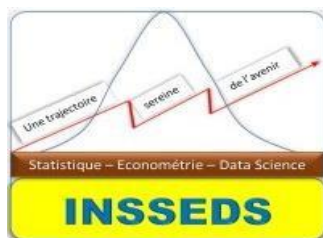


MINISTRE DE L'ENSEIGNEMENT
SUPERIEUR ET DE RECHERCHE
SCIENTIFIQUE



Institut Supérieur de
Statistique D'Econométrie

REPUBLIQUE DE COTE
D'IVOIRE



Union-Discipline-Travail

MASTER 1

STATISTIQUES – ECONOMETRIE – DATA SCIENCE

MINI PROJET

ANALYSE STATISTIQUES ECONOMETRIQUES

MODÉLISATION DES ACCIDENTS AUTOMOBILES

Nom: YOBO

Prénom(s): BAYE GUY ANGE HENOC

Enseignant – Encadreur

AKPOSSO DIDIER MARTIAL



Accident

Avant-propos

Le présent projet a pour objectif d'analyser un jeu de données issu du secteur de l'assurance IARD (assurance de dommages), plus précisément dans le domaine de l'assurance automobile. Le jeu de données étudié, intitulé `assurance_auto_makani.csv`, contient les informations de 374 393 clients d'une société d'assurance. Ces données permettent d'examiner différents profils d'assurés ainsi que les sinistres déclarés, dans le but de mieux comprendre les facteurs associés aux accidents automobiles.



Avant-propos	3
INTRODUCTION	6
PARTIE I : ANALYSE DESCRIPTIVE ET EXPLORATOIRE DES DONNEES	7
A) APPROCHE METHODOLOGIQUE DES DONNEES	7
A.1) Information sur le jeu de donnée	7
A.2) Détection et apurement des valeurs Manquantes et Aberrantes / extrême	7
B) STATISTIQUES DESCRIPTIVES (PARAMETRES)	9
B.1) Graphiques des quantitatives	9
C) Relation entre les variables et la Variable cible « ACCIDENT »	11
C.1) distribution des variables qualitatives	11
C.2) Distribution des variables quantitatives	12
PARTIE II : MODELISATION DE L'ACCIDENT PAR RÉGRESSION LOGISTIQUE	13
1) RÉGRESSION LOGISTIQUE	13
2) CRÉATION DU MODÈLE DE RÉGRESSION LOGISTIQUE	14
3) CALCUL DES ODD RATIO ET DES EFFETS MARGINAUX	15
3.1) Ordre ratio (ODD RATIO)	15
4) CALCULER LE TMC : taux de mauvais classement et Matrice de Confusion	17
5) INDICATEURS DE PERFORMANCES	17
6) PREDICTION DE LA PROBABILITE DE FAIRE UN ACCIDENT	19
PARTIE III : MODELISATION POISSONNIENE DE "frequence"	20
A) ☞ Test d'adéquation à la loi de Poisson	21
A.1) Estimation des paramètres par la méthode du maximum de vraisemblance	21
A.2) Tests d'adéquation	21
B) Construction du modèle BINOMIALE NEGATIVE	22
B.1) Spécification du modèle	22
B.2) ♦ AIC et ♦ BIC	23
B.3) Vérification des hypothèses du modèle	24
DASHBOARD	25
✅ CONCLUSION	26

LISTE DES FIGURES

Figure 1 extrait du jeu de donnée	7
Figure 2 Valeur manquantes	8
Figure 3 Valeurs abberentes.....	8
Figure 4 Valeurs abberente winzorisé	9
Figure 5 variable quali	12
Figure 6 variable quanti	12
Figure 7 Frequence.....	20
Figure 8 distribution des frequences	20
Figure 9 Residus	24



INTRODUCTION

Dans le secteur de l'assurance automobile, la compréhension fine du risque associé à chaque assuré est essentielle pour adapter les tarifs, prévenir la sinistralité et optimiser la gestion du portefeuille client. Le jeu de données **assurance_auto_makani.csv**, utilisé dans cette étude, regroupe les informations de **374 393 assurés** d'une compagnie d'assurance IARD, incluant les caractéristiques personnelles des clients, celles de leurs véhicules, ainsi que les circonstances des éventuels accidents. L'objectif principal de ce projet est de modéliser le risque d'accident automobile à partir des variables disponibles, à travers deux axes : (1) estimer la **probabilité de survenance** d'un accident selon le profil de l'assuré, et (2) prédire le **nombre d'accidents** pour chaque assuré. Cette double approche vise à améliorer la segmentation des assurés selon leur niveau de risque, en vue d'optimiser les stratégies tarifaires et préventives.

Les résultats attendus sont : une modélisation fiable de la variable accident (oui/non) pour détecter les facteurs explicatifs majeurs, une estimation du nombre d'accidents via une régression adaptée aux données de comptage, une meilleure compréhension des variables influençant le risque automobile (liées au véhicule, au comportement de l'assuré et aux conditions de circulation), ainsi que des recommandations concrètes sur l'ajustement des primes d'assurance en fonction des profils de risque.

La méthodologie suivie se divise en deux grandes étapes. D'abord, un **prétraitement des données** comprenant le nettoyage, la gestion des valeurs manquantes, la transformation des variables (notamment la **binarisation** de la variable « fréquence » en une variable cible « accident »), l'encodage des variables catégorielles et, si besoin, la normalisation. Ensuite, deux **analyses statistiques** sont menées : une **régression logistique** pour modéliser la probabilité de survenance d'un accident, et une **régression de Poisson** pour modéliser le nombre d'accidents. La qualité de ces modèles sera évaluée à l'aide d'indicateurs appropriés (AUC pour la logistique, déviance ou RMSE pour la Poisson), afin d'assurer la robustesse et l'interprétabilité des résultats.

PARTIE I : ANALYSE DESCRIPTIVE ET EXPLORATOIRE DES DONNEES.

A) APPROCHE METHODOLOGIQUE DES DONNEES

L'approche méthodologique des données englobe l'organisation, la collecte, l'analyse et l'interprétation des données dans le cadre d'une étude ou d'une recherche. Elle repose sur un ensemble de principes, de techniques et de processus visant à traiter les données de manière systématique et rigoureuse, afin d'obtenir des résultats fiables et pertinents.

A.1) Information sur le jeu de donnée

Âge	Sexe	Âge Véhicule	Marques	Couleur	Carburant	Sièges	Portes	Année Fabric.	Transmission	Accident	Gravité	Trajet	Lumière	Météo	Route	Manceuvre	Fréquence	Prime annuelle
22	0<5	NISSAN	gris	Gazoil	21	4	2007	man	Oui	3	1..
22	(>5	NISSAN	argent	Gazoil	21	4	2008	man	Oui	2	1...
22	1>5	NISSAN	argent	Gazoil	21	4	2000	man	Oui	3	3..
22	1<5	NISSAN	gris	Gazoil	21	4	2010	man	Oui	3	3...
22	0<5	NISSAN	bleu	Gazoil	21	4	2006	man	Oui	3	1..
22	1<5	NISSAN	gris	Gazoil	21	4	2007	man	Oui	3	4...
47	0>5	TOYOTA	noir	Gazoil	21	4	2007	man	Non	0	0..
24	(<5	TOYOTA	noir	Gazoil	21	4	2010	man	Oui	2	4...
27	0<5	TOYOTA	noir	Gazoil	21	4	2010	man	Oui	3	4..
45	(<5	MERCEDES	noir	Gazoil	21	4	2010	man	Non	0	0...
28	1<5	SUZUKI	noir	Gazoil	21	4	2007	man	Oui	3	4..
29	1<5	SUZUKI	noir	Gazoil	21	4	2007	man	Oui	3	4...

Figure 1 extrait du jeu de donnée

A.2) Détection et apurement des valeurs Manquantes et Aberrantes / extrême

Dans cette section, nous allons chercher à identifier visuellement les éventuelles valeurs manquantes dans notre jeu de données, puis à les traiter. Ces valeurs peuvent provenir d'erreurs de mesure, de saisie, de calcul ou bien de valeurs extrêmes réelles présentes dans les données. Les valeurs atypiques peuvent avoir un impact majeur sur les résultats des analyses statistiques, en faussant des indicateurs comme la moyenne ou l'écart-type, mais aussi en influençant les tests d'hypothèse. Il est donc crucial de détecter et de traiter ces valeurs extrêmes avant de procéder à toute analyse statistique.

A.2.1) Visualisation des valeurs manquantes manquantes



Figure 2 Valeur manquantes

Notre jeu de donnée ne présente aucunes valeurs manquantes

A.2.2) Visualisation des valeurs aberrantes

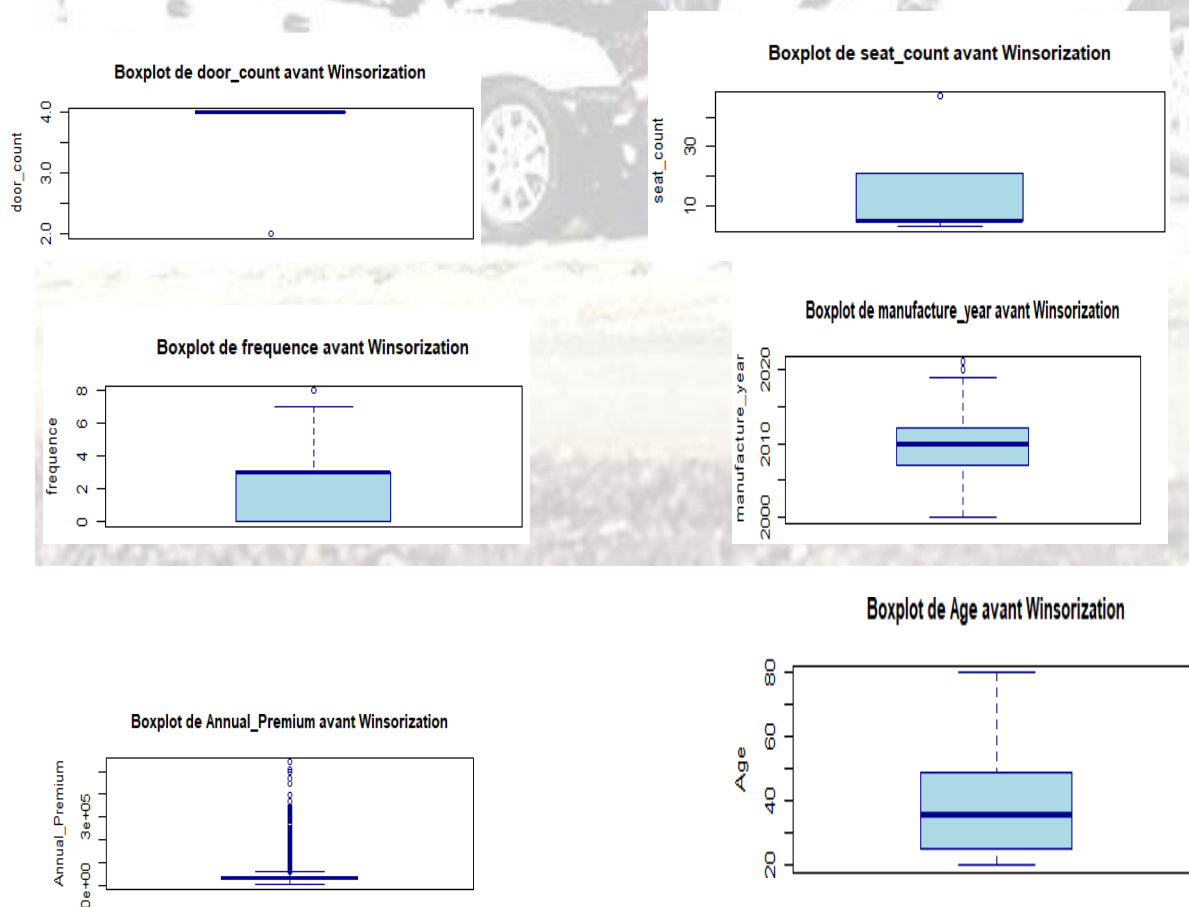
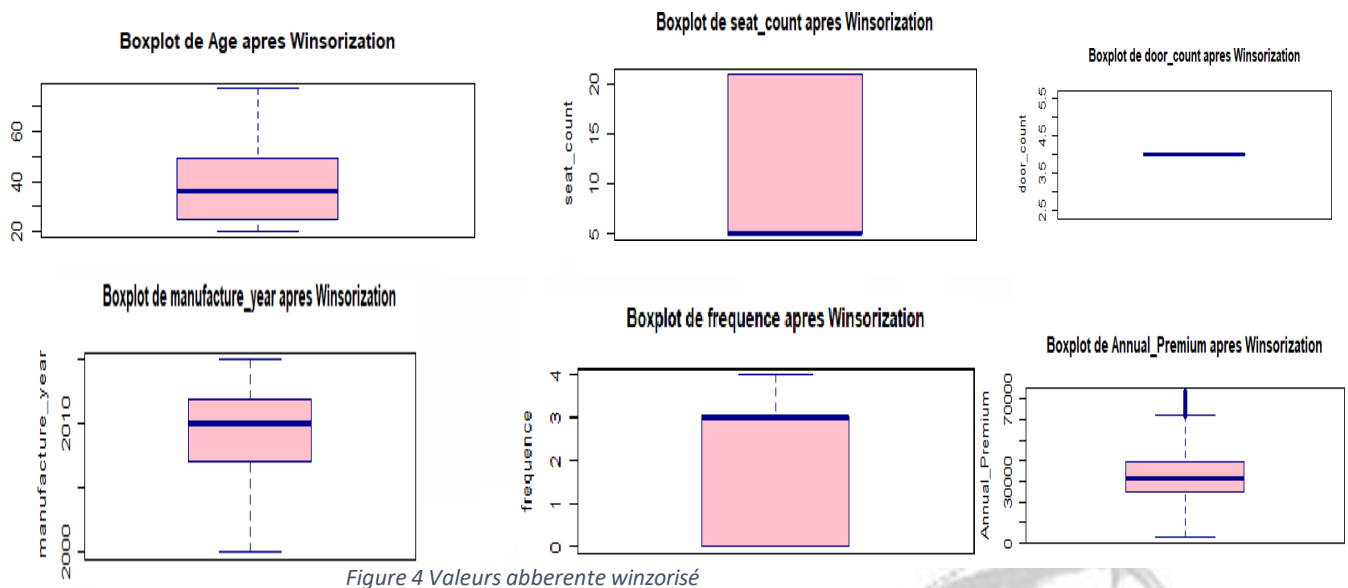


Figure 3 Valeurs aberrantes

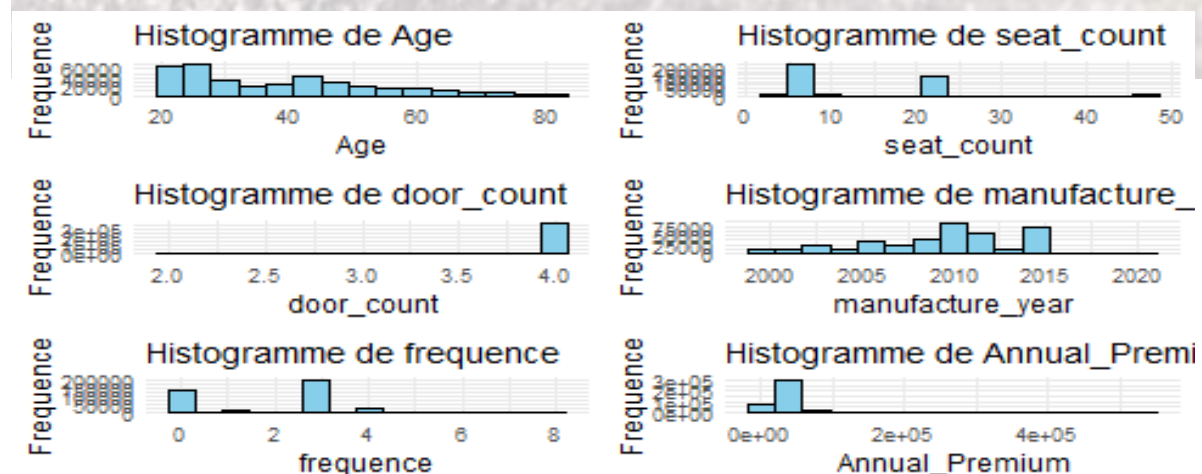
A.2.3) Visualisation des valeurs aberrantes winzorisées



B) STATISTIQUES DESCRIPTIVES (PARAMETRES)

	mean	sd	min	max	median	IQR
Age	38.49	15.17	20	77	36	24.00
seat_count	11.05	7.71	5	21	5	16.00
door_count	4.00	0.00	4	4	4	0.00
manufacture_year	2009.31	4.13	2000	2015	2010	5.00
frequence	1.91	1.52	0	4	3	3.00
Annual_Premium	30486.02	15786.23	2630	72852	31692	14901.75

B.1) Graphiques des quantitatives



◆ Age

- **Distribution fortement asymétrique à droite (positive)** : la majorité des assurés sont jeunes (entre 20 et 40 ans).
- Très peu d'assurés ont plus de 60 ans.

◆ seat_count (Nombre de sièges)

- La distribution est **fortement concentrée autour de 21 sièges**, ce qui est surprenant et pourrait indiquer une erreur ou une codification particulière (ex. : véhicule collectif ou erreur d'unité).
- Ce point mérite une **vérification dans les données sources**.

◆ door_count (Nombre de portes)

- La **majorité des véhicules ont 4 portes**.
- Quelques rares cas de véhicules à 2 ou 3 portes.

◆ manufacture_year (Année de fabrication)

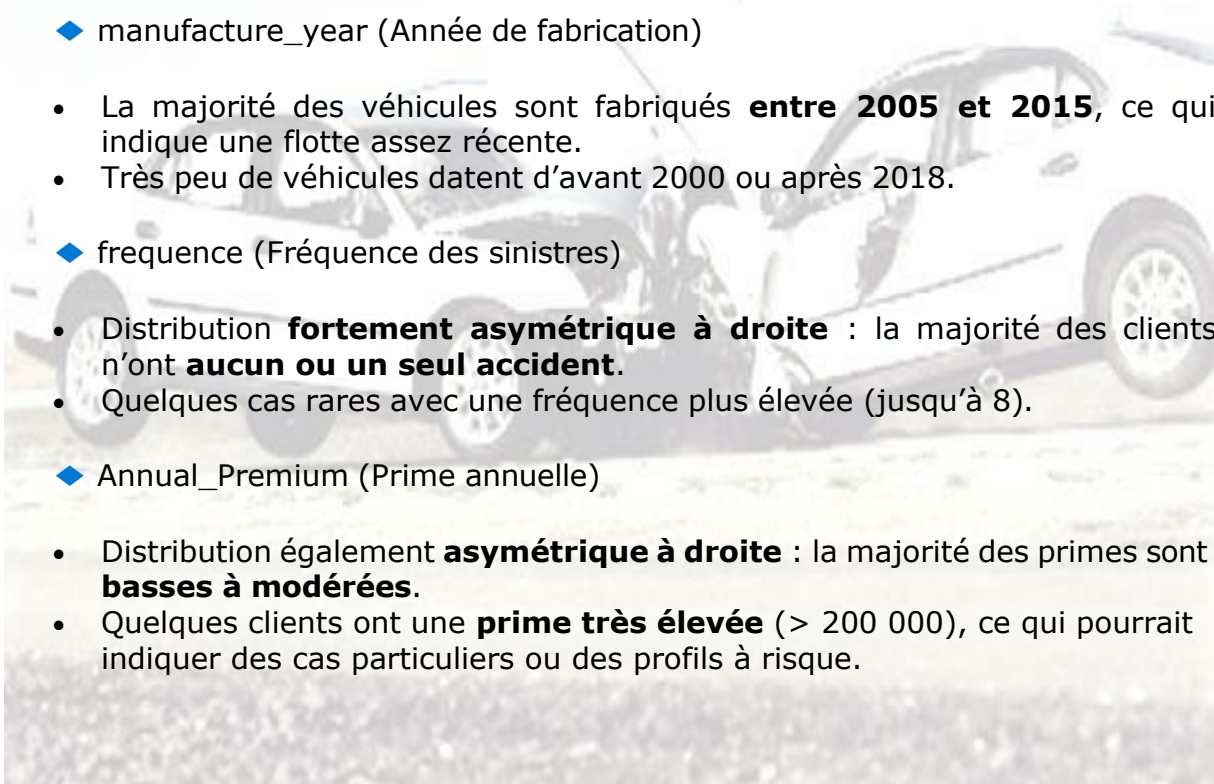
- La majorité des véhicules sont fabriqués **entre 2005 et 2015**, ce qui indique une flotte assez récente.
- Très peu de véhicules datent d'avant 2000 ou après 2018.

◆ frequence (Fréquence des sinistres)

- Distribution **fortement asymétrique à droite** : la majorité des clients n'ont **aucun ou un seul accident**.
- Quelques cas rares avec une fréquence plus élevée (jusqu'à 8).

◆ Annual_Premium (Prime annuelle)

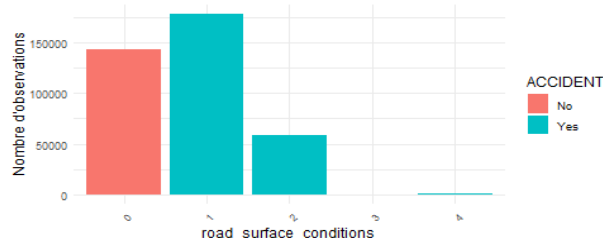
- Distribution également **asymétrique à droite** : la majorité des primes sont **basses à modérées**.
- Quelques clients ont une **prime très élevée** (> 200 000), ce qui pourrait indiquer des cas particuliers ou des profils à risque.



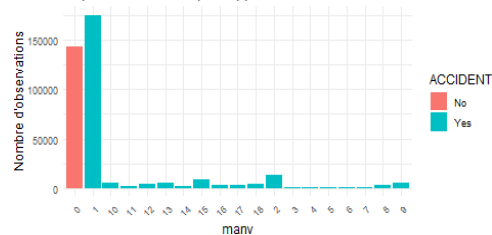
C) Relation entre les variables et la Variable cible « ACCIDENT »

C.1) distribution des variables qualitatives

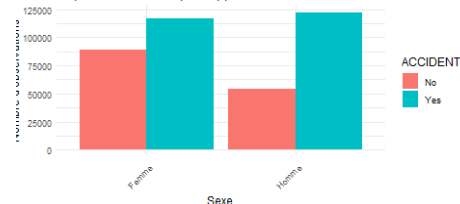
Repartition de road_surface_conditions par rapport a ACCIDENT



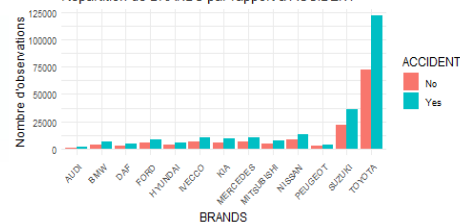
Repartition de manv par rapport a ACCIDENT



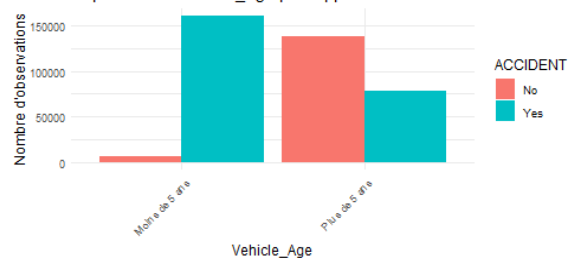
Repartition de Sexe par rapport a ACCIDENT



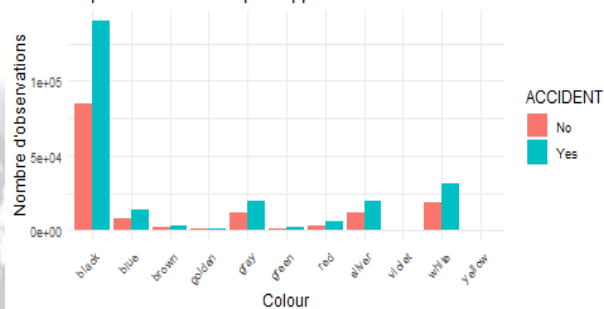
Repartition de BRANDS par rapport a ACCIDENT



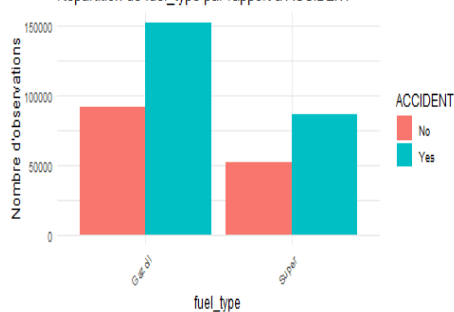
Repartition de Vehicle_Age par rapport a ACCIDENT



Repartition de Colour par rapport a ACCIDENT



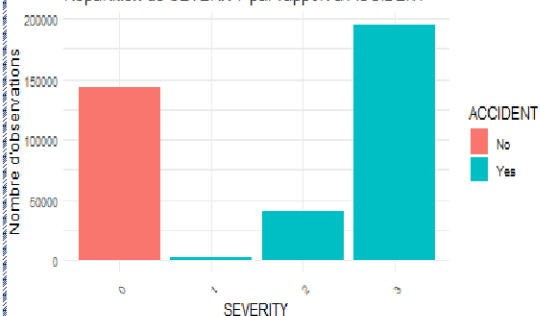
Repartition de fuel_type par rapport a ACCIDENT



Repartition de transmission par rapport a ACCIDENT



Repartition de SEVERITY par rapport a ACCIDENT



Repartition de trajet par rapport a ACCIDENT

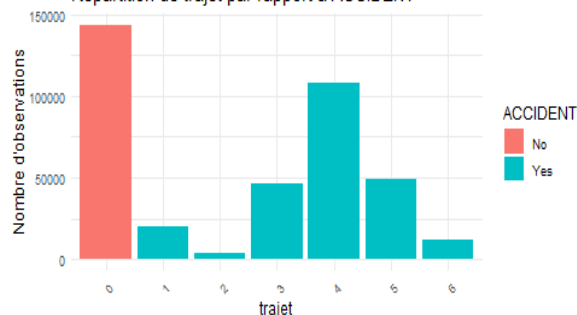
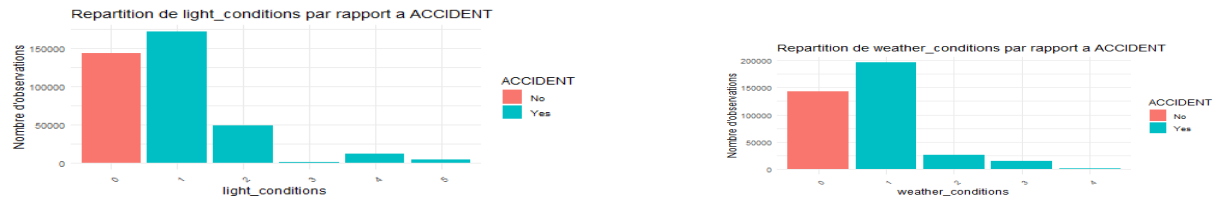


Figure 5 variable quali



C.2) Distribution des variables quantitatives

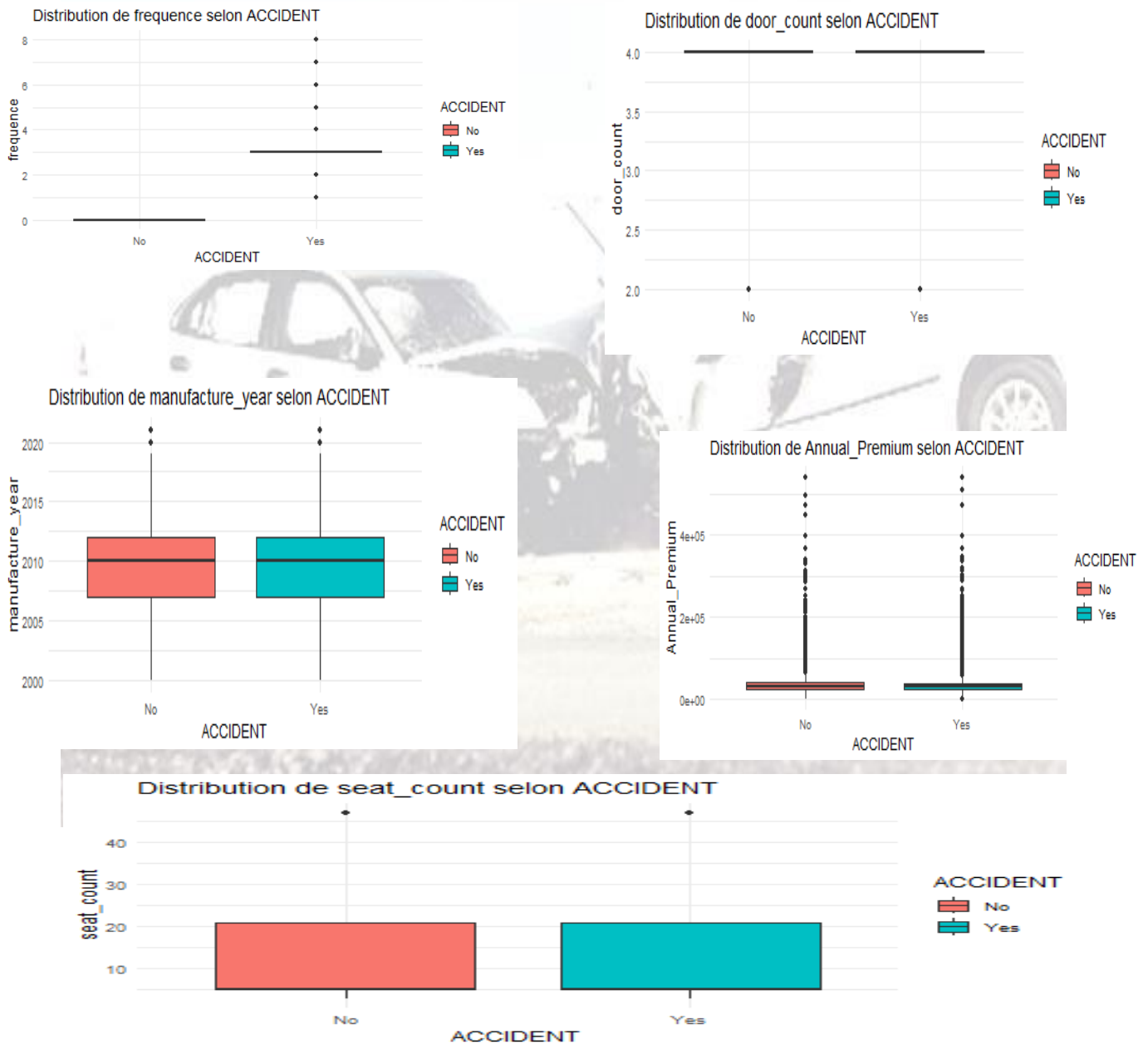


Figure 6 variable quanti

PARTIE II : MODELISATION DE L'ACCIDENT PAR RÉGRESSION LOGISTIQUE

1) RÉGRESSION LOGISTIQUE

La **régression logistique** est une méthode statistique utilisée pour modéliser une **variable dépendante binaire** (c'est-à-dire qui prend deux valeurs possibles, comme 0 ou 1, oui ou non, succès ou échec) en fonction d'une ou plusieurs **variables indépendantes**.

📌 Principe de base :

Contrairement à la **régression linéaire** qui prédit des valeurs continues, la **régression logistique** prédit la **probabilité** qu'un événement se produise.

Par exemple :

Est-ce qu'un email est un spam ou non ?
 Est-ce qu'un patient est malade ou pas ?
 Est-ce qu'un client fera un accident ou non ?

📌 Formule mathématique

La régression logistique utilise la **fonction sigmoïde (logistique)** pour transformer une valeur réelle en une probabilité entre 0 et 1 :

$$P(Y = 1 | X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \dots + \beta_n X_n)}}$$

- $P(Y = 1 | X)$: probabilité que la sortie soit 1
- β_0 : l'ordonnée à l'origine (intercept)
- β_i : coefficients associés aux variables explicatives X_i

F En pratique :

- **Entrée** : données avec une colonne "cible" (0 ou 1) et des variables explicatives
- **Sortie** : probabilité d'appartenir à la classe 1

2) CRÉATION DU MODÈLE DE RÉGRESSION LOGISTIQUE

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	2.7925215	0.0147438	189.403	<2e-16 ***
Age	-0.0033663	0.0003701	-9.096	<2e-16 ***
SexeHomme	0.1557593	0.0084317	18.473	<2e-16 ***
Vehicle_AgePlus de 5 ans	-3.2880390	0.0144896	-226.924	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Interprétation intelligente et synthétique

Ce modèle nous dit : quels profils sont plus ou moins susceptibles d'avoir un accident, toutes choses égales par ailleurs.

3 p Intercept : 2.79 - Quand toutes les variables explicatives sont à zéro (ce qui veut dire "âge = 0", sexe = "Femme", véhicule de moins de 5 ans), la log-odds d'accident est ****2.79****. - En pratique, l'intercept seul n'est pas très interprétable ici, mais il sert de point de départ au modèle.

$\#$ Age : -0.00337 (très significatif, $p \leq 2e-16$)** - Chaque ****année supplémentaire**** réduit ****l'odds (la cote)**** d'avoir un accident d'environ ****0.34%**** :

$$e^{-0.00337} \approx 0.9966$$

- Autrement dit, ****les conducteurs plus âgés sont légèrement moins susceptibles d'avoir un accident****, selon le modèle.

- Effet modéré mais statistiquement très fort.

● 3. ****SexeHomme : +0.156**** - Les hommes ont une ****odds d'accident environ 17% plus élevée**** que les femmes :

$$e^{0.156} \approx 1.17$$

- Autrement dit, ****à âge et ancienneté de véhicule constants****, les hommes sont ****statistiquement plus à risque**** que les femmes.

4. ****Vehicle_AgePlus de 5 ans : -3.29**** les véhicules de plus de 5 ans ont une ****odds d'accident presque 96% plus faible**** que ceux de moins de 5 ans :


$$e^{-3.29} \approx 0.037$$

NB : Ce résultat est ****contre-intuitif**** ! Il mérite probablement une investigation

- Est-ce un biais dans les données ?
- Les véhicules plus anciens sont-ils conduits plus prudemment (propriétaires plus âgés, trajets courts) ?
- Variable mal codée ?

Qualité du modèle

- ****Null deviance**** : 510,987 → modèle sans prédicteurs
- ****Residual deviance**** : 353,611 → modèle avec prédicteurs
- → Donc, ****le modèle améliore bien l'ajustement****, et l'AIC (353,619) est relativement bas.

 **Résumé** Le modèle logistique révèle que les hommes ont un risque d'accident significativement plus élevé que les femmes, et que ce risque diminue légèrement avec l'âge du conducteur. Contre toute attente, les véhicules de plus de 5 ans sont associés à un risque beaucoup plus faible d'accident, un résultat qui mérite une attention particulière quant à l'interprétation (biais ou effet indirect possible). Le modèle est statistiquement très significatif et explique une part substantielle de la variabilité.

3) CALCUL DES ODD RATIO ET DES EFFETS MARGINAUX

3.1) Ordre ratio (ODD RATIO)

Estimation des paramètres

PARAMETRES			
(Intercept)	Age	SexeHomme	Vehicle_AgePlus de 5 ans
2.792521454	-0.003366282	0.155759306	-3.288039006

Variable	Coefficient	Effet sur la probabilité d'accident	Interprétation
(Intercept)	2.7925	Base (quand toutes les variables = 0)	Log-cote initiale relativement élevée de survenue d'un accident
Âge	-0.0034	Diminue légèrement	Plus l'assuré est âgé, moins il a de probabilité d'avoir un accident
Sexe : Homme	0.1558	Augmente légèrement	Les hommes ont une probabilité légèrement plus élevée d'avoir un accident
Âge du véhicule > 5 ans	-3.2880	Diminue fortement	Les véhicules plus anciens sont associés à une probabilité beaucoup plus faible

Intervalles de confiance

confint(model)	2.5 %	97.5 %
(Intercept)	2.763661275	2.821456842
Age	-0.004091805	-0.002641122
SexeHomme	0.139231665	0.172283426
vehicle_AgePlus de 5 ans	-3.316475540	-3.259676443

Interpretation :

Variable	Coefficient	Intervalle de confiance (95 %)	Interprétation
(Intercept)	2.7925	[2.7637 ; 2.8215]	L'ordonnée à l'origine est bien estimée ; l'intervalle est étroit, donc bonne précision
Âge	-0.0034	[-0.0041 ; -0.0026]	Effet négatif significatif : l'âge diminue la proba d'accident, et l'intervalle ne contient pas 0
Sexe : Homme	0.1558	[0.1392 ; 0.1723]	Effet positif significatif : être un homme augmente la proba d'accident, intervalle exclut 0
Âge du véhicule > 5 ans	-3.2880	[-3.3165 ; -3.2597]	Effet fortement négatif : les véhicules anciens sont beaucoup moins accidentogènes

Calcul des ordres ratio

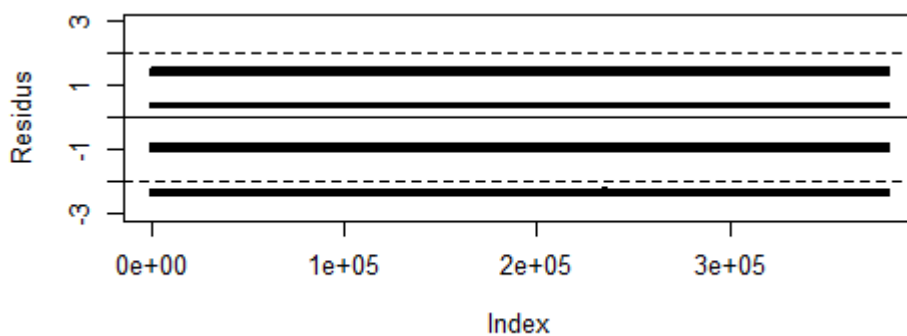
ODD_RATIO			
(Intercept)	Age	SexeHomme	vehicule_AgePlus de 5 ans
16.32212344	0.99663938	1.16854491	0.03732698
Variable	Odds Ratio (OR)	Interprétation	
(Intercept)	16.32	Pour les modalités de référence (âge = 0, femme, véhicule < 5 ans), les chances d'accident sont élevées.	
Âge	0.9966	Chaque année de plus diminue légèrement les chances d'accident (0.34 % en moins environ par an).	
Sexe : Homme	1.17	Les hommes ont 16.8 % de chances en plus d'avoir un accident comparé aux femmes.	
Âge du véhicule > 5 ans	0.037	Les véhicules de plus de 5 ans ont 96.3 % de chances en moins d'avoir un accident.	

Effets marginaux

summary(margins_model1)							
factor	AME	SE	z	p	lower	upper	
Age	-0.0005	0.0001	-9.1007	0.0000	-0.0006	-0.0004	
SexeHomme	0.0239	0.0013	18.3920	0.0000	0.0214	0.0265	
vehicule_AgePlus de 5 ans	-0.5734	0.0017	-333.6835	0.0000	-0.5767	-0.5700	

Variable	AME	Interprétation
Âge	-0.0005	Chaque année supplémentaire diminue la probabilité d'accident de 0,05 point de pourcentage.
Sexe : Homme	0.0239	Être un homme augmente la probabilité d'accident de 2,39 points de pourcentage par rapport aux femmes.
Âge du véhicule > 5 ans	-0.5734	Avoir un véhicule de plus de 5 ans réduit la probabilité d'accident de 57,34 points de pourcentage.

Analyse des résidus



4) CALCULER LE TMC : taux de mauvais classement et Matrice de Confusion

Calculer le taux de mauvais classement à partir de la matrice de confusion

	accident	no accident
0	9600	139281
1	156753	76520

5) INDICATEURS DE PERFORMANCES

AUC: 0.8150445

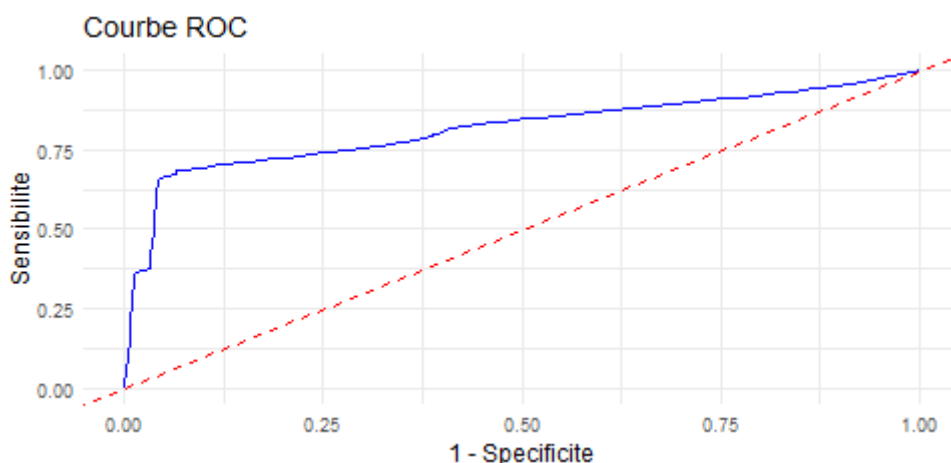
◆ Qu'est-ce que l'AUC ?

L'**AUC** (Area Under the Curve) est l'aire sous la courbe ROC. Elle mesure la **capacité du modèle à distinguer entre deux classes** : ici, les assurés ayant eu un **accident** (classe 1) et ceux n'en ayant **pas eu** (classe 0).

Interprétation de ton AUC = 0.815

- Cela signifie que **dans 81,5 % des cas**, ton modèle **prédit correctement** qu'un individu ayant eu un accident a un score de risque supérieur à celui d'un individu n'en ayant pas eu.
- Plus précisément, si tu prends **au hasard un assuré accidenté et un non-accidenté**, ton modèle a **81,5 % de chances** d'assigner un **score de probabilité plus élevé** à l'accidenté.

Calculer la courbe ROC



Interpretation :

. ****Compréhension des Axes**** - ****Axe X (1 - Spécificité)**** : Représente le taux de faux positifs. Une valeur proche de 0 indique une faible probabilité de faux positifs, tandis qu'une valeur proche de 1 indique une forte probabilité. - ****Axe Y (Sensibilité)**** : Représente le taux de vrais positifs, c'est-à-dire la capacité du modèle à identifier correctement les positifs. Une valeur proche de 1 indique une bonne performance.

. ****Interprétation de la Courbe**** - ****Forme de la Courbe**** : La courbe commence à (0,0) et finit à (1,1). Plus la courbe est proche du coin supérieur gauche du graphique, meilleure est la performance du modèle. - ****Zone sous la courbe (AUC)**** : Bien que non visible sur l'image, l'AUC (Area Under the Curve) est une mesure clé. Elle varie de 0 à 1. Une AUC proche de 1 indique un excellent modèle, tandis qu'une AUC de 0.5 indique un modèle aléatoire.

. ****Comparaison avec la Diagonale Rouge**** - La ligne rouge diagonale représente un modèle aléatoire. Si la courbe ROC est au-dessus de cette ligne, cela signifie que le modèle a une meilleure capacité de discrimination qu'un modèle aléatoire.

. ****Analyse des Points Clés**** - ****Points de Coupure**** : Chaque point sur la courbe correspond à un seuil de classification. Il est important d'examiner ces seuils pour trouver un équilibre entre sensibilité et spécificité selon le contexte d'application. - ****Performance à Différents Seuils**** : Évalue comment la sensibilité et la spécificité changent avec différents seuils. Cela peut influencer le choix de seuil basé sur les coûts des faux positifs et des faux négatifs.

6) PREDICTION DE LA PROBABILITE DE FAIRE UN ACCIDENT*Affichage des parametres estimes*

PARAMETRES			
(Intercept)	Age	SexeHomme	vehicule_AgePlus de 5ans
2.792521454	-0.003366282	0.155759306	-3.288039006

Calcul des probabilités

PROBABILITE_PREDITE	MODALITE_PREDITE
	<dbl> <chr>
1	0.938 accident
2	0.361 no accident
3	0.398 no accident
4	0.947 accident
5	0.938 accident
6	0.947 accident
7	0.361 no accident
8	0.361 no accident
9	0.398 no accident
10	0.361 no accident

Après avoir exploré la régression logistique, qui est particulièrement adaptée à la **modélisation de variables binaires**, il est naturel de s'intéresser à un autre type de modèle linéaire généralisé : la **régression de Poisson**.

Alors que la régression logistique permet d'estimer la probabilité d'occurrence d'un événement (succès/échec), la **régression de Poisson** s'intéresse à un autre type de variable dépendante : **les variables de comptage**.

En d'autres termes, lorsque la variable cible ne se limite plus à deux états (0 ou 1), mais représente un **nombre entier de fois** qu'un événement se produit (par exemple, nombre d'appels à un service client, nombre de visites sur un site, nombre de cas dans une épidémie), la régression de Poisson devient un outil statistique adapté et puissant.

PARTIE III : MODELISATION POISSONNIENE DE "frequence"

Représentation graphique de la fréquence

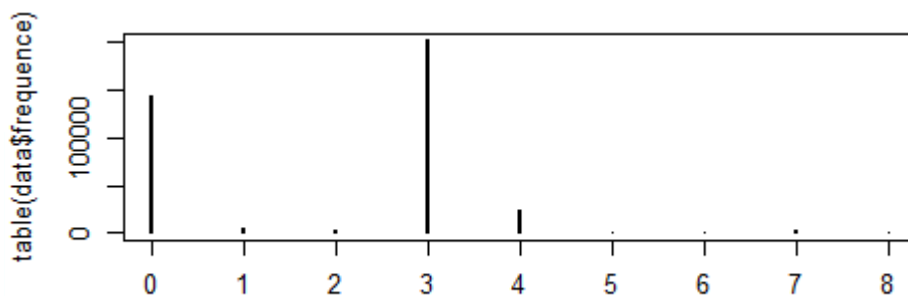


Figure 7 Frequence

la distribution des fréquences

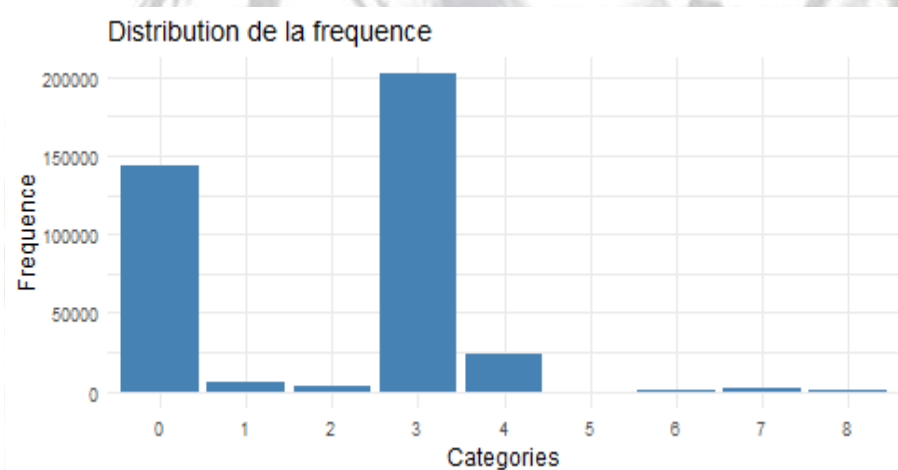
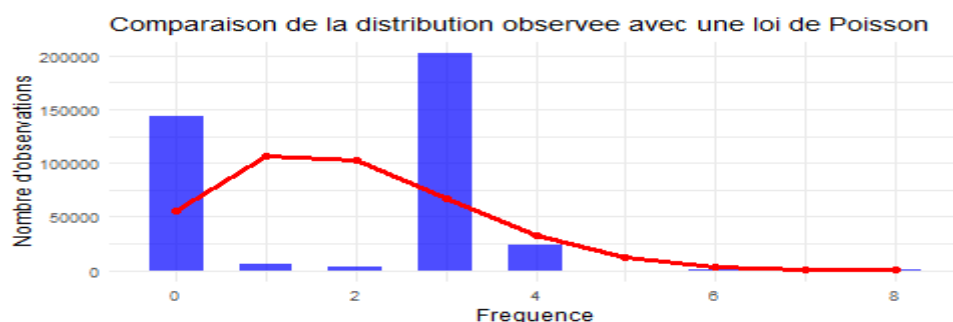


Figure 8 distribution des frequences

Calculer la moyenne empirique (λ estimé) :

Lambda estime : 1.934712



Calcul des fréquences théoriques (même n et même lambda)

Chi-squared test for given probabilities

```
data: obs_df$obs
X-squared = 629971, df = 8, p-value < 2.2e-16
mean(data$frequence)
[1] 1.934712
> var(data$frequence)
[1] 2.516623
```

NB: "Dans ce cas, la moyenne est relativement proche de la variance, ce qui permet d'envisager l'utilisation d'un modèle de régression de Poisson. En revanche, si la variance est nettement supérieure à la moyenne, on parle alors de phénomène de surdispersion, ce qui peut rendre le modèle de Poisson inadapté."

A) 🔍 Test d'adéquation à la loi de Poisson**Hypothèses du test :**

- H_0 (hypothèse nulle) : la variable suit une *loi de Poisson*
- H_1 (hypothèse alternative) : la variable **ne suit pas** une loi de Poisson.

Règle de décision :

- Si la **p-value** $< 0,05$: on **rejette H_0** , la variable ne suit **pas** une loi de Poisson.
- Si la **p-value** $\geq 0,05$: on **ne rejette pas H_0** , on peut considérer que la variable **suit** une loi de Poisson.

A.1) Estimation des paramètres par la méthode du maximum de vraisemblance

```
> summary(fpois)
Fitting of the distribution ' pois ' by maximum likelihood
Parameters :
      estimate Std. Error
lambda 1.934712 0.002250033
Loglikelihood: -721787.7  AIC: 1443577  BIC: 1443588
```

A.2) Tests d'adéquation

```
Chi-squared statistic: 623997.5
Degree of freedom of the Chi-squared distribution: 5
Chi-squared p-value: 0
Chi-squared table:
      obscounts theocounts
<= 0 143473.0000 55208.1946
```

```

<= 1 5408.0000 106811.9717
<= 2 3530.0000 103325.2163
<= 3 202303.0000 66634.8547
<= 4 23983.0000 32229.8178
<= 7 3132.0000 17603.8560
> 7 325.0000 340.0889
Goodness-of-fit criteria
                                1-mle-pois
Akaike's Information Criterion 1443577
Bayesian Information Criterion 1443588

```

NB : p.value inférieur à 5%, les données suivent une loi de poisson

B) Construction du modèle BINOMIALE NEGATIVE

B.1) Spécification du modèle

`summary(fnb)`

Call:

```
glm.nb(formula = frequence ~ Age + Sexe + Vehicle_Age + fuel_type +
  manufacture_year, data = data, init.theta = 16352.57622,
  link = log)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	4.0683591	0.5669911	7.175	7.21e-13	***
Age	-0.0020675	0.0001558	-13.266	< 2e-16	***
SexeHomme	0.0373507	0.0023577	15.842	< 2e-16	***
Vehicle_AgePlus de 5 ans	-0.9201941	0.0044765	-205.563	< 2e-16	***
fuel_typeSuper	-0.0243423	0.0024696	-9.857	< 2e-16	***
manufacture_year	-0.0014609	0.0002823	-5.176	2.27e-07	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for Negative Binomial(16352.58) family taken to be 1)

Null deviance: 726760 on 382153 degrees of freedom

Residual deviance: 558676 on 382148 degrees of freedom

AIC: 1275492

Number of Fisher Scoring iterations: 1

Theta: 16353

Std. Err.: 3520

warning while fitting theta: nombre limite d'iterations atteint

2 x log-likelihood: -1275478

Interprétation du modèle fnb

Voici un résumé des points clés à commenter :

 Significativité des variables

Toutes les variables sont hautement significatives ($p < 0.001$), ce qui indique qu'elles ont un **effet statistique clair** sur la fréquence.

🔍 Interprétation des coefficients (en log)

- (Intercept) = 4.068 → C'est le log de la fréquence attendue pour une personne de référence (valeurs de base pour toutes les variables).
- Age = -0.00207 → À chaque année supplémentaire, la fréquence attendue diminue légèrement.
- SexeHomme = +0.0373 → Les hommes ont une fréquence légèrement plus élevée que les femmes.
- Vehicle_AgePlus de 5 ans = -0.92 → Les véhicules plus anciens ont une fréquence beaucoup plus faible.
- fuel_typeSuper = -0.0243 → L'utilisation de carburant "Super" est associée à une légère baisse de fréquence.
- manufacture_year = -0.00146 → Les véhicules plus récents ont une fréquence un peu plus faible.

📌 Tous les coefficients s'interprètent sur l'échelle logarithmique (log-linéaire).

⚠️ Theta élevé = surdispersion importante

- Theta $\approx 16\,353$ (écart-type ≈ 3520) indique une **forte dispersion**.
- Le message "nombre limite d'itérations atteint" est un **warning** : le modèle a eu du mal à estimer le paramètre de dispersion, ce qui suggère une complexité importante dans la structure des données. Ça mérite peut-être un ajustement ou une validation plus fine.

B.2) ◆ AIC et ◆ BIC

-AIC = Akaike Information Criterion

```
> aic_val
[1] 1275492
```

BIC = Bayesian Information Criterion

```
> bic_val
[1] 1275568
```

📊 Interprétation de tes valeurs

- L'**AIC = 1 275 492** et le **BIC = 1 275 568** sont **relativement proches**, ce qui est **bon signe** : le modèle n'est **pas surparamétré**.
 - Ces valeurs **ne sont pas interprétées seules**, mais **comparées à d'autres modèles**.
- 👉 Un modèle alternatif avec **AIC ou BIC plus bas** serait **préférable**.

B.3) Vérification des hypothèses du modèle

Résidus simulés pour vérifier la distribution, l'hétéroscédasticité, l'indépendance

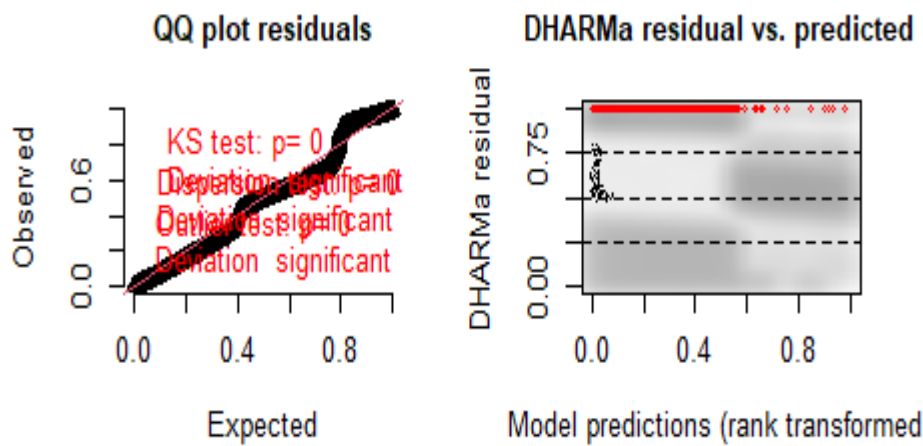


Figure 9 Residus

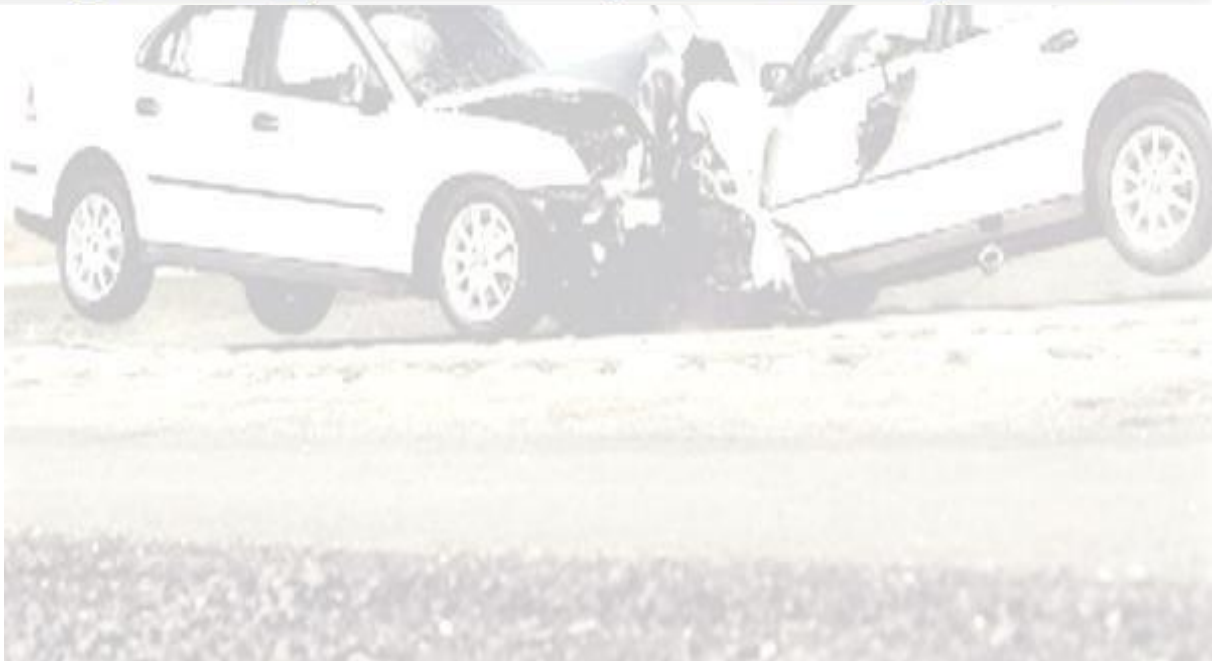
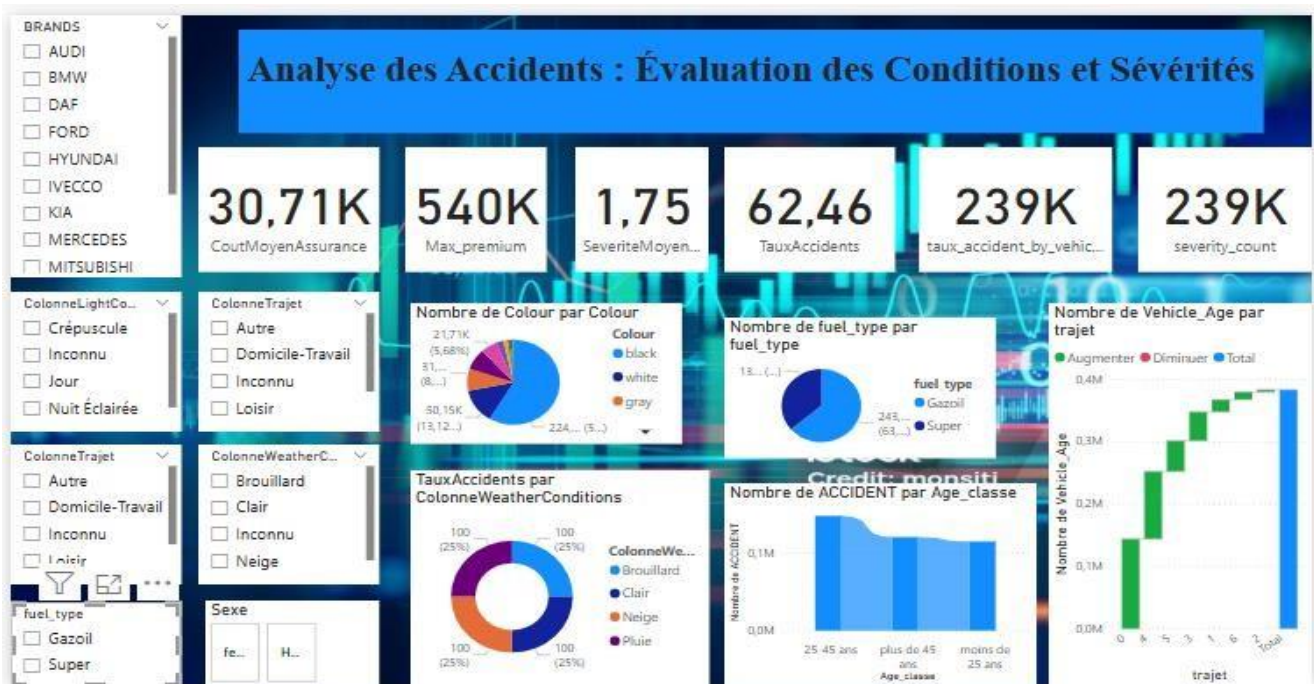
Multicolinéarité via VIF

```
print(vif_vals)
age      Sexe      vehicle_Age      fuel_type manufacture_year
3.255026  1.026810  3.266908      1.035347      1.035339
```

Interprétation

- Toutes les valeurs sont **inférieures à 5**, ce qui signifie qu'il n'y a **pas de multicolinéarité préoccupante** dans ton modèle.
- Les variables Age et Vehicle_Age ont des VIF un peu plus élevés (~3.26), mais restent dans une zone **acceptable**. Cela peut simplement refléter un certain **recoupement logique** entre l'âge du conducteur et celui du véhicule, sans pour autant invalider le modèle.
- Les autres variables (Sexe, fuel_type, manufacture_year) ont des VIF très proches de 1 → **aucune corrélation notable** avec les autres variables.

DASHBOARD





CONCLUSION

Rappel de la problématique L'objectif principal de cette étude était de modéliser le risque d'accident automobile à partir d'un jeu de données fourni par une compagnie d'assurance IARD. L'analyse s'appuie sur un ensemble de 374 393 observations décrivant le profil des assurés et les caractéristiques de leurs véhicules. Deux axes ont été explorés : - La modélisation de la **probabilité de survenance d'un accident** à l'aide d'une **régression logistique**, après binarisation de la variable cible fréquence. - La modélisation du **nombre d'accidents** via une **régression de Poisson**, plus adaptée aux données de type "comptage".

Résumé des principaux résultats obtenus - Régression logistique : Les variables telles que l'**âge**, le **sexe**, et l'**ancienneté du véhicule** se sont révélées significatives dans la prédiction du risque d'accident. Par exemple : - Les **hommes** présentent une probabilité d'accident légèrement plus élevée que les femmes. - Les **véhicules plus anciens** sont associés à un **risque moindre** dans les données, ce qui peut refléter des comportements plus prudents ou d'autres facteurs non observés. - L'**âge de l'assuré** est **négativement corrélé** à la probabilité d'accident : plus l'assuré est âgé, moins il est susceptible d'avoir un accident.

- **Régression de Poisson :** L'écart modéré entre la moyenne et la variance du nombre d'accidents justifie l'usage du modèle de Poisson. Les résultats montrent que plusieurs variables explicatives influencent significativement la fréquence des sinistres.
- Les indicateurs AIC et BIC ont permis de comparer les modèles et de sélectionner les spécifications les plus efficaces tout en évitant le surajustement.
- L'**AUC** de la courbe ROC montre une bonne capacité discriminante du modèle logistique.

Recommandations - Utiliser le **modèle logistique** pour identifier les profils à risque élevé afin de cibler les campagnes de prévention ou ajuster les primes. - Utiliser le **modèle de Poisson** pour anticiper le **volume de sinistres**, ce qui peut guider les décisions actuarielles et budgétaires. - Intégrer ces modèles dans un outil de **scoring client** pour une tarification plus fine et un meilleur pilotage des risques.

Limites de l'analyse - Binarisation de la variable fréquence : une simplification qui peut entraîner une perte d'information. - Certaines variables potentiellement explicatives n'étaient pas présentes dans le jeu de données (ex : historique de conduite, bonus/malus, type de contrat). - **Surdispersion possible** dans les données de comptage qui pourrait justifier l'usage de modèles alternatifs (ex : régression négative binomiale).

Perspectives et approfondissements - Tester d'autres modèles comme : - **Régression binomiale négative** pour mieux traiter la surdispersion. - **Arbres de décision ou forêts aléatoires** pour capter d'éventuelles non-linéarités et interactions complexes. - Intégrer des **variables temporelles** (ex : saison, année) pour modéliser les tendances ou effets calendaires. - Croiser avec d'autres données externes (ex : géolocalisation, trafic routier, zones à risques) pour enrichir les modèles. - Mettre en place un **système de scoring automatisé** dans les processus de souscription ou de gestion des contrats.

