

# INTERPRÉTABILITÉ DES MODÈLES DE TARIFICATION EN ACTUARIAT

## Application à l'assurance automobile

Franklin FEUKAM KOUHOUE

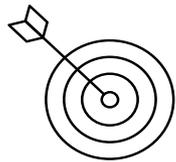
# SOMMAIRE

- 1 • MOTIVATION DE LA RECHERCHE D'INTERPRÉTABILITÉ EN TARIFICATION
- 2 • APPROCHES ET OUTILS D'INTERPRÉTABILITÉ
- 3 • CAS D'APPLICATION : Assurance automobile

Introduction



**Contexte:** En France, 83%. C'est la part des assureurs qui considèrent que l'IA va profondément modifier les processus internes et la relation client. (Source, [ACPR](#), 2022).

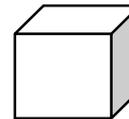


**Problématique:**

	$x_1$	$x_2$	$y$
Assuré 1	...	...	...
Assuré 2		...	...
...	...	...	...
Assuré n	...	...	...

$x_1$  = âge de l'assuré  
 $x_2$  = ancienneté du véhicule  
 $y$  = fréquence de sinistres

Glass-Box



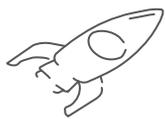
$$\hat{y} = f(x_1, x_2) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

Black-Box



$$\hat{y} = g(x_1, x_2)$$

*g opaque*



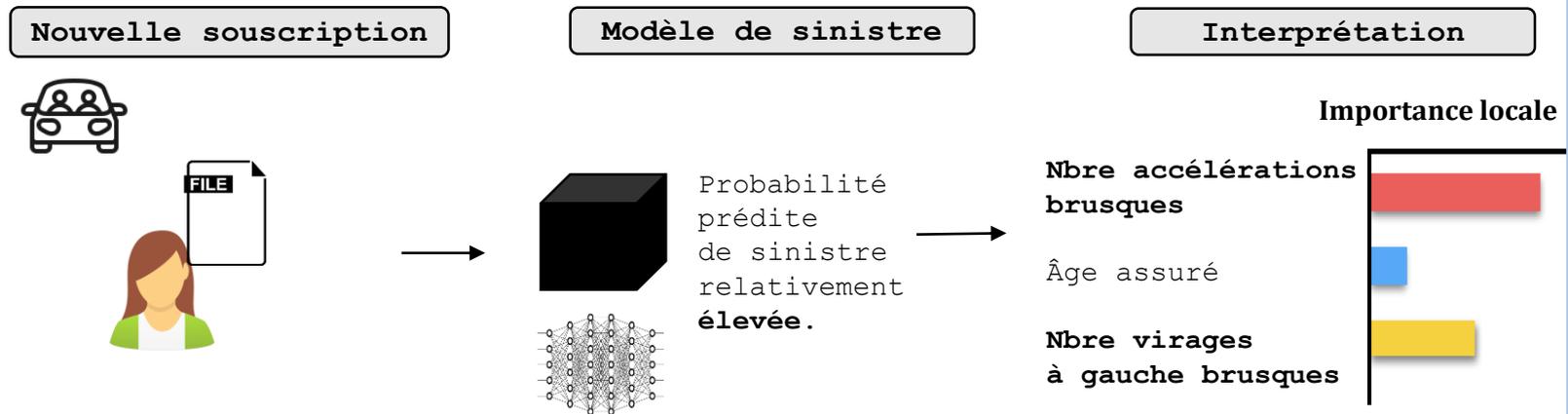
**Apports des travaux:** Dans cette présentation, nous illustrons comment les techniques d'**Explainable AI** peuvent améliorer la précision de la tarification en assurance tout en restant conformes et transparentes.



1 • Motivation de la recherche d'interprétabilité en tarification

**Pour l'Assureur & l'Assuré**

Comprendre les raisons de la hausse de la prime d'assurance.



**Observations**

En se servant de méthodes d'interprétabilité (par contrefactuelles), l'assureur peut formuler des stratégies d'amélioration du comportement de conduite de l'assuré. Ceci afin de réduire sa fréquence prédite de sinistre par le modèle lors de la prochaine souscription, et par ricochet la prime à payer.

2 • Approches et outils d'interprétabilité

Approches d'interprétabilité

Dans la littérature, nous distinguons de **deux grandes approches d'interprétabilité**:

- Interprétabilité dite **basée sur le modèle (IBM)**
- Interprétabilité dite **post-hoc**

Approche IBM	Approche post hoc
<ul style="list-style-type: none"><li>• Parcimonie</li><li>• Simulable</li><li>• Modulaire</li></ul>	<ul style="list-style-type: none"><li>• Méthodes Globales vs Méthodes Locales</li><li>• Méthodes spécifiques au modèle vs Méthodes indépendantes du modèle</li><li>• Méthodes dépendantes des données vs indépendantes des données</li></ul>

Remarque

Dans le mémoire nous nous sommes focalisés sur l'approche d'interprétabilité dite *post hoc*, et sur les outils qui sont indépendants des données et indépendants du modèle à expliquer.

3 • Cas d'application: assurance automobile

Etapes phares de la mise en oeuvre

Nous présentons ici les principales étapes de la mise en œuvre de ces méthodes d'interprétabilité dans notre contexte de tarification automobile.

01. TÂCHES À EFFECTUER

- Assurance automobile
- Prédire la fréquence de sinistre

02. DONNÉES DISPONIBLES

- Base de données synthétiques pour la télématique des conducteurs.
- Données classiques et télématiques
- Score de crédit
- [\[So et al., 2021\]](#), Université du Connecticut

03. MODÉLISATION

- GLM (Poisson)
- LocalGLMnet (Poisson)
- Forêt aléatoire

$$x \mapsto m(\mu(x)) = \beta_0 + \langle \beta(x), x \rangle = \beta_0 + \sum_{j=1}^p \beta_j(x)x_j$$

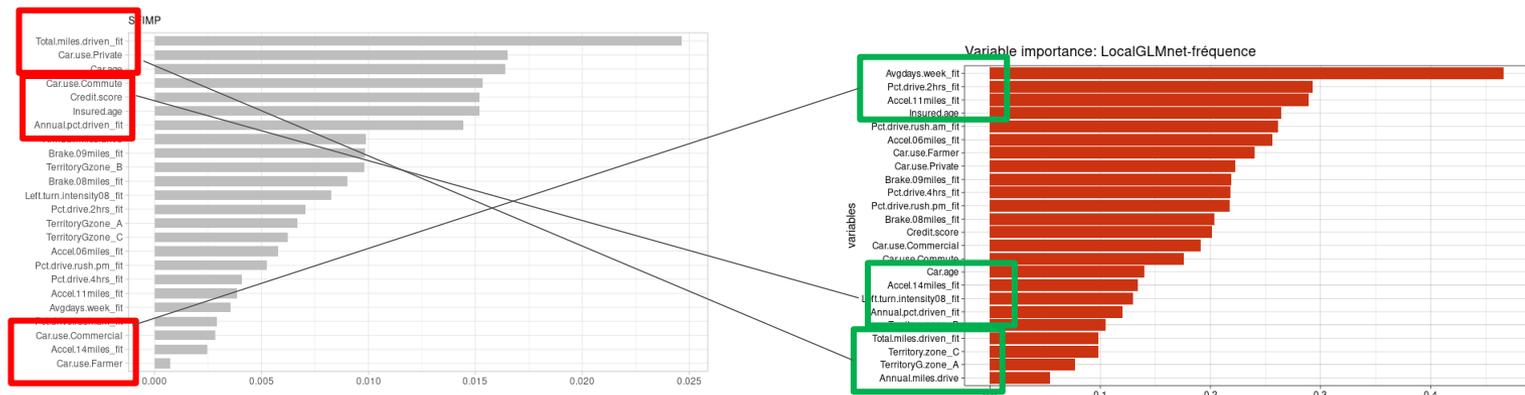
04. INTERPRÉTATION

- GLM Poisson
- LocalGLMnet Poisson

### 3 • Cas d'application: assurance automobile

#### Importance des caractéristiques

Un premier élément que nous regardons lorsque nous souhaitons interpréter un modèle est l'importance globale des caractéristiques.



#### Observation

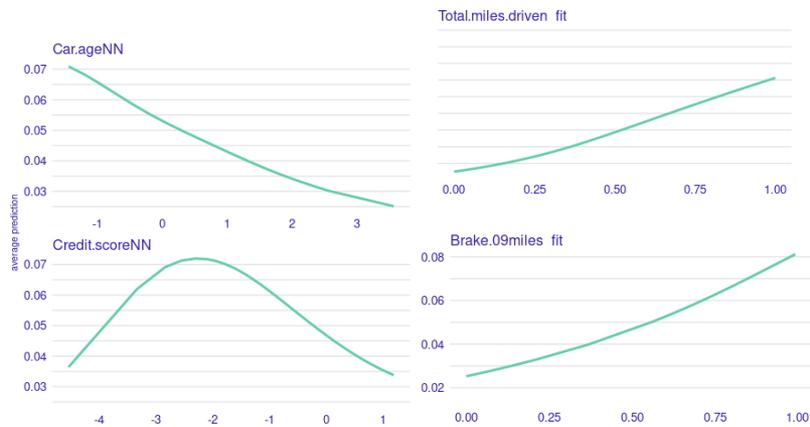
Le fait que les résultats des deux approches ne coïncident pas nécessairement est dû au fait qu'ils ne reposent pas sur le même principe de fonctionnement sous-jacent.

### 3 • Cas d'application: assurance automobile

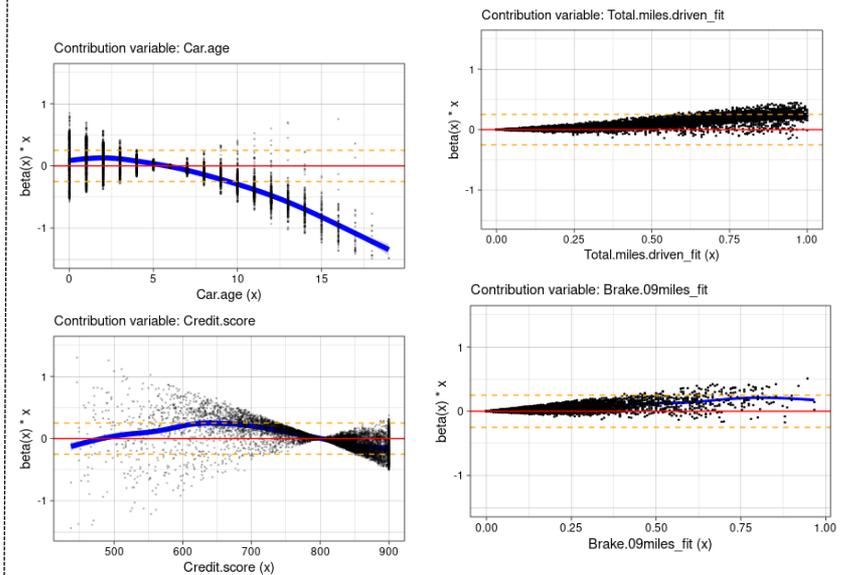
#### Effets des caractéristiques

Un deuxième élément que nous regardons lorsque nous souhaitons interpréter un modèle est l'effet global des caractéristiques.

#### *ALE (Post hoc)*



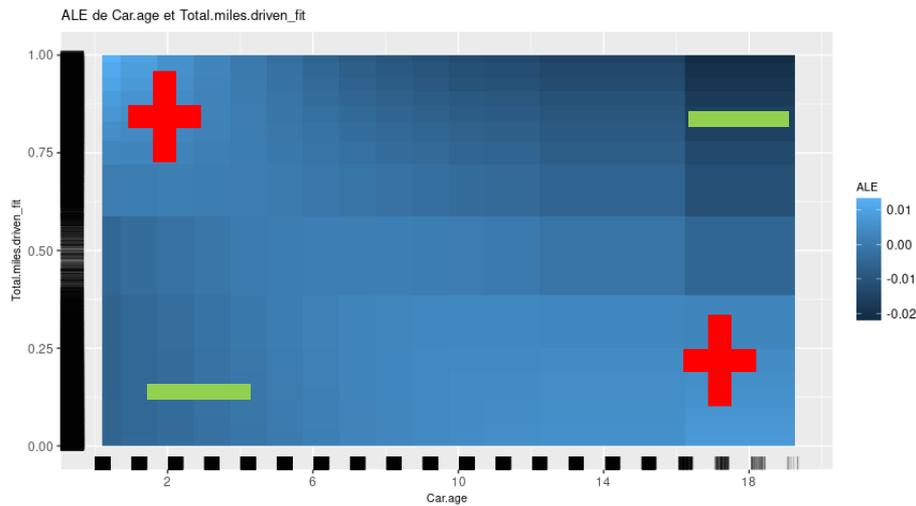
#### Contribution (IBM)



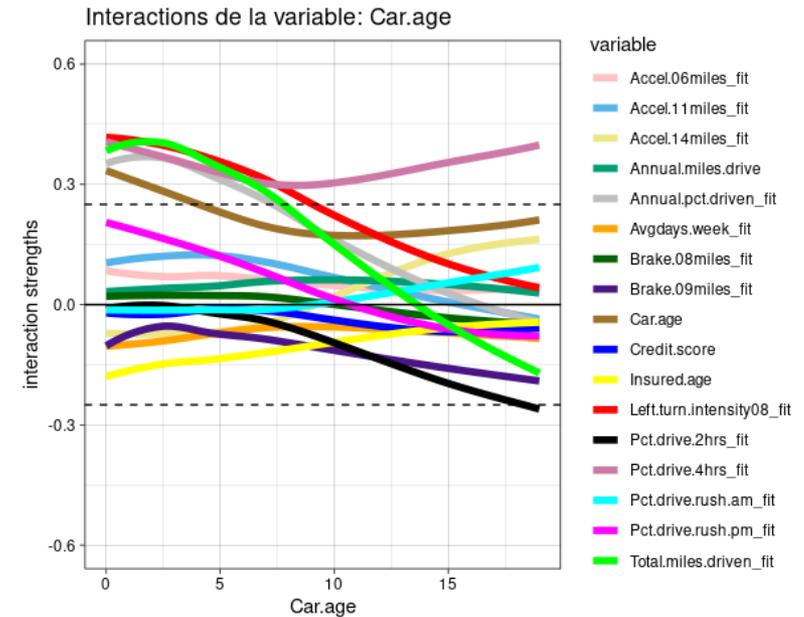
### 3. CAS D'APPLICATION: RÉSULTATS

#### Interprétation globale: Interactions entre les caractéristiques (3/3)

ALE-2D (Post hoc)



Interactions (IBM)



$$\nabla\beta_j(x) = \left( \frac{\partial}{\partial x_1} \beta_j(x), \dots, \frac{\partial}{\partial x_p} \beta_j(x) \right)^T \in \mathbb{R}^p$$

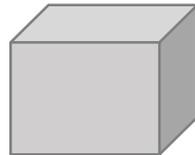
### 3. CAS D'APPLICATION: RÉSULTATS

#### Interprétation locale: illustration (1/5)



<i>Insured.age</i> : 20 ans	<i>Car.age</i> : 0 an; <i>Credit.score</i> : 580	<i>Annual.miles.drive</i> : 12427.2
<i>Car.use</i> : Commute	<i>TerritoryG</i> : zone_C	<i>Total.miles.driven_fit</i> : 0.92
<i>Annual.pct.driven_fit</i> : 0.63	<i>Pct.drive.rush.am_fit</i> : 0.11	<i>Left.turn.intensity08_fit</i> : 0.96
<i>Pct.drive.rush.pm_fit</i> : 0.95	<i>Brake.08miles_fit</i> : 0.95	<i>Brake.09miles_fit</i> : 0.83
<i>Pct.drive.4hrs_fit</i> : 0.90	<i>Pct.drive.2hrs_fit</i> : 0.98	<i>Accel.14miles_fit</i> : 0.85
<i>Accel.11miles_fit</i> : 0.87	<i>Accel.06miles_fit</i> : 0.96	<i>Avgdays.week_fit</i> : 1

LocalGLMnet  
fréquence



Fréquence prédite= **2.77**  
(La plus élevée du jeu  
de données test)

???

Client (e)



Régulateur



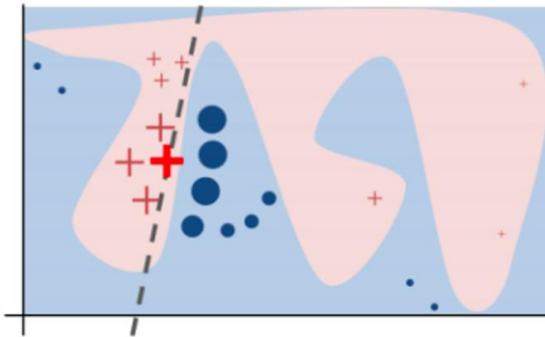
Assureur



### 3. CAS D'APPLICATION: RÉSULTATS

#### Interprétation locale: illustration (2/5)

LIME (*Post hoc*): illustration



LIME (*Post hoc*): définition

explicateur linéaire parcimonieux

$$\xi(x) = \operatorname{argmin}_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

mesure de l'infidélité  $\mathcal{L}(f, g, \pi_x)$     modèle interprétable  $g$     mesure de complexité  $\Omega(g)$   
 une instance  $x$     modèle à expliquer  $f$     mesure de proximité pour définir le voisinage local  $\pi_x$

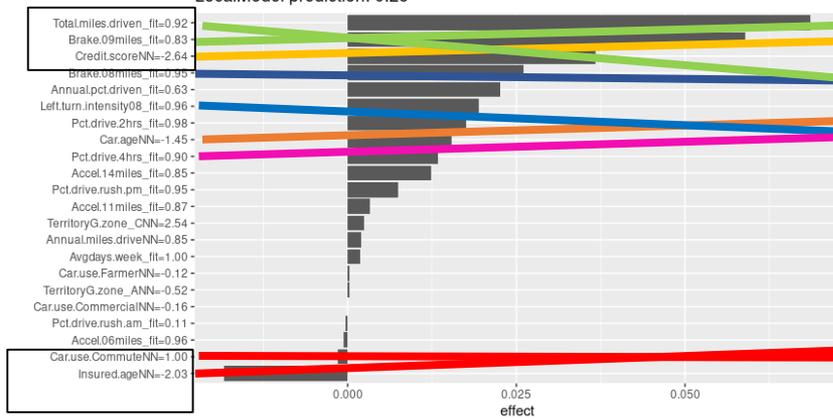
Intuition: Construire un modèle de substitution locale interprétable

## 3. CAS D'APPLICATION: RÉSULTATS

### Interprétation locale: illustration (3/5)

**LIME (Post hoc)**

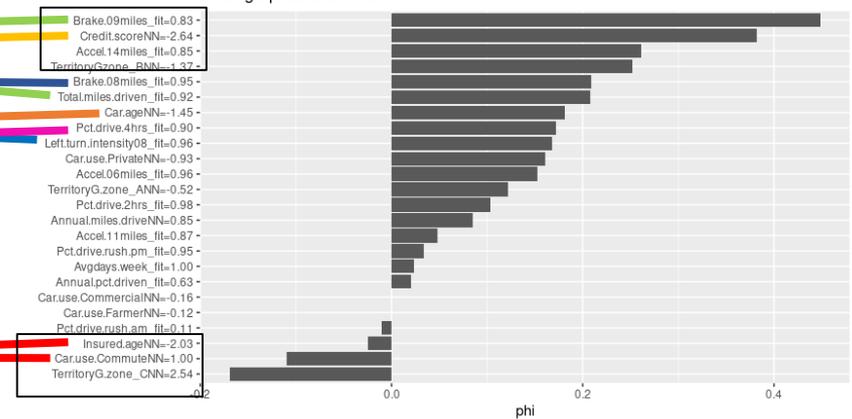
Actual prediction: 2.77  
LocalModel prediction: 0.25



- **Total.miles.driven\_fit**, **Brake.xxmiles** et **Credit.scoreNN** parmi celles qui contribuent le plus positivement.

**SHAP (Post hoc)**

Actual prediction: 2.77  
Average prediction: 0.05

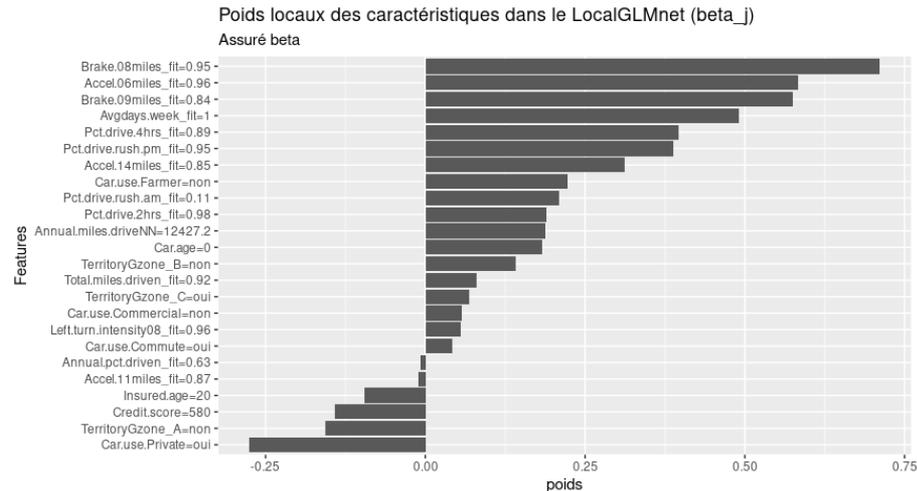


- **Insured.ageNN** contribue **négativement**.

## 3. CAS D'APPLICATION: RÉSULTATS

### Interprétation locale: illustration (4/5)

#### Interprétation Basée sur le Modèle (IBM)

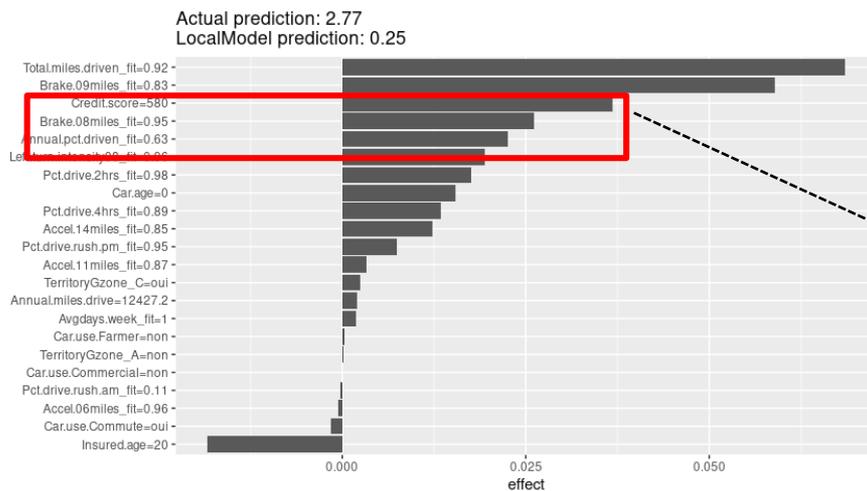


Les interprétations post-hoc sont globalement cohérentes avec l'interprétation basée sur le modèle, à quelques petites exceptions près.

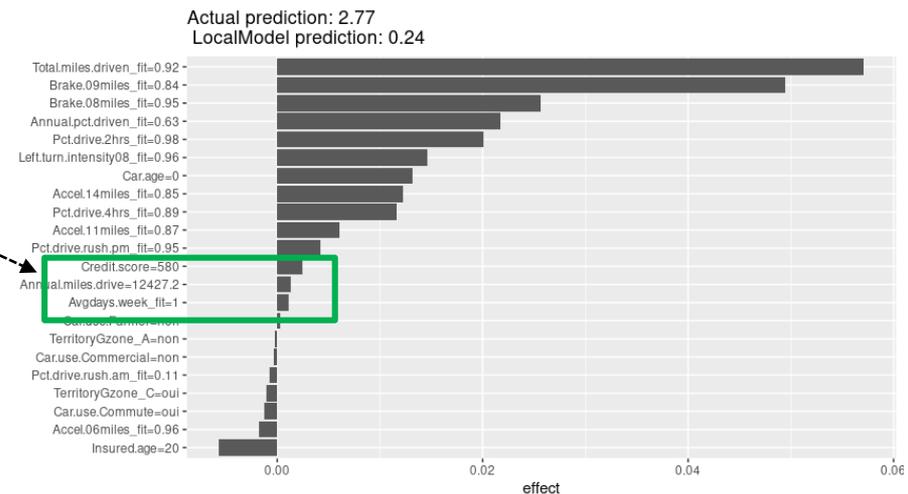
### 3. CAS D'APPLICATION: RÉSULTATS

#### Interprétation locale: limites de la méthode LIME (5/5)

Explications LIME *initiales*



Explications LIME *attaquées*



$$e(x) = \begin{cases} f(x), & \text{si } is\_OOD(x) \leq 0.8 \\ \psi(x), & \text{si } is\_OOD(x) > 0.8 \end{cases}$$

En dehors du poids de la variable **Credit.score** qui a été modifié, les explications sont quasiment restées inchangées par ailleurs.

## 3. CAS D'APPLICATION: RÉSULTATS

### Ingénierie des caractéristiques: principe (1/2)



#### Observation

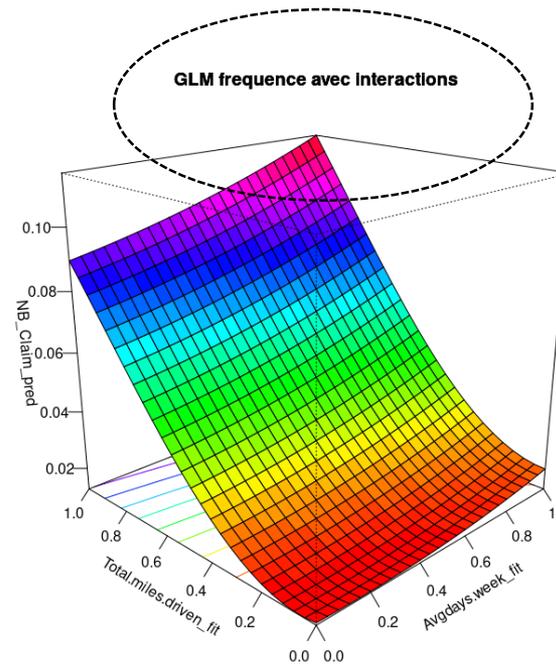
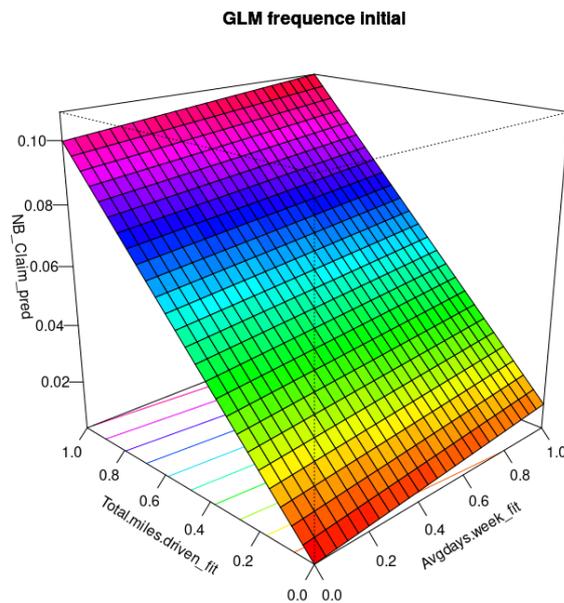
L'interprétabilité peut permettre de réaliser de l'ingénierie des caractéristiques

## 3. CAS D'APPLICATION: RÉSULTATS

### Ingénierie des caractéristiques: résultats (2/2)

#### Cas pratique

Sur ces graphiques nous illustrons un cas pratique où nous nous servons des résultats de l'interprétation de notre modèle hybride LocalGLMnet pour enrichir le modèle GLM initial.



## LIMITES ET PERSPECTIVES DE L'ÉTUDE

- Mettre en œuvre les explications par contrefactuelle;
- Il est difficile d'évaluer la qualité des interprétations et de les comparer objectivement les unes par rapport aux autres: pas de métriques pertinentes;
- Vers des modèles Hybrides (ce qui fait l'objet d'importantes de recherches scientifiques en ce moment, pas seulement en assurance);
- Automatisation des processus d'interprétation en compagnie d'assurance.

Merci de votre attention !