

# Rapport mi-projet Fairness pour l'IA

Erwan LEMATTRE, Yannis CHUPIN

5 mars 2024

# Table des matières

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Découverte du jeu de données</b>	<b>3</b>
2.1	La base de données . . . . .	3
2.2	Répartition des données . . . . .	3
2.2.1	Catégories de véhicule . . . . .	3
2.2.2	Gravité de l'accident . . . . .	3
<b>3</b>	<b>Préparation des données</b>	<b>4</b>
<b>4</b>	<b>Analyse des données</b>	<b>4</b>
4.1	Analyse univarié des données . . . . .	5
4.1.1	Les accidents mortels . . . . .	5
4.1.2	Les piétons . . . . .	5
4.1.3	L'âge . . . . .	6
4.1.4	Le genre des conducteurs . . . . .	6
4.1.5	Le type de collision . . . . .	7
4.2	Analyse bivariée des données . . . . .	7
4.2.1	Le genre du conducteur . . . . .	7
<b>5</b>	<b>Apprentissage</b>	<b>7</b>
5.1	Split . . . . .	7
5.2	One Hot Encoding . . . . .	8
5.3	Random Tree Classifier . . . . .	8
5.4	GaussianNB . . . . .	9
<b>6</b>	<b>Audit du modèle</b>	<b>9</b>
6.1	Génération des contrefactuels avec Dice . . . . .	9
6.2	Expliquer le modèle avec les valeurs de Shapley . . . . .	9
<b>7</b>	<b>Conclusion</b>	<b>10</b>
<b>A</b>	<b>Les données de notre dataset</b>	<b>11</b>

# 1 Introduction

Ce projet a pour objectif d’analyser les accidents de la circulation routière afin de pouvoir dire à partir des données d’un véhicule accidenté si l’accident est mortel ou non. Les données sont des données libres mises à disposition par le *Ministère de l’intérieur et des Outre-Mer*. Le jeu de donnée correspond aux accidents de 2005 à 2022 en France. Nous allons dans une première partie analyser ces données afin d’extraire les informations utiles à l’apprentissage et de pouvoir repérer d’éventuelles sources de biais pour notre modèle.

Vous pouvez retrouver le code sur le GitHub du projet. Le fichier `main.ipynb` contient le code principale que nous allons suivre tout au long de ce rapport. Le fichier `utils.py` contient toutes les fonctions auxiliaires que nous utilisons dans le fichier principale.

## 2 Découverte du jeu de données

### 2.1 La base de données

La base de données est composée de plusieurs tables : *usagers*, *vehicules*, *lieux* et *caracteristiques*. Nous avons joint ces quatres parties pour obtenir un dataframe contenant une cinquantaine de colonnes. On peut retrouver l’ensemble des attributs dans l’annexe A

### 2.2 Répartition des données

Afin de pouvoir conserver les données utiles pour l’apprentissage nous avons analysé la répartition des différentes données dans notre dataframe. Nous avons ainsi pu faire différentes observations.

Voici quelques-unes d’entre elles qui nous sont ensuite utiles pour la préparation des données.

#### 2.2.1 Catégories de véhicule

La base de données nous donne beaucoup de catégories différentes. Nous avons cependant pu remarquer que la majorité des véhicules sont dans seulement 5 catégories.

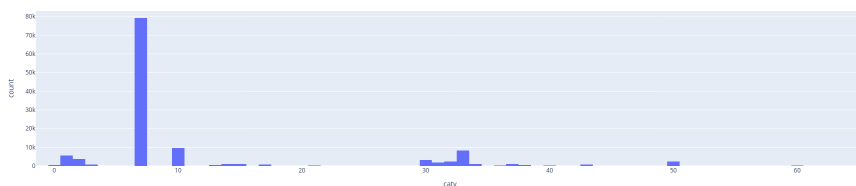


FIGURE 1 – Répartition des catégories de véhicules

#### 2.2.2 Gravité de l’accident

En affichant l’effectif d’accidents mortels nous avons pu remarquer qu’ils ne représentent qu’une infime partie des accidents. Le peu de données sur ces accidents ne nous permet pas l’apprentissage d’un modèle. C’est la raison pour laquelle nous avons décidé de nous intéresser non pas à la mortalité à l’échelle d’une personne mais plutôt à l’échelle d’un accident. Nous nous mettons pour cela au niveau d’un véhicule car cela nous permet de conserver plus d’informations

(à l'échelle d'un accident on aurait dû enlever trop d'informations pour conserver seulement les attributs plus généraux à l'accident).

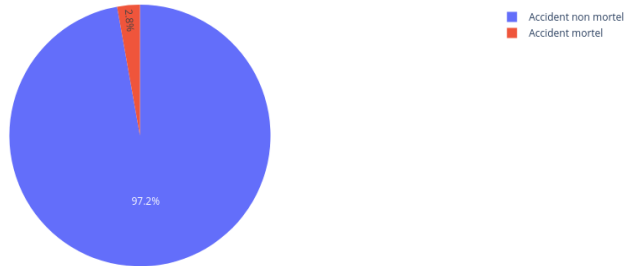


FIGURE 2 – Proportion d'accidents mortels

### 3 Préparation des données

À partir des observations précédentes, nous avons supprimé les attributs moins intéressants pour l'apprentissage et nous avons modifié certains attributs afin d'en extraire les informations intéressantes.

Les attributs supprimés sont : *voie*, *v1*, *v2*, *pr*, *pr1*, *lartpc*, *larroul*, *num\_veh*, *occutc*, *adr*, *senc*, *etatp*, *actp*, *manv*, *jour*, *com*, *hrmn*, *motor*, *place*, *vosp*, *locp*.

Nous avons effectué les modifications suivantes :

- Création d'un attribut *mortal* qui vaut 1 si le véhicule est impliqué dans un accident mortel, 0 sinon.
- À partir de l'attribut *sexe* nous avons créé un attribut *sexe\_conducteur* qui garde seulement le sexe du conducteur du véhicule.
- Création d'un attribut *piéton* qui vaut 1 si un piéton est impliqué dans l'accident, sinon 0.
- Nous avons utilisé l'année de naissance et l'année de l'accident pour récupérer l'âge du conducteur.
- L'attribut *vma* a été découpé en 4 catégories de vitesse.
- Pour les attributs *catv* et *vatr* nous avons gardé les valeurs les plus représentées dans la base de données.

Nous avons également réduit les valeurs de certains attributs. Par exemple pour des attributs avec des valeurs telles que *Non-renseigné*, *Autre* . . . Nous avons regroupé ces valeurs en une seule valeur. L'objectif était ici de simplifier en réduisant les catégories mais également d'améliorer les performances de notre modèle.

### 4 Analyse des données

Une fois nos données préparées, nous avons pu les visualiser. Nous allons montrer dans les deux prochaines parties les observations intéressantes que nous avons pu faire lors de l'analyse de

notre dataset.

## 4.1 Analyse univarié des données

### 4.1.1 Les accidents mortels

Une donnée intéressante à observer est la proportion de véhicules impliqués dans un accident mortel. C'est en effet la valeur que nous voulons prédire.

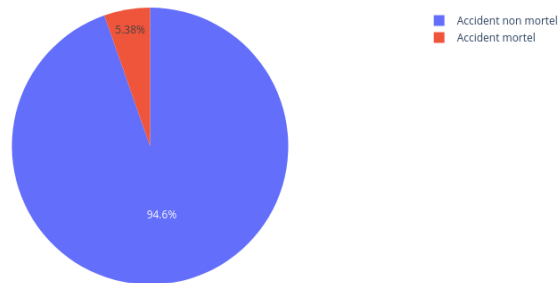


FIGURE 3 – Proportion des véhicules impliqués dans un accident mortel

Nous pouvons remarquer sur la figure 3 que le fait de s'intéresser aux véhicules impliqués dans un accident mortel et non plus aux personnes victimes permet de doubler le pourcentage. Même si cette proportion reste faible, cela va nous permettre d'avoir plus de données dans la catégorie mortel lors de l'apprentissage et par conséquent d'avoir un meilleur modèle.

### 4.1.2 Les piétons

Nous nous sommes ensuite intéressé aux accidents dans lesquels un piéton est impliqué. La figure 4 nous montre qu'un peu moins de 10% des accidents impliquent un piéton.

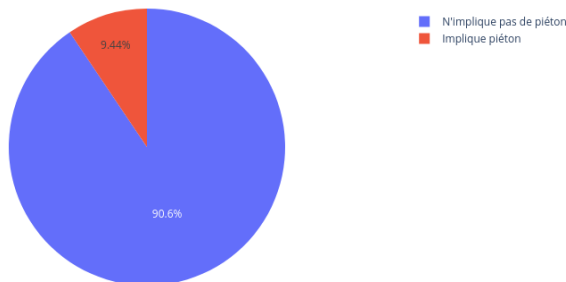


FIGURE 4 – Proportion des accidents avec piéton

### 4.1.3 L'âge

Nous pouvons visualiser l'âge des conducteurs via une boîte à moustache. La figure 5 nous montre la répartition de l'âge des conducteurs. Lors du prétraitement des données, les valeurs aberrantes ont été enlevée. On retrouve donc logiquement des âges contenus entre 0 et 100 ans. L'âge médian des conducteurs est 33 ans avec le premier quartile à 21 et le troisième quartile à 49 ans. Même si on peut imaginer que des valeurs sont fausses (il y a des conducteurs de moins de 16 ans), les valeurs sont tout de même assez cohérentes par rapport à ce que l'on pourrait imaginer de la répartition de l'âge des conducteurs.

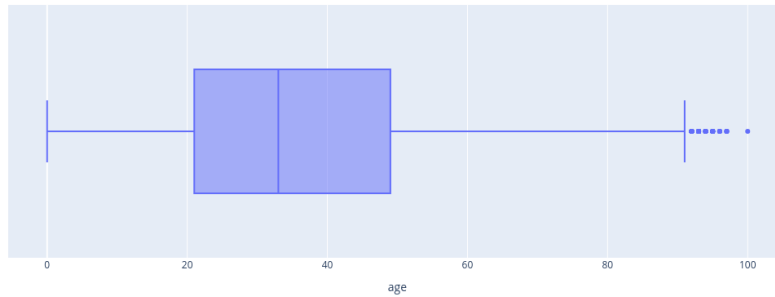


FIGURE 5 – Âge des conducteurs

### 4.1.4 Le genre des conducteurs

Le genre des conducteurs est assez intéressant à analyser. Sur la figure 6 nous pouvons remarquer différence importante entre le nombre de femme au volant (indice 0) et le nombre d'hommes au volant (indice 1). Cette différence pourrait être une source de biais pour notre modèle. En effet le fait qu'il y ait beaucoup plus de données d'accident avec des hommes ne signifie pas qu'il y a plus de chance d'avoir un accident si on est un homme. Cela signifie peut-être que la proportion d'homme au volant est plus élevée et donc qu'il y a plus d'accident avec un homme au volant car il y a plus d'homme au volant. Le risque ici est que notre modèle associe une homme à un accident mortel car il y a beaucoup plus d'accidents mortels avec un homme au volant.

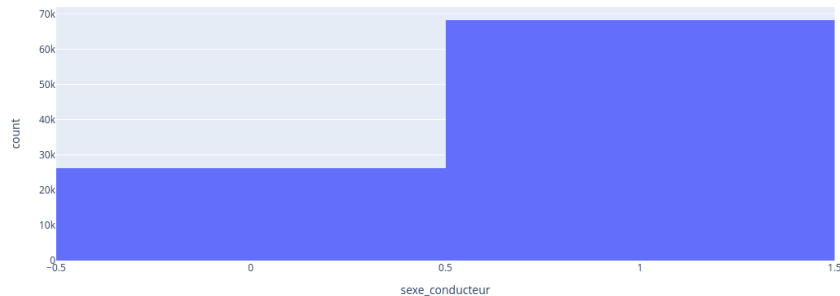


FIGURE 6 – Genre des conducteurs

#### 4.1.5 Le type de collision

La figure 7 montre la répartition des différents types de collisions dans notre dataset. On peut remarquer que tous les types de collisions sont plutôt bien représentés dans notre dataset. C'est un attribut qui pourra être assez intéressant pour l'apprentissage.

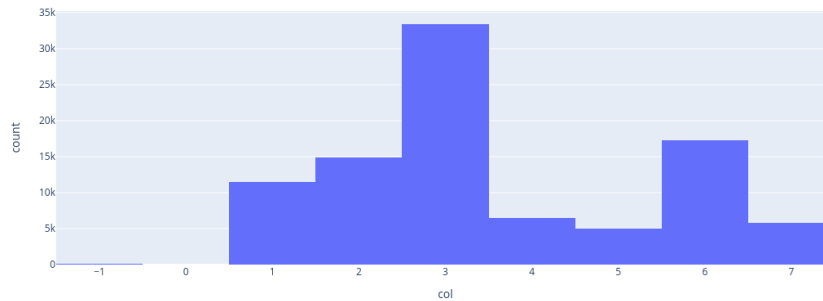


FIGURE 7 – Les types de collision

## 4.2 Analyse bivariée des données

Nous allons dans cette partie donner quelques exemples intéressants obtenus lors de l'analyse bivariée. On peut retrouver l'ensemble des graphiques observés dans le fichier `main.ipynb`.

### 4.2.1 Le genre du conducteur

Un attribut qu'il est intéressant d'analyser est le genre du conducteur. En effet il peut être source de biais s'il y a un déséquilibre entre homme et femme. On retrouve globalement la même proportion dans la corrélation que l'on soit homme ou femme. Les hommes étant beaucoup plus représentés dans le dataset, la proportion d'hommes est logiquement plus élevée. On peut cependant faire quelques remarques.

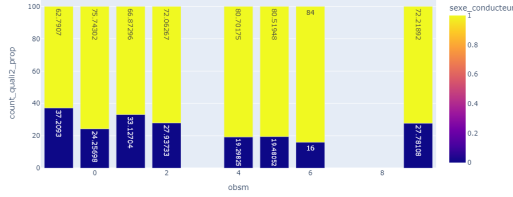
La figure 8a nous montre une proportion de femme moins élevée quand *obsm* vaut 6. La proportion de femmes est deux fois plus élevée quand *obsm* vaut 1. Ceci pourrait biaiser notre modèle.

Sur la figure 8b, on remarque également une proportion différente de femmes en fonction du type de véhicule. On pourrait expliquer cela par le fait que certains véhicules sont dans la réalité plus utilisés par les hommes, par exemple pourrait imaginer qu'il y a plus d'hommes qui conduisent des motos. Il faudra être vigilant car cela peut être source de biais. Il se peut que notre modèle associe une moto à un homme. Et que dans le cas d'une femme sur une moto le résultat soit soit forcément un accident mortel, soit forcément un accident non mortel.

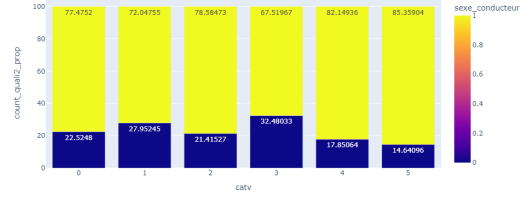
## 5 Apprentissage

### 5.1 Split

Séparer nos données en deux sets demandait de tenir compte d'une spécificité. Nous nous plaçons du point de vue d'un véhicule. Mais c'est bien souvent plusieurs véhicules qui sont



(a) Variation de la proportion de femme en fonction de la présence d'un obstacle dans l'accident



(b) Variation de la proportion de femme en fonction de la présence d'un obstacle dans l'accident

impliqués dans un accident. Il nous a donc fallu adapter le split pour que les véhicule d'un même accident soit dans le même set.

Par ailleurs, nous avons pu constater que ne pas adopter cette approche, provoquait un overfit lors du test. Typiquement, un trop grand nombre de mort était trouvés avec succès.

## 5.2 One Hot Encoding

La base de donnée nous fournit pour la plus part des attributs des données qui sont convertibles en entiers. Par exemple une catégorie de véhicule (*catv*) est dénomée par un chiffre. L'attribut cependant reste catégoriel. Dans les faits, Le seul attribut quantitatif qu'il nous reste est l'âge du conducteur. Le reste est soit binaire, soit qualitatif.

Tous les attributs qualitatifs ont donc été transformés en OneHot. Au départ nous souhaitions optimiser le nombre de colonnes en faisant les one hot par nous même. Typiquement, preonons surf (l'état de la surface de la route). Il est inutile de créer une colonne Non-renseigné, Autre ou normale. Nous ne sommes intéressés que par les états spécifiques de la route.

Cet approche cependant était source d'ennuis puisqu'elle empêchait la réalisation de l'audit. Nous avons donc utilisés la méthode standard pour nos oneHot.

## 5.3 Random Tree Classifier

Une fois les données prêtes, nous avons utilisés le classifieur d'arbres aléatoires pour obtenir un modèle. C'est le premier que nous avons utilisé. Pour ce qui est des résultats, l'accuracy score est de 1 pour l'entraînement et de 0.9014 pour le teste.

Pour ce qui est de la matrice de confusion, on remarque que 27687 véhicules ont été classés à raison en tant que accidents non létaux. Tandis que 291 accidents mortels ont été trouvés. Toutefois 1413 accidents mortel ont été resécés à tord comme étant non létaux et 1648 véhicule ont subit le sort inverse.

En résumé, on remaque que le modèle réussit très bien à classifier les accidents qui n'ont pas conduit à la mort de quelqu'un. Cemendant, il lui apparaît bien plus difficile de trouver les véhicules impliqués accidents qui sont mortels. On peut attribuer ceci au fait que ces accidents létaux ne représentent que 5.38% des accidents rétertoriés dans la base.



## 5.4 GaussianNB

# 6 Audit du modèle

## 6.1 Génération des contrefactuels avec Dice

Dice nous a permis d'établir quels attributs influent sur la prédiction de notre modèle. Les résultats sont cependant assez variants d'une exécution à l'autre. On peut tout de même retrouver lesquels sont fréquemment impliqués dans le changement des exemples contrefactuels. Parmi eux on peut noter que *âge* est souvent représenté. Le taux d'accidents (notamment mortels) étant plus élevé chez les jeunes et les personnes âgées, cela paraît assez cohérent. L'attribut *obsm* est également souvent représenté dans les exemples contrefactuels. On a pu remarquer que l'attribut *col* est peu représenté dans les résultats contrairement à l'hypothèse que nous avons pu faire lors de l'analyse.

Enfin les attributs *dep* et *sexe\_conducteur* sont également représentés. Ils pourraient être sources de biais, notamment l'attribut *dep*, le modèle pourrait associer un département à une prédiction d'accident mortel.

## 6.2 Expliquer le modèle avec les valeurs de Shapley

Afin d'analyser la contribution des différents attributs dans notre modèle, nous avons utilisé les valeurs de Shapley. Nous avons d'abord essayé de calculer la valeur exacte avec la fonction `ShapleyValues` du module *ShapKit*. Cependant le nombre d'attributs de notre dataset est trop important et le calcul trop long. Il a donc fallu calculer une approximation de la valeur de Shapley avec la fonction `MonteCarloShapley`. Nous devons cependant être prudents car le calcul de l'approximation nous donne une tendance pour notre modèle mais certainement pas une valeur exacte. On peut le constater en changeant le nombre d'itérations qui produit un résultat différent.

Avec `n_iter=1000` on obtient le diagramme en cascade de la figure 9. Ce diagramme nous montre en partant d'un score référence de 1 quels attributs contribuent à obtenir un score de 0 (l'utilisation d'un arbre de décision nous donne forcément une valeur binaire 0 ou 1 d'où le passage de 1 à 0).

Les attributs qui contribuent le plus sont ici *trajet*, *agg*, *circ*, *mois* et *obsm*. On aurait pu ici d'attendre à plus de contribution d'attributs comme *lum*, *vma* ou encore *col* (qui ici a une contribution mineure et dans le sens contraire des autres).

Une donnée intéressante ici est la valeur pour l'attribut *sexe\_conducteur*. On nous donne une contribution nulle de cet attribut ce qui peut nous encourager dans l'idée que le sexe n'est pas un biais pour notre modèle. Ce résultat étant une approximation, il est possible que la véritable valeur ne soit pas complètement nulle mais du moins assez faible.

Enfin, ce résultat nous montre que seule une petite partie des attributs contribue au résultat. On aurait voulu une contribution un peu plus répartie entre les attributs car ici le changement d'une petite partie des attributs pourrait faire basculer le résultat or on sait bien que les 5 attributs qui contribuent ici le plus ne permettent pas à eux seuls d'expliquer ce changement.

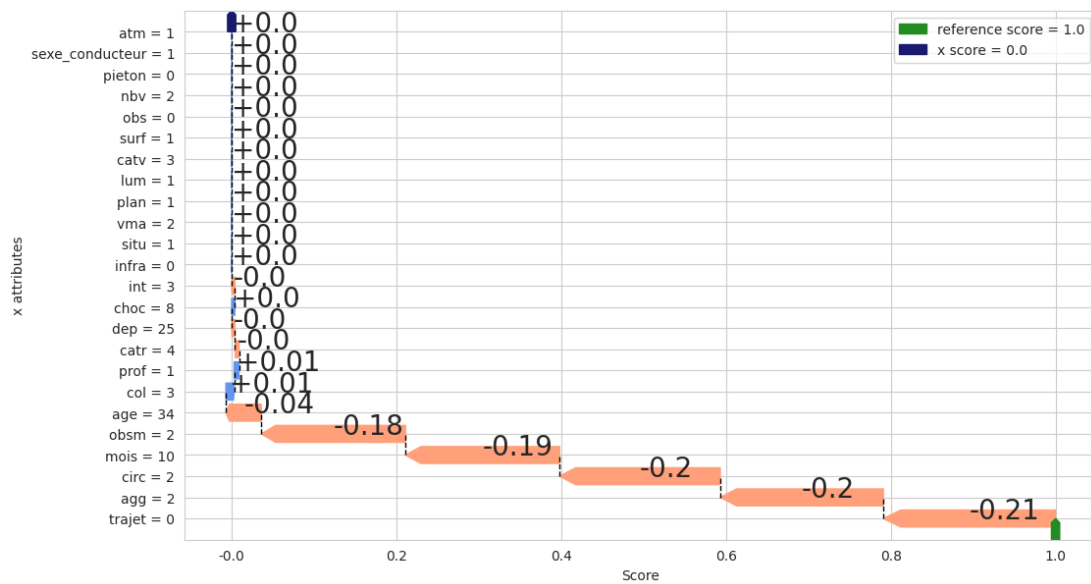


FIGURE 9 – Une approximation des valeurs de Shapley

## 7 Conclusion

## A Les données de notre dataset

Attribut	Description
<i>Num_Acc</i>	Numéro d'identifiant de l'accident
<i>jour mois</i>	Jour de l'accident, mois de l'accident
<i>an</i>	Année de l'accident.
<i>hrmn</i>	Heure et minutes de l'accident.
<i>lum</i>	Lumière : conditions d'éclairage dans lesquelles l'accident s'est produit.
<i>dep</i>	Département : Code INSEE.
<i>com</i>	Commune : Le numéro de commune est un code donné par l'INSEE.
<i>agg</i>	Localisation dans une agglomération ou non.
<i>int</i>	Si intersection : type d'intersection.
<i>atm</i>	Conditions atmosphériques.
<i>col</i>	Type de collision.
<i>adr</i>	Adresse postale : variable renseignée pour les accidents survenus en agglomération.
<i>lat</i>	Latitude.
<i>long</i>	Longitude.
<i>catr</i>	Catégorie de route (type de route).
<i>voie</i>	Numéro de la route.
<i>circ</i>	Régime de circulation.
<i>nbv</i>	Nombre total de voies de circulation.
<i>vosp</i>	Signale l'existence d'une voie réservée.
<i>prof</i>	Profil en long décrit la déclivité de la route à l'endroit de l'accident.
<i>plan</i>	Tracé en plan de la route.
<i>surf</i>	Etat de la surface de la route.
<i>infra</i>	Aménagement — Infrastructure s'il y en a.
<i>situ</i>	Situation géographique de l'accident.
<i>vma</i>	Numéro d'identifiant de l'accident.
<i>id_vehicule</i>	Identifiant du véhicule, foreign key.
<i>catv</i>	Catégorie du véhicule accidentée.
<i>obs</i>	Type d'obstacle heurté.
<i>obsm</i>	Type d'obstacle mobile heurté.
<i>choc</i>	Point de choc initial.
<i>manv</i>	Manoeuvre principale avant l'accident.
<i>catu</i>	Catégorie d'utilisateur (Conducteur, passager, piéton)
<i>grav</i>	Gravité de l'accident pour l'utilisateur.
<i>sexe</i>	Sexe du conducteur de la voiture (calculé dans le notebook).
<i>trajet</i>	Motif du déplacement au moment de l'accident.
<i>mortal</i>	Indique si le véhicule est impliqué dans un accident mortel. Attribut calculé dans le notebook.