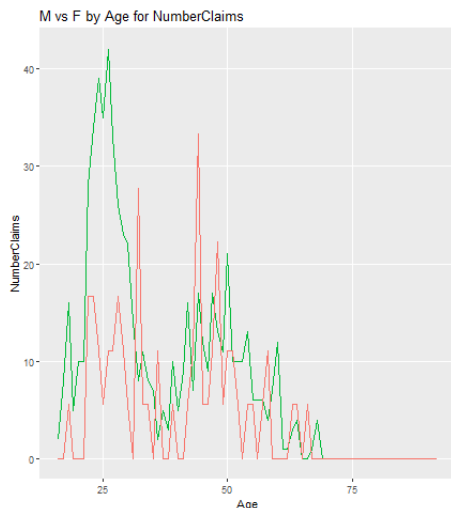


Figure 11: Gender claims by Age

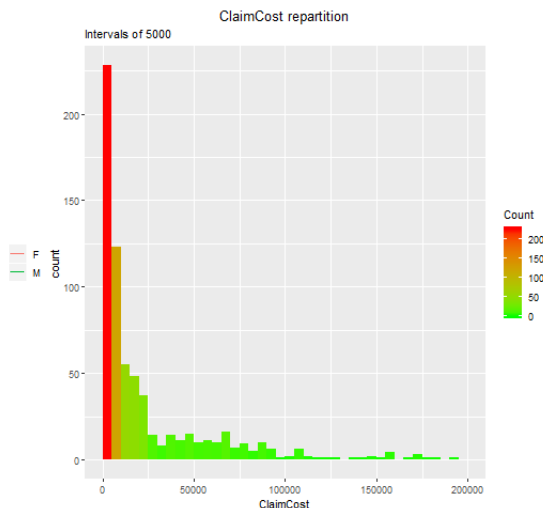


Figure 12: Distribution ClaimCost

La figure 11 a pour but de mieux comprendre cette disparité entre les hommes et les femmes dans la déclaration des sinistres. Une comparaison par âge permet de mettre en évidence une information importante. En effet, si les hommes semblent bien causer plus d'accidents que les femmes jusqu'à la trentaine, cette tendance s'égalise passé ce cap. Enfin, la figure 12 nous informe de la distribution de la variable ClaimCost, représentant le coût total des sinistres d'un client. Il est important de noter qu'il s'agit bien de la somme de tous les sinistres déclarés du client, ainsi si une personne a reçu une compensation de l'assurance pour plusieurs sinistres, la variable contient la somme du coût de tous ces sinistres déclarés. Il semble que la grande majorité des clients ne dépassent pas 5000€ de compensation venant de l'assurance (Le graphique ne montre que les valeurs non-nulles, ainsi les clients n'ayant jamais déclaré de sinistres ne sont pas comptés dans l'intervalle  $[0, 5000]$ ). La distribution de cette variable montre bien qu'assez peu de clients dépassent les 20.000€ de compensation, mais certains peuvent monter à près de 200.000€.

## 2.3 Assurances

Dans cette sous-section seront expliquées les particularités de l'analyse de données dans le cadre de l'assurance, et plus précisément dans le cas qui nous intéresse ici.

### 2.3.1 Principe

Un contrat entre l'assureur et l'assuré implique le paiement d'une prime d'assurance par le second, en échange de compensations lorsqu'un sinistre est déclaré. Le fait de regrouper un grand nombre de clients va avoir un effet stabilisateur de variance pour l'assureur. En effet il est difficile de prédire précisément le nombre d'accidents qu'un assuré aura, cependant plus le nombre d'assurés sera élevé



et plus le nombre (total) de sinistres sera prévisible. Ceci est du à la loi des grands nombres qui implique :

$$\bar{X} \rightarrow \mu \text{ when } n \rightarrow \infty$$

La moyenne de l'échantillon  $\bar{X}$  converge vers le moyenne de la population  $\mu$  lorsque la taille de l'échantillon augmente.

Dans le cadre de l'assurance, les variables explicatives sont aussi appelées facteur d'évaluations, les deux termes seront utilisés indifféremment au cours de ce travail.

### 2.3.2 Ratio

Il est important de comprendre ce que l'on cherche à expliquer avec un modèle, aussi il faut savoir que cette analyse se base sur un ratio. Un ratio  $Y$  est un rapport entre une réponse  $X$  et une exposition<sup>4</sup>  $v$ .

$$Y = X/v$$

Il existe différents ratios mais ici ne sera présenté que celui utilisé dans l'analyse à venir : La fréquence des réclamations. Il s'agit du rapport entre le nombre de réclamations et la durée du contrat, nous obtenons ainsi une fréquence qui sera la variable prédite par nos modèles.

### 2.3.3 Distribution du ratio

Pour la création d'un modèle cherchant à prédire la fréquence des réclamations c'est souvent la distribution Poisson qui est utilisée. Une loi poisson de paramètre  $\lambda$  implique :

$$Y \sim \text{Poisson}(\lambda)$$

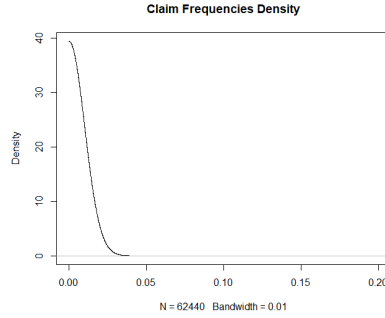
$$E[Y] = \lambda$$

$$\text{Var}[Y] = \lambda$$

Afin de vérifier que la distribution de la fréquence des sinistres suit bien une loi de poisson celle-ci est affichée. En effet la plupart des valeurs se situent près de 0, les valeurs supérieures à 1 sont présentes mais extrêmement rares et ne sont donc pas affichées afin d'observer plus précisément les valeurs types.

---

<sup>4</sup>L'exposition peut être le nombre de sinistres ou la durée du contrat.



La loi de poisson fait partie de la famille des modèles à dispersion exponentielle (ED), une catégorie regroupant un grand nombre de lois statistiques connues comme la loi normale, gamma ou encore binomiale. Une variable aléatoire  $Y$  suit une dispersion exponentielle si sa fonction de distribution suit la forme suivante :

$$f_{Y_i}(y; \theta_i; \phi) = \exp \left\{ \frac{y\theta_i - a(\theta_i)}{\phi/v_i} \right\} c(y, \phi, v_i)$$

- $\theta_i$ : paramètre dépendant de  $i$ .
- $\phi$ : paramètre identique pour tous les  $i$ . C'est le paramètre de dispersion
- $c(\cdot)$  est indépendante du paramètre  $\theta_i$ .

## 2.4 Preprocessing

Le preprocessing des données est une étape très importante lorsque nous travaillons avec ce type de données et celui-ci n'est pas à négliger. Il consiste à préparer les données pour les algorithmes, afin que ceux-ci puissent fonctionner de manière optimale. Dans ce travail, plusieurs méthodes de preprocessing ont été utilisées.

### 2.4.1 Variables binaires (Dummy variables)

Certaines de nos variables sont catégorielles, et plus particulièrement sont des variables non-ordinales. Cela signifie que les différentes catégories ne peuvent pas être ordonnées, il s'agit purement de l'assignation d'un client à un groupe donné. Ce type de variable ne peut pas être encodé tel quel et doit être transformé en une série de variables binaires. La conversion d'une variable catégorielle en variable binaire consiste à créer  $n - 1$  nouvelles variables dites "binaires" et qui prendront pour valeur 1 si le client appartient à cette catégorie, 0 sinon. Par exemple si  $x_{ij}$  représente la variable associée au client  $i$  et à la variable  $j$ , nous pouvons écrire :

$$x_{ij} = \begin{cases} 0 & \text{if } x_i \notin j \\ 1 & \text{if } x_i \in j \end{cases}$$

Les variables Gender, Zone, Class, et BonusClass peuvent être transformées

de la sorte. Ce travail sera effectué avec le package `caret`<sup>5</sup> de R. La fonction incluse dans ce package crée  $n$  variables binaires, la dernière sera donc supprimée car celle-ci ne représente aucune information. Prenons pour exemple la variable `Gender`. Si  $x_{iGender} = 1$  alors le client est une femme. Cependant nous ne devons pas créer une deuxième variable binaire puisque si  $x_{iGender} = 0$  nous pouvons en déduire que le client est un homme. Ceci explique pourquoi nous n'avons besoin que de  $n - 1$  variables binaires pour représenter toute l'information de  $n$  catégories appartenant à une variable catégorielle.

### 2.4.2 Normalisation

Les données quantitatives ne doivent pas être transformées en variables binaires (bien que celles-ci peuvent être converties en intervalles puis en variables binaires, ce ne sera pas le cas dans le cadre de ce travail). Cependant cela ne signifie pas qu'aucune méthode de preprocessing ne doit être appliquée à ces données. La normalisation des données quantitatives aide certains algorithmes à converger plus rapidement et peut même parfois améliorer leur précision. Cette normalisation des données est particulièrement importante pour les réseaux de neurones, car l'optimisation de ceux-ci est basée sur les résultats des fonctions d'activation<sup>6</sup> utilisées. Celles-ci varient rapidement entre 0 et 1 (ou -1 et 1) et sans la normalisation des données leurs réponses aux valeurs extrêmes seraient localisées dans les extrémités, ce qui peut empêcher l'algorithme de converger dans les cas les plus graves.

Plusieurs méthodes peuvent être utilisées afin de procéder à la normalisation des données. Cependant nous ne décrirons que celle qui sera utilisée dans le cadre de ce travail : la normalisation Min Max.

Elle consiste (pour une variable  $j$ ) à convertir l'ensemble des données sur l'intervalle  $[0,1]$  en se servant de  $Max(x_j)$ <sup>7</sup> et  $Min(x_j)$ <sup>8</sup> par la formule :

$$z_{ij} = \frac{x_{ij} - \min(x_j)}{\max(x_j) - \min(x_j)}$$

Cette transformation sera appliquée aux variables `OwnersAge` et `VehiculeAge`.

### 2.4.3 Outliers

Lors de l'encodage des données, il arrive que certaines valeurs des variables quantitatives soient assez extrêmes et semblent fort distantes des autres observations. Dans ce cas nous pouvons soit garder la valeur si celle-ci semble ajouter de l'information au modèle ou la supprimer dans le but d'une meilleure généralisation par celui-ci. Dans la section 2.2 une donnée a particulièrement attiré notre attention comme outlier potentiel : une moto de 99 ans. Pour celle-ci nous pouvons entrevoir deux possibilités :

- Il peut simplement s'agir d'une erreur d'encodage, auquel cas l'observation doit être supprimée pour éviter de biaiser le modèle.

---

<sup>5</sup>Classification And REgression Tools : Package de machine learning pour R.

<sup>6</sup>Voir chapitre 3

<sup>7</sup>Valeur maximale prise par la variable  $j$ .

<sup>8</sup>Valeur minimale prise par la variable  $j$ .

- Il peut également s'agir d'une vraie moto de type "ancêtre", cependant ce genre de véhicule ne se conduit pas de la même manière qu'un véhicule utilisé au quotidien, et on peut supposer que les risques d'accidents sont bien moindres. Dans ce cas-ci la suppression de la donnée n'est pas obligatoire mais peut aider à une meilleure généralisation du modèle.

Puisqu'aucun moyen simple de vérifier cette hypothèse n'est disponible, celle-ci sera supprimée de la base de donnée car la perte d'information est moindre.

La suppression de trop d'information peut cependant nuire au modèle, les données extrêmes mais non-aberantes seront donc conservées.

#### 2.4.4 Durée de contrat

Pour le bon fonctionnement du modèle, certaines autres données ont du être modifiées car celles-ci entraînaient des erreurs dans la fonction de déviance utilisée pour la calibration des poids des réseaux de neurones. Les observations incriminées sont celles ayant une durée de contrat nulle. Ces valeurs nulles de durée de contrat sont peu cohérentes (impossible d'avoir un sinistre déclaré au cours du contrat si celui-ci n'est pas encore en vigueur) et elles seront donc changées par la valeur correspondant à un jour ( $1/365 \simeq 0.00274$ ). Certaines de ces observations n'ont cependant pas une valeur nulle pour la variable NbClaims, ce qui signifie que des sinistres auraient été déclarés dans des contrats dont la durée est nulle. Ce cas de figure n'étant pas possible, elles seront supprimées car il s'agit probablement d'erreur d'encodage. De plus changer la valeur de la durée du contrat par 0.00274 dans ce cas reviendrait à considérer une fréquence de sinistre de  $\simeq 365$  soit un sinistre par jour ce qui pourrait grandement biaiser le modèle.

## 3 Algorithmes

Dans cette section, les différents algorithmes utilisés lors des analyses seront présentés à des fins de compréhension du travail. ....

Afin d'effectuer un travail de prédiction en matière de régression ou de classification, il est nécessaire d'avoir un outil puissant à disposition : l'apprentissage automatique ou "machine learning". Il s'agit d'un ensemble de méthodes basées sur des outils statistiques permettant la création de modèles pouvant être utilisés pour la prédiction. (réf wiki)

### 3.1 Explication d'un modèle

#### 3.1.1 Définitions et apprentissage

Avant de pouvoir rentrer dans les détails, il est important de définir les notations qui seront utilisées, de définir ce qu'est un modèle et comment celui-ci se construit. Dans notre cas seuls les modèles supervisés<sup>9</sup> seront abordés. Deux éléments sont nécessaires pour la création d'un modèle, une base de données et un algorithme supervisé (choix du modèle). La base de données est composée de  $N$  entrées des variables explicatives et de la variable expliquée  $y_i$  où  $i$  représente l'individu et où  $x_{ij}$  représente la valeur prise pour la variable  $j$  pour

---

<sup>9</sup> Algorithmes nécessitant une base de données pour l'apprentissage.

l'individu  $i$ . Le modèle va se servir des valeurs prises par les  $x_{ij}$  pour tenter d'expliquer  $y_i$ . Pour ce faire l'algorithme doit subir une phase d'apprentissage, au cours de laquelle celui-ci va recevoir des entrées provenant de la base de données et se servir de celles-ci pour "apprendre" et se calibrer sur le problème que l'on cherche à résoudre. Cette phase d'apprentissage est différente d'un algorithme à un autre. Une fois cette phase d'apprentissage terminée le modèle ainsi créé peut être utilisé pour la prédiction (après que celui-ci ait été validé, voir sous-sections suivantes). Les valeurs prédites par l'algorithmes seront notées  $\hat{y}_i$ .

### 3.1.2 Évaluation

Pour l'apprentissage, les algorithmes supervisés se basent sur un critère mesurant l'erreur d'un modèle : la fonction de perte. Ce critère permet la comparaison entre différents modèles et un choix optimal parmi ceux-ci, en sélectionnant celui qui minimise l'erreur commise dans les prédictions. Il existe un grand nombre de critères pouvant servir pour l'évaluation des modèles mais le choix de celui-ci doit se faire en fonction du problème adressé et de la base de données. Le choix de ce critère est important car il peut impacter l'efficacité des modèles créés.

Le critère peut être défini comme  $f(y, \hat{y})$  une fonction dépendant des valeurs prises par  $y$  et  $\hat{y}$  mesurant la différence entre les prédictions et la réalité d'une certaine manière. ([https://en.wikipedia.org/wiki/Loss\\_function](https://en.wikipedia.org/wiki/Loss_function)) Une fonction de perte intuitive est donnée en guise d'exemple par la racine de l'erreur quadratique moyenne (RMSE).

$$RMSE(\theta, \hat{\theta}) = \sqrt{MSE(\theta, \hat{\theta})} = \sqrt{E((\hat{\theta} - \theta)^2)}$$

Le RMSE est sûrement un des critères les plus utilisés aujourd'hui (y compris pour les réseaux de neurones), cependant utiliser celui-ci reviendrait à considérer le ratio de la fréquence des réclamations comme étant distribué selon une loi gaussienne. Les résultats s'en trouvent alors biaisés car les propriétés statistiques de la variable expliquée ne sont pas prises en compte. Afin de palier à ce problème, un autre critère statistique sera utilisé : la déviance de Poisson. La déviance est utilisée la plupart du temps pour des tests d'hypothèses et elle joue un rôle important dans l'analyse des distributions de la famille exponentielle. Elle peut être vue comme une forme de distance entre deux modèles (d'où son utilisation en temps que critère). Notons  $D^*$  la déviance mise à l'échelle, un test entre les maximum de vraisemblance du modèle entraîné et du modèle saturé<sup>10</sup>. Notons aussi  $l(\hat{y}_i)$  le log-vraisemblance de  $\hat{y}_i$  et  $l(y_i)$  le log-vraisemblance de  $y_i$ . (wiki deviance)

$$D^*(y_i, \hat{y}_i) = 2 \left( l(y_i) - l(\hat{y}_i) \right)$$

En utilisant la fonction de distribution des ED vue précédemment, le résultat suivant est obtenu :

<sup>10</sup>Modèle où les  $\hat{y}_i$  sont définis comme les véritables valeurs  $y_i$ , le maximum de vraisemblance de ce modèle est le meilleur que nous pouvons obtenir du fait que les  $y_i$  sont parfaitement prédits.

- $a'(\theta) = e^\theta$
- $h(y) = \ln(y)$

$$D^*(y_i, \hat{y}_i) = \frac{2}{\phi} v_i(y_i h(y_i)) - a(h(y_i)) - y_i h(\hat{y}_i) + a(h(\hat{y}_i))$$

La déviance non mise à l'échelle peut ensuite être obtenue en multipliant cette formule par le paramètre de dispersion  $\phi$  (voir section 2.3.3). La dérivation de cette déviance non mise à l'échelle appliquée à la loi de poisson résulte en la déviance de Poisson non-mise à l'échelle, c'est ce critère qui sera utilisé pour l'apprentissage de nos modèles.

$$D(y_i, \hat{y}_i) = 2v_i(y_i \ln(y_i) - y_i \ln(\hat{y}_i) - y_i + \hat{y}_i), \quad \forall y_i > 0$$

La base de donnée utilisée portant sur des événements rares, beaucoup de cas  $y_i = 0$  seront étudiés, en modifiant la formule ci dessus la formule adéquate pour ce cas particulier est donnée par :

$$2v_i \hat{y}_i, \quad \forall y_i = 0$$

### 3.1.3 Overfit, AIC et BIC

## 3.2 Réseaux de neurones