

# Interprétabilité des Modèles de Tarification en Actuariat

## *Application à l'Assurance Automobile*

Franklin FEUKAM KOUHOUE

*Actuaire Associé*

27 mars 2025

1

**I**NTRODUCTION

1 • CONTEXTE

De grands défis sociétaux

▶ **Changement Climatique**

- Augmentation de la **fréquence** des évènements climatiques **extrêmes** (inondations, tempêtes, grêle).

▶ **Hyperpersonnalisation des services**

- Attente des assurés : une offre tarifaire et de garanties adaptées à **leurs besoins** et à leur **profil de risque réel**.

Des réflexions émergentes

Tackling Climate Change with Machine Learning

David Rolnick<sup>1\*</sup>, Priya L. Donti<sup>2</sup>, Lynn H. Kaack<sup>3</sup>, Kelly Kochanski<sup>4</sup>, Alexandre Lacoste<sup>5</sup>, Kris Sankaran<sup>6,7</sup>, Andrew Slavov<sup>8</sup>, Nikola Milojevic-Dupont<sup>10,11</sup>, Natasha Jaques<sup>12</sup>, Anna Waldman-Brown<sup>12</sup>, Alexandra Luccioni<sup>6,7</sup>, Tegan Maharaj<sup>6,8</sup>, Evan D. Sherwin<sup>9</sup>, S. Karthik Mukkavilli<sup>6,7</sup>, Konrad P. Körding<sup>1</sup>, Carla Gomes<sup>13</sup>, Andrew Y. Ng<sup>14</sup>, Demis Hassabis<sup>15</sup>, John C. Platt<sup>16</sup>, Felix Creutzig<sup>10,11</sup>, Jennifer Chayes<sup>17</sup>, Yoshua Bengio<sup>6,7</sup>

<sup>1</sup>University of Pennsylvania, <sup>2</sup>Carnegie Mellon University, <sup>3</sup>ETH Zürich, <sup>4</sup>University of Colorado Boulder, <sup>5</sup>Element AI, <sup>6</sup>Mila, <sup>7</sup>Université de Montréal, <sup>8</sup>École Polytechnique de Montréal, <sup>9</sup>Harvard University, <sup>10</sup>Mercator Research Institute on Global Commons and Climate Change, <sup>11</sup>Technische Universität Berlin, <sup>12</sup>Massachusetts Institute of Technology, <sup>13</sup>Cornell University, <sup>14</sup>Stanford University, <sup>15</sup>DeepMind, <sup>16</sup>Google AI, <sup>17</sup>Microsoft Research



The Insurance Market in the Era of Digital Transitions  
Relationships Between Insurers, Big Tech, and Insurtech

**AUTHORS** Arthur Charpentier, Université du Québec à Montréal  
Raphaël Suire, Nantes University, IAE—Graduate school of Management, France  
**SPONSOR** Actuarial Innovation and Technology Strategic Research Program Steering Committee



Point d'Attention

Le Manque d'Interprétabilité : un obstacle

Original Research Article



Artificial intelligence and personalization of insurance: Failure or delayed ignition?

Big Data & Society  
January-March 2025  
© The Author(s) 2025  
Article reuse guidelines:  
sagepub.com/journals-permissions  
DOI: 10.1177/20539517241291817  
journals.sagepub.com/home/bds

Arthur Charpentier<sup>1</sup> and Xavier Vamparys<sup>2</sup>



1 • PLAN DE LA PRESENTATION

**2. Motivation de la recherche d'Interprétabilité**

**3. Méthodes d'Interprétabilité explorées**

**4. Résultats du cas d'application**

# 2

## **M**OTIVATION DE LA RECHERCHE D'INTERPRÉTABILITÉ

2 • MOTIVATION DE LA RECHERCHE D'INTERPRETABILITE

Pour L'ACTUAIRE

L'Actuaire est le concepteur et le **garant technique** du modèle d'évaluation du risque (modèle de calcul de la **prime d'assurance**)

Les Raisons

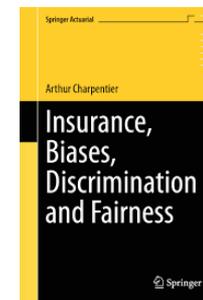
- ▶ **Comprendre le fonctionnement du Modèle**
  - **Pédagogie interne:** « *Ce que l'on conçoit bien s'énonce clairement, Et les mots pour le dire arrivent aisément.* » **Nicolas Boileau** (1636 - 1711).
  - **Mieux maîtriser** le risque (assurance indicielle).
- ▶ **Valider et contrôler les éventuelles risques Modèle**
  - Détecter des éventuelles **biais d'équité** ou **biais de discrimination tarifaire** dus à l'utilisation des variables sensibles, telles que le **genre** ou l'**origine ethnique**, pour segmenter les assurés.
  - Détecter les autres éventuelles **incohérences** et **zones d'instabilité** dans le fonctionnement du modèle.

Pour aller plus loin

IA, biais et équité  
(en actuariat et en assurance)

**Arthur Charpentier**  
avec Laurence Barry, Marie-Pier Côté, Olivier Côté,  
Agathe Fernandes-Machado, Ewen Gallic, François Hu, Philipp Ratz  
(et Ana Patrón Piñerez, Mulah Moriah, etc)

Marseille, Février 2025



2 • MOTIVATION DE LA RECHERCHE D'INTERPRETABILITE

Pour L'ASSUREUR

L'Assureur (ou Compagnie d'assurance) porte la responsabilité globale des décisions à l'issue du processus de tarification.

Les Raisons

- ▶ **S'assurer de la conformité réglementaire**
  - Exigence de **transparence** de l'ACPR.
  - Le règlement européen sur l'IA (AI Act)
- ▶ **Améliorer les stratégies de gestion du risque**
  - Comprendre **les raisons d'une hausse de la prime** de l'assuré. Puis fournir à l'assuré des recommandations pour améliorer son profil de risque.
  - Identifier des **nouvelles corrélations** (nouveaux facteurs de risques) peu intuitives.
- ▶ **Identifier les profils de risques émergents**

Pour aller plus loin



Gouvernance des algorithmes d'intelligence artificielle dans le secteur financier

Ce document de réflexion s'inscrit dans le cadre des travaux menés par l'ACPR sur l'intelligence artificielle (IA) depuis 2018. En mars 2019, après un premier rapport et une première consultation publique, l'ACPR a lancé des travaux exploratoires avec quelques acteurs du secteur financier afin d'éclairer les enjeux d'explicabilité et de gouvernance de l'IA - au sens essentiellement de Machine Learning (ML). Composés d'entretiens et d'ateliers techniques, ils couvraient trois domaines : la lutte contre le blanchiment et le financement du terrorisme (LCB-FT), les modèles internes et en particulier le scoring de crédit, et la protection de la clientèle. Deux axes d'étude en sont ressortis : ceux de l'évaluation et de la gouvernance des algorithmes d'IA.



3

**M**ÉTHODES D'INTERPRÉTABILITÉ EXPLORÉES

3 • METHODES D'INTERPRETABILITE EXPLOREES

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

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

**IBM**

- ▶ **Parcimonie**
- ▶ **Modularité**
- ▶ **Simulabilité**

**Post hoc**

- ▶ **Méthodes globales vs Méthodes locales**
- ▶ **Méthodes spécifiques au modèle vs Méthodes indépendantes du modèles**
- ▶ **Méthodes dépendantes des données vs Méthodes indépendantes des données**

**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**.

# 4

## RÉSULTATS DU CAS D'APPLICATION

4 • 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

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

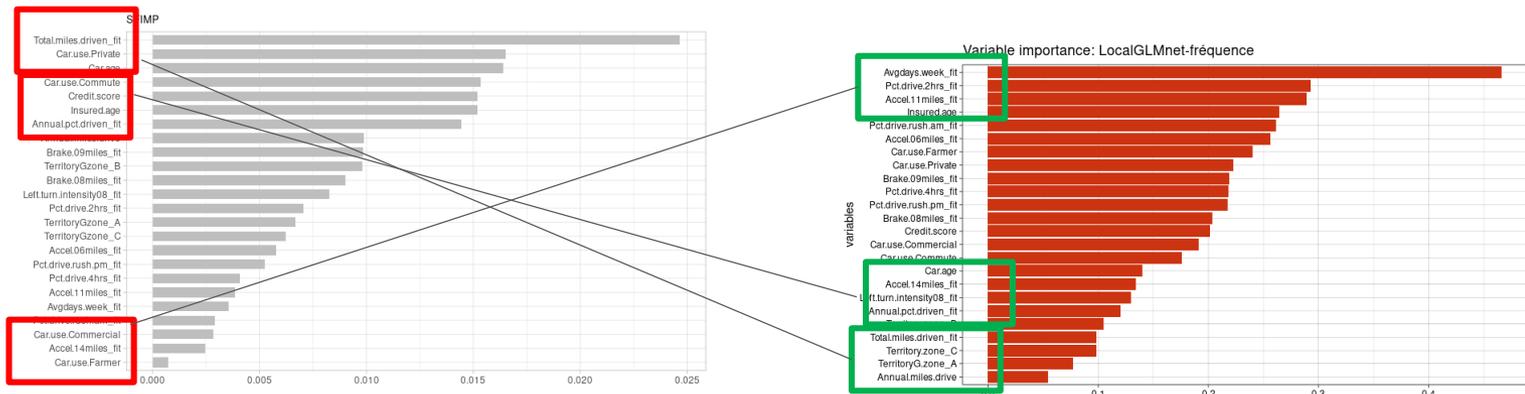
04. INTERPRÉTATION

- GLM Poisson
- LocalGLMnet Poisson

4 • 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

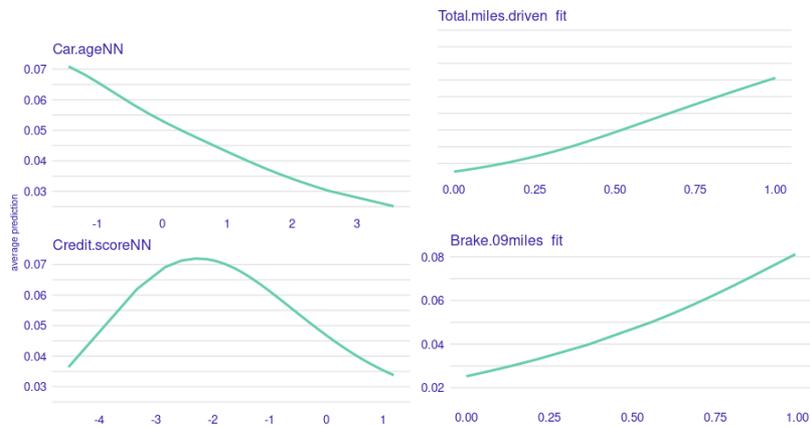
Le fait que les résultats des deux approches ne coïncident pas nécessairement est du au fait qu'ils ne reposent pas sur le même principe de fonctionnement sous-jacent.

4 • CAS D'APPLICATION: assurance automobile

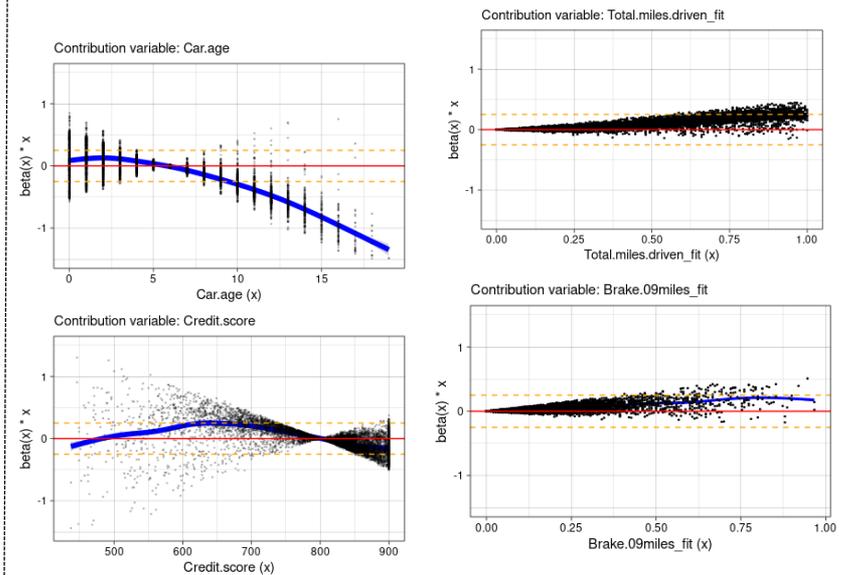
Effets des caractéristiques

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

ALE (Post hoc)



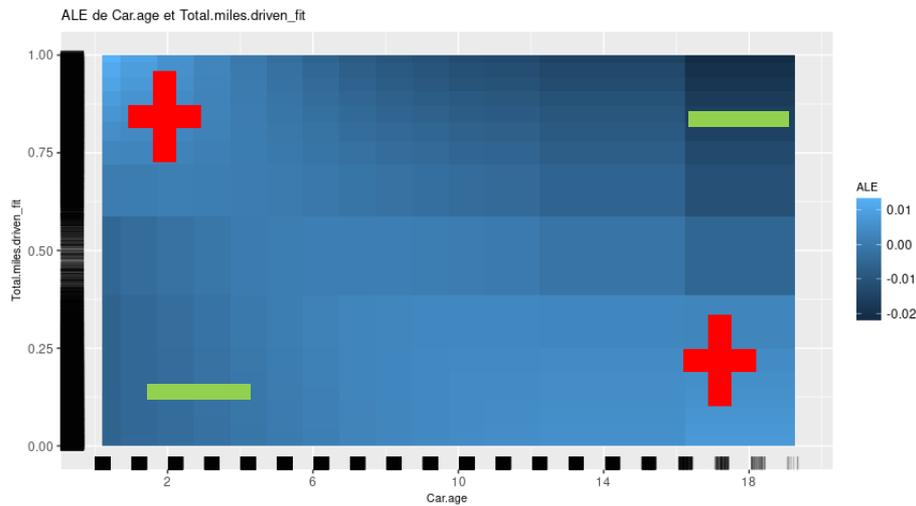
Contribution (IBM)



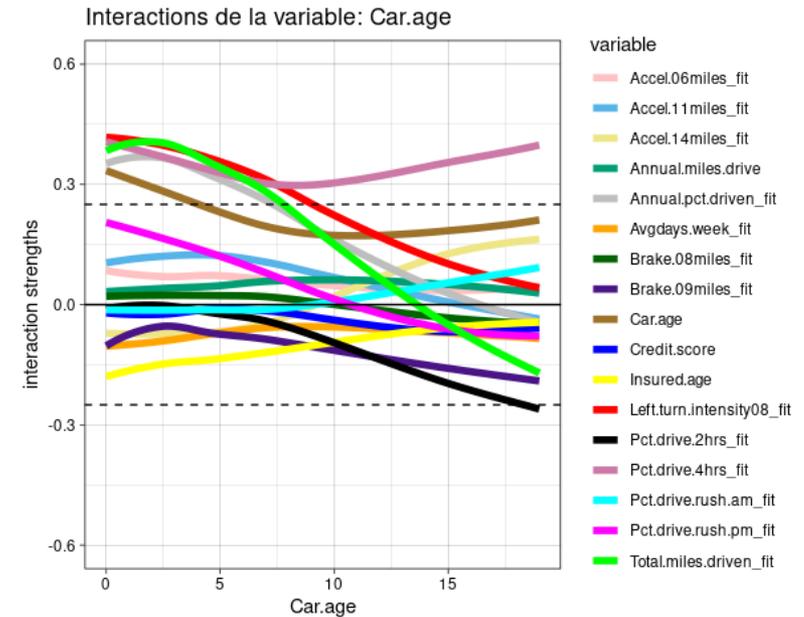
## 4 • CAS D'APPLICATION: ASSURANCE AUTOMOBILE

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

ALE-2D (Post hoc)



Interactions (IBM)



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

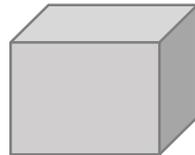
## 4 • CAS D'APPLICATION: ASSURANCE AUTOMOBILE

### 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



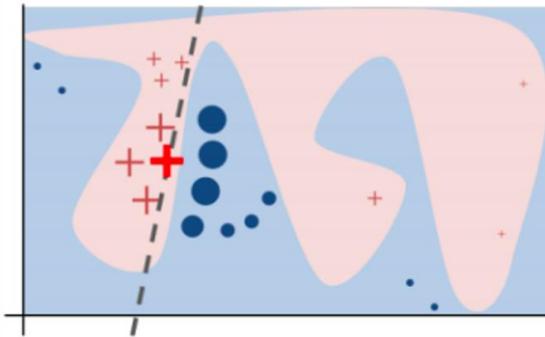
Assureur



## 4 • CAS D'APPLICATION: ASSURANCE AUTOMOBILE

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

#### LIME (*Post hoc*): illustration



**Intuition: Construire un modèle de substitution locale interprétable**

#### LIME (*Post hoc*): définition

explicateur linéaire parcimonieux

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

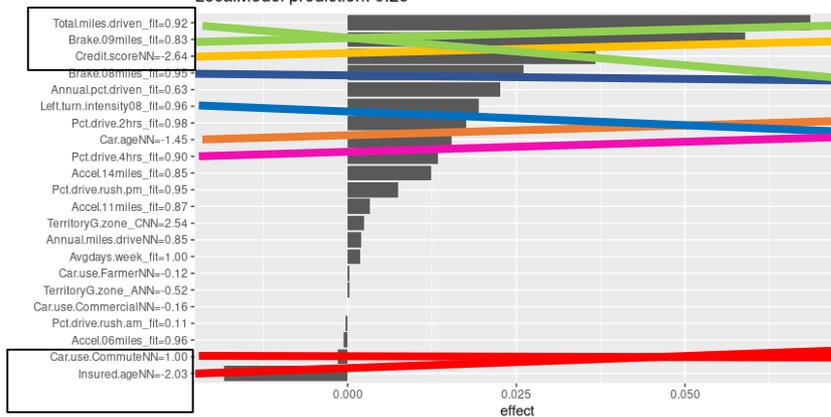
mesure de l'infidélité  $\mathcal{L}(f, g, \pi_x)$     modèle interprétable  $G$     mesure de complexité  $\Omega(g)$   
 une instance  $x$     modèle à expliquer  $f$     mesure de proximité pour définir le voisinage local  $\pi_x$

# 4 • CAS D'APPLICATION: ASSURANCE AUTOMOBILE

## Interprétation locale: illustration (3/5)

**LIME (Post hoc)**

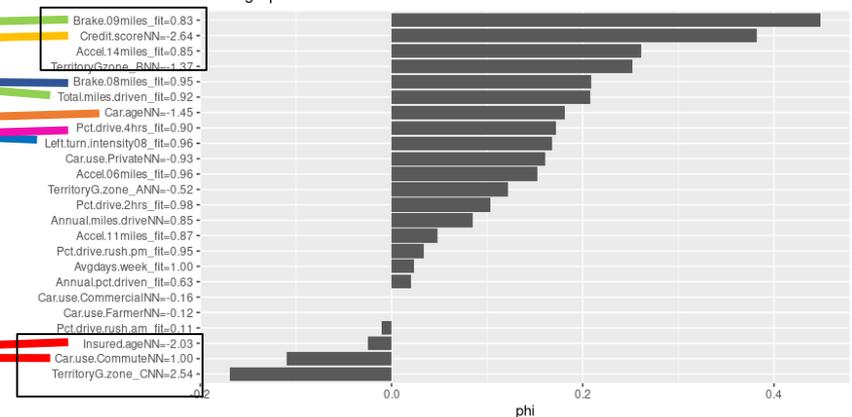
Actual prediction: 2.77  
LocalModel prediction: 0.25



- **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

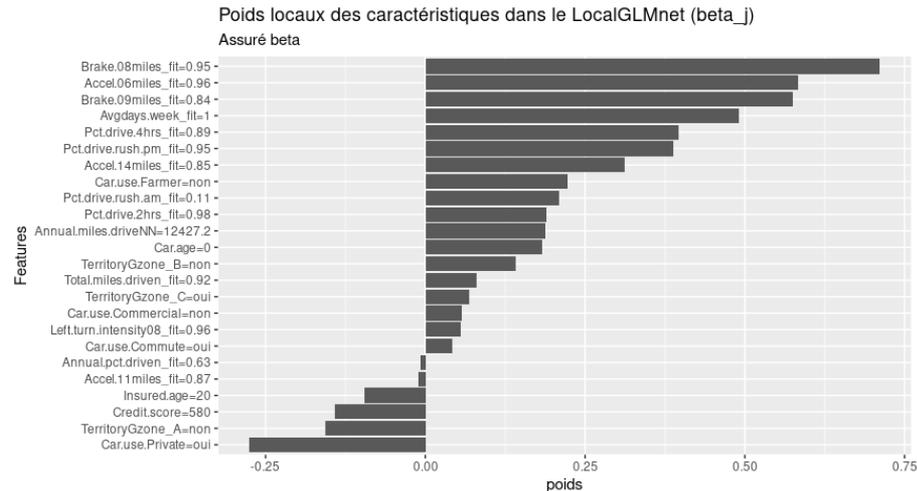


- **Insured.ageNN** contribue **négativement**.

## 4 • CAS D'APPLICATION: ASSURANCE AUTOMOBILE

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

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

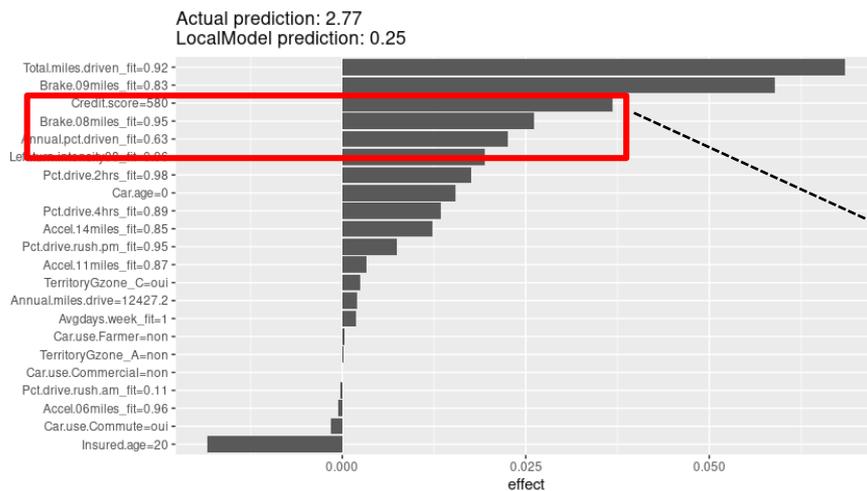


Les interprétations post-hoc sont globalement cohérentes avec l'interprétation basée sur le modèle, à quelques petites exceptions près.

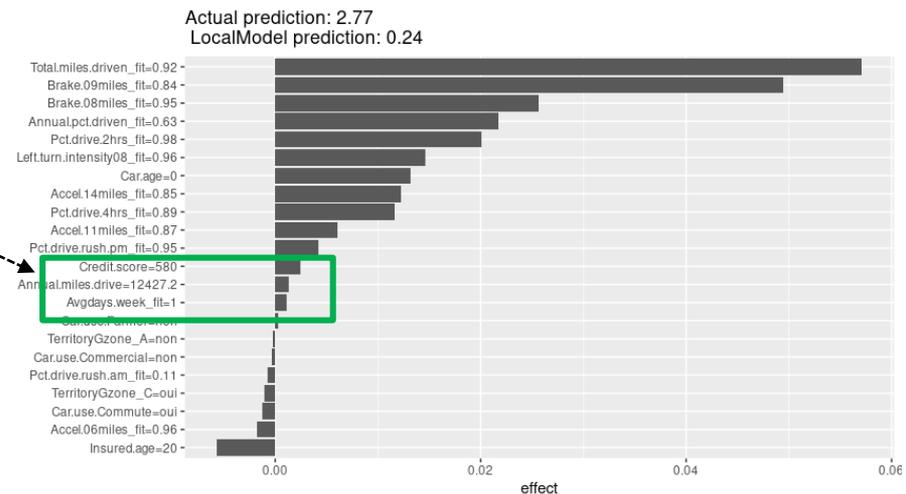
## 4 • CAS D'APPLICATION: ASSURANCE AUTOMOBILE

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

Explications LIME *initiales*



Explications LIME *attaquées*

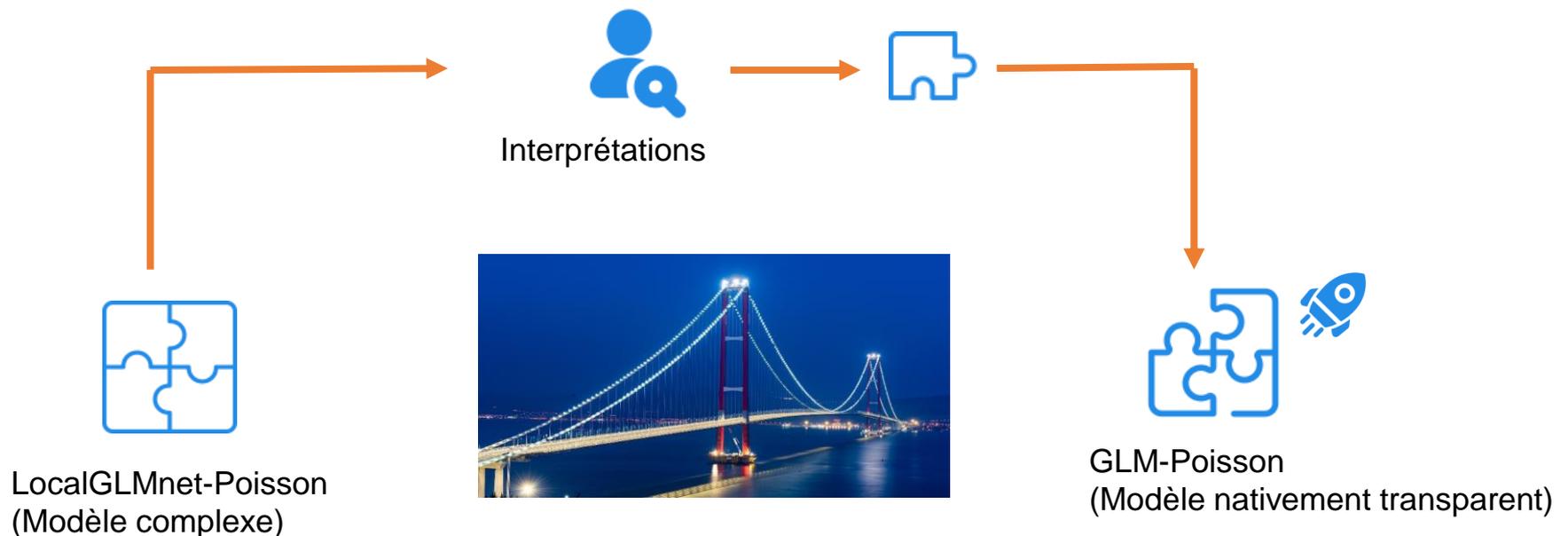


$$e(x) = \begin{cases} f(x), & \text{si } is\_OOD(x) \leq 0.8 \\ \psi(x), & \text{si } is\_OOD(x) > 0.8 \end{cases}$$

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

## 4 • CAS D'APPLICATION: ASSURANCE AUTOMOBILE

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



#### Observation

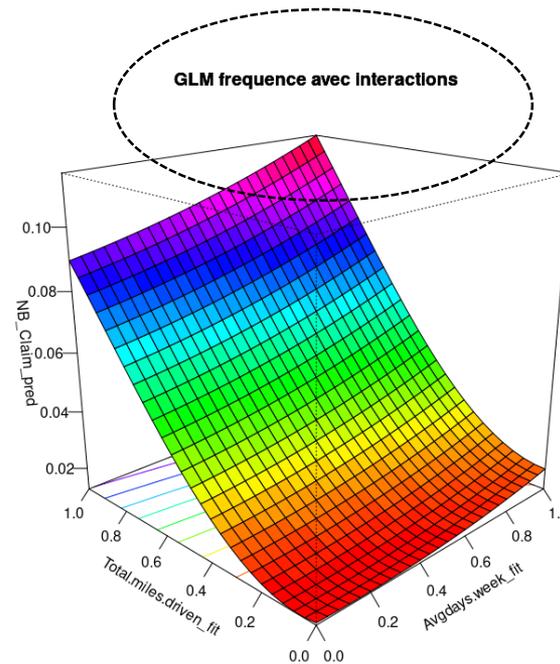
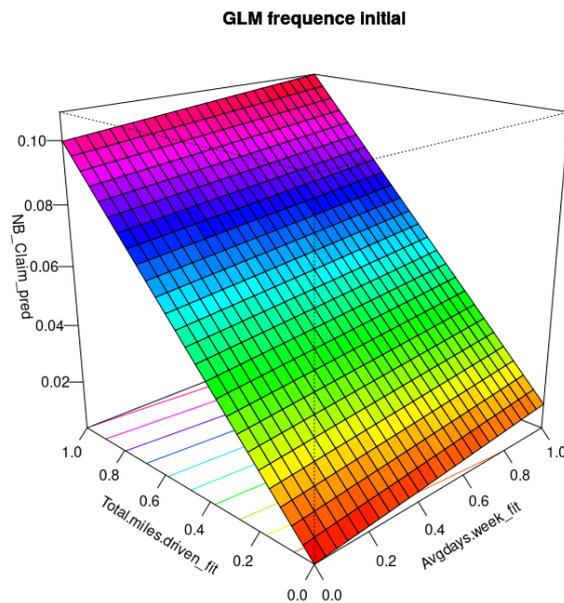
L'interprétabilité peut permettre de réaliser de l'ingénierie des caractéristiques

## 4 • CAS D'APPLICATION: ASSURANCE AUTOMOBILE

### 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.

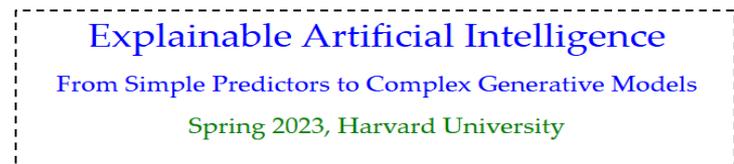


5

**C**ONCLUSION

### 5 • CONCLUSION

- ❑ L'interprétabilité **n'est pas un luxe technique**, mais une **exigence éthique** pour que la complexité algorithmique reste au service de l'intelligence humaine notamment dans des **domaines sensibles** comme celui des **assurances**.
- ❑ Mettre en place des explications par des **méthodes contrefactuelles**.
- ❑ La communauté scientifique se tourne davantage vers des **modèles hybrides**.



**THANK YOU**

**QUESTION ?**