



Classification de la gravité des accidents de la route

Présenté par :

BATTAL IKRAM

DAOUIRI SARA

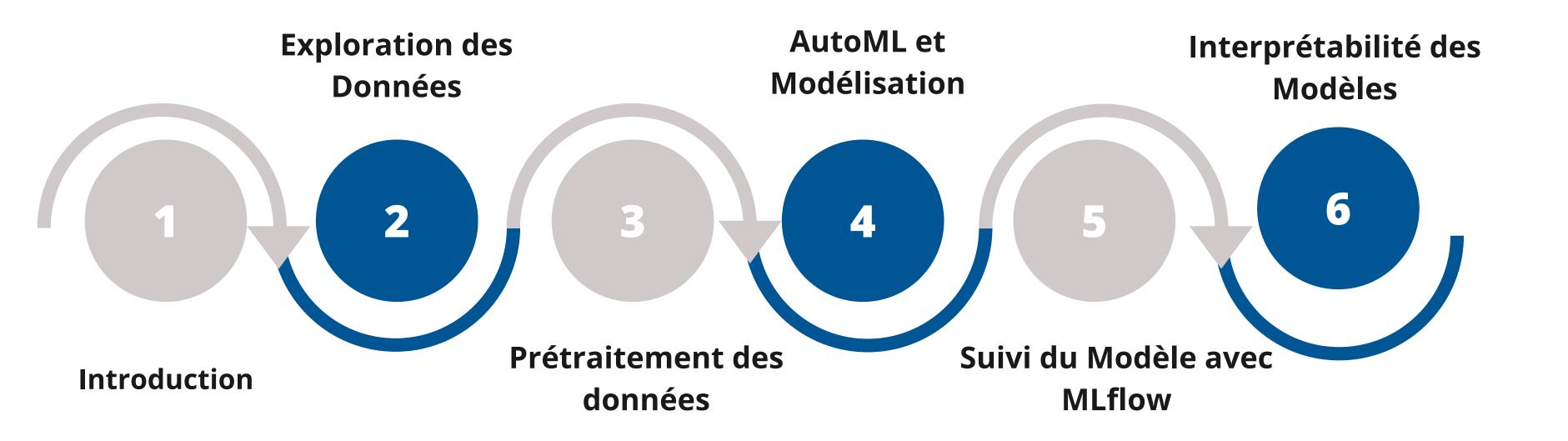
BENHABBACH DOHA

ECH-CHARFI CHAIMAE

Encadré par:

Feda Almuhisen

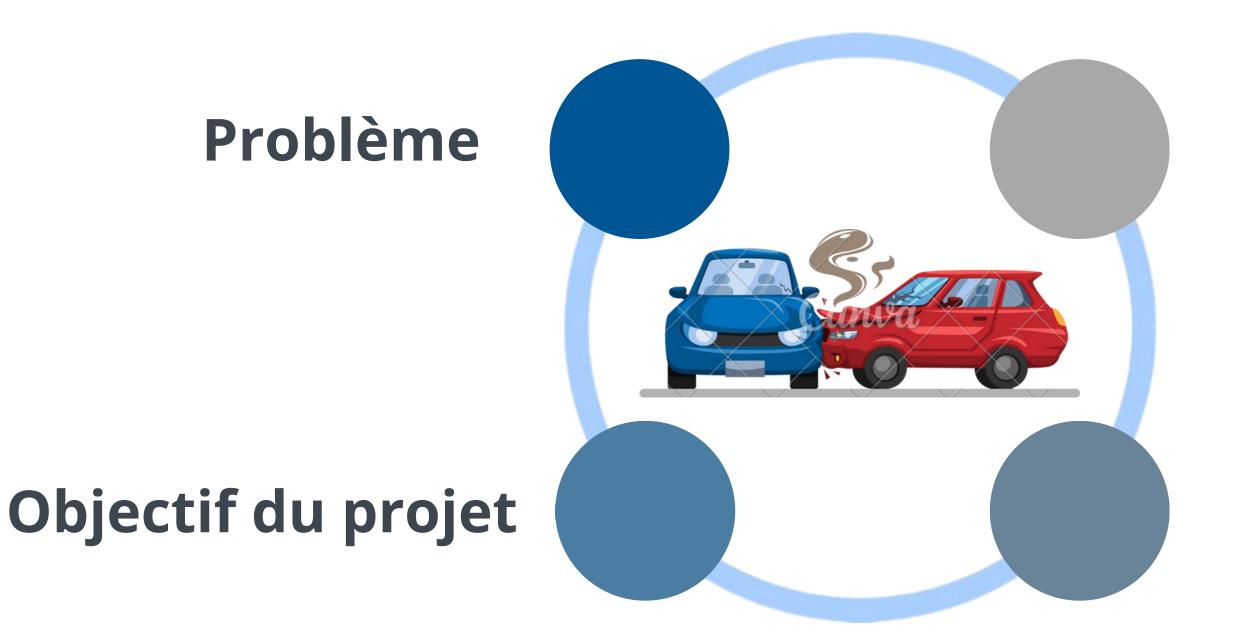
Plan



01 Introduction

01 Introduction

Contexte du Projet



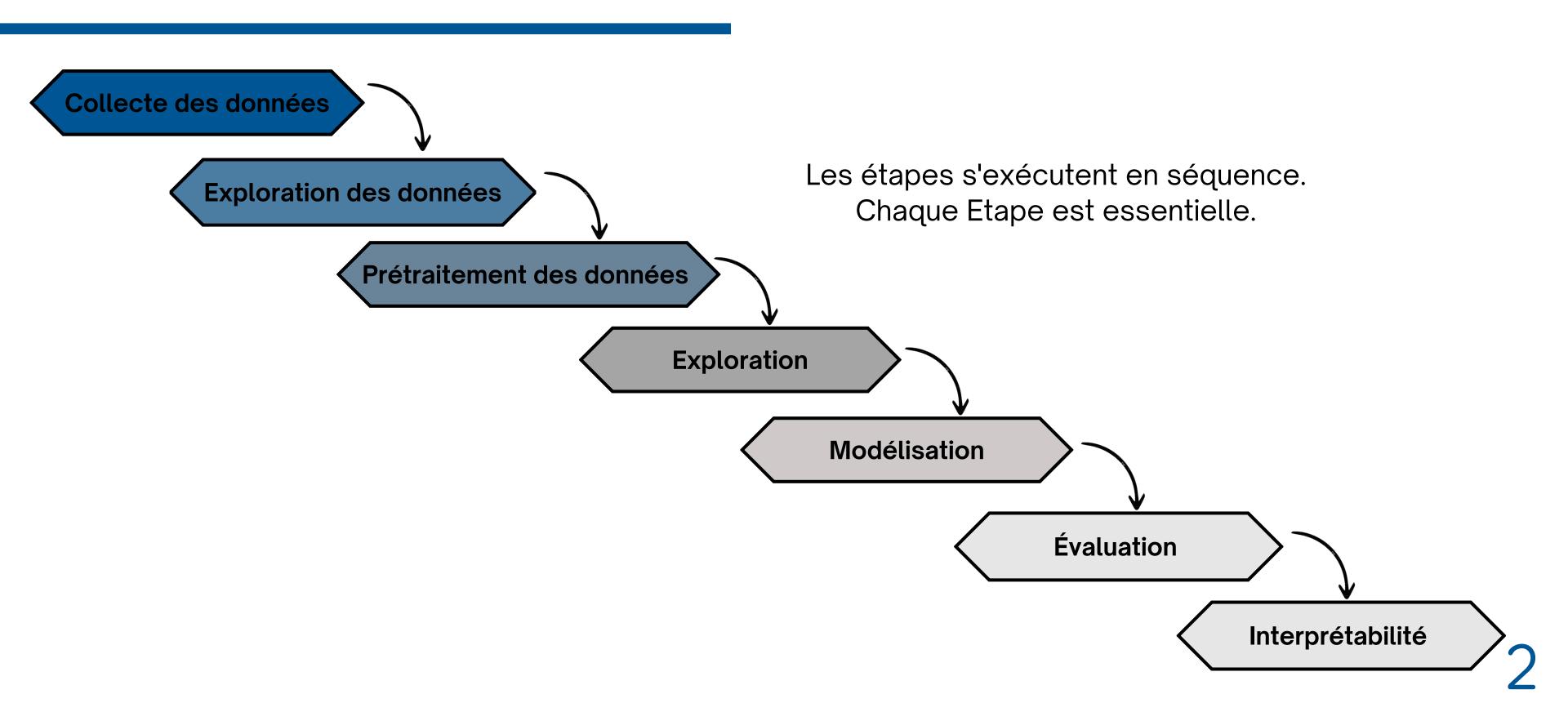
Source des données

Méthodologie

1

01 Introduction

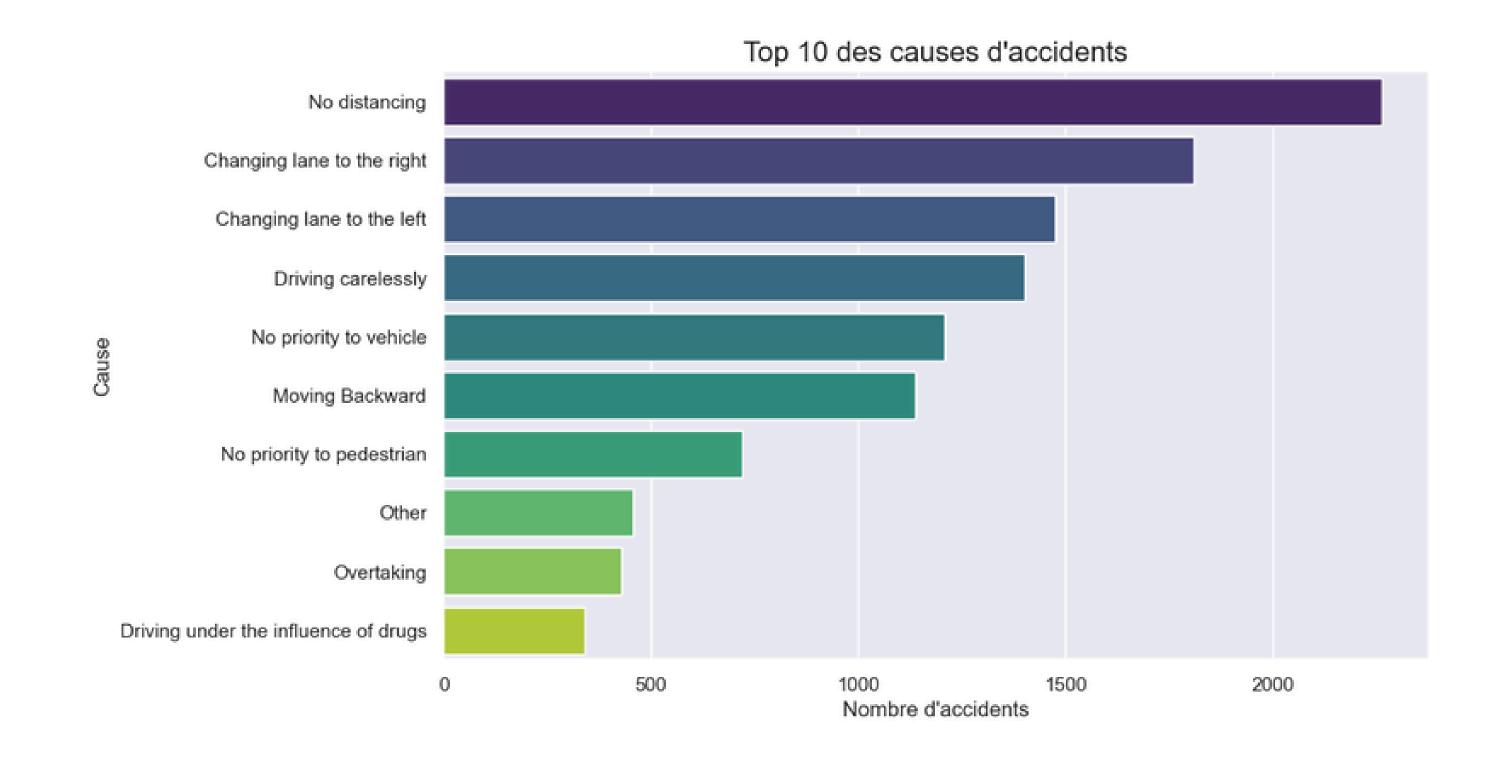
Méthodologie du Projet



Les causes des accidents

L'âge des conducteurs

BoxPlot

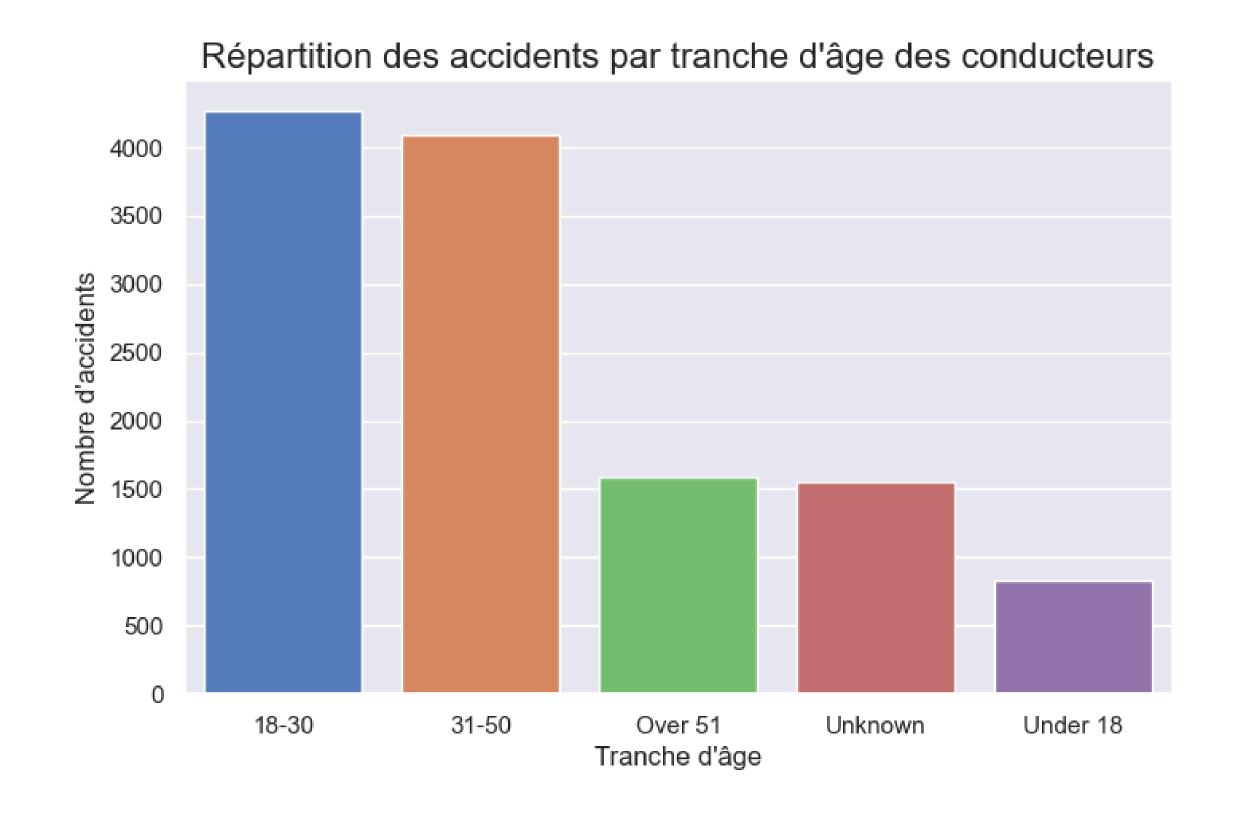


Les causes des accidents



L'âge des conducteurs

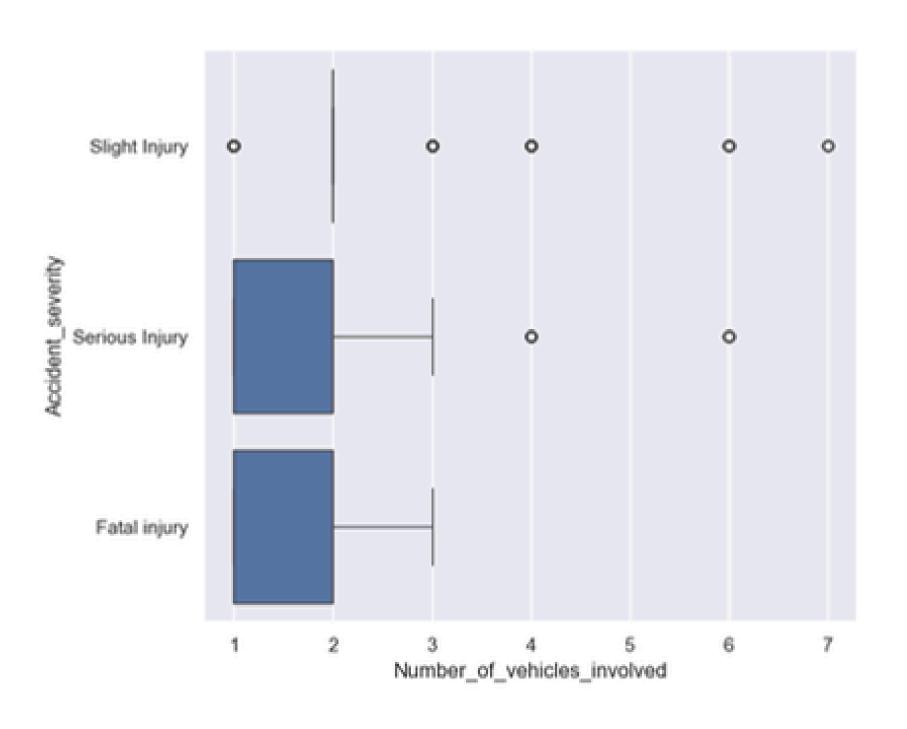
Boxplot

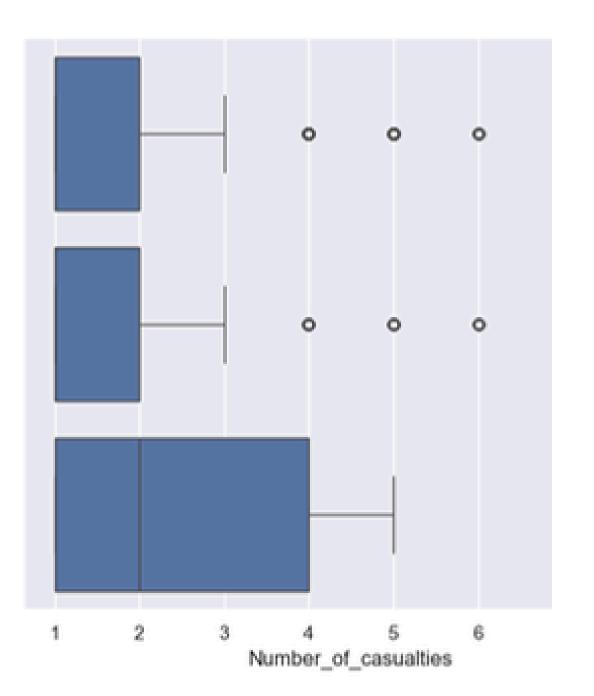


Les causes des accidents

L'âge des conducteurs







Standardisation

Valeurs manquantes

Encodage

Matrice de corrélation



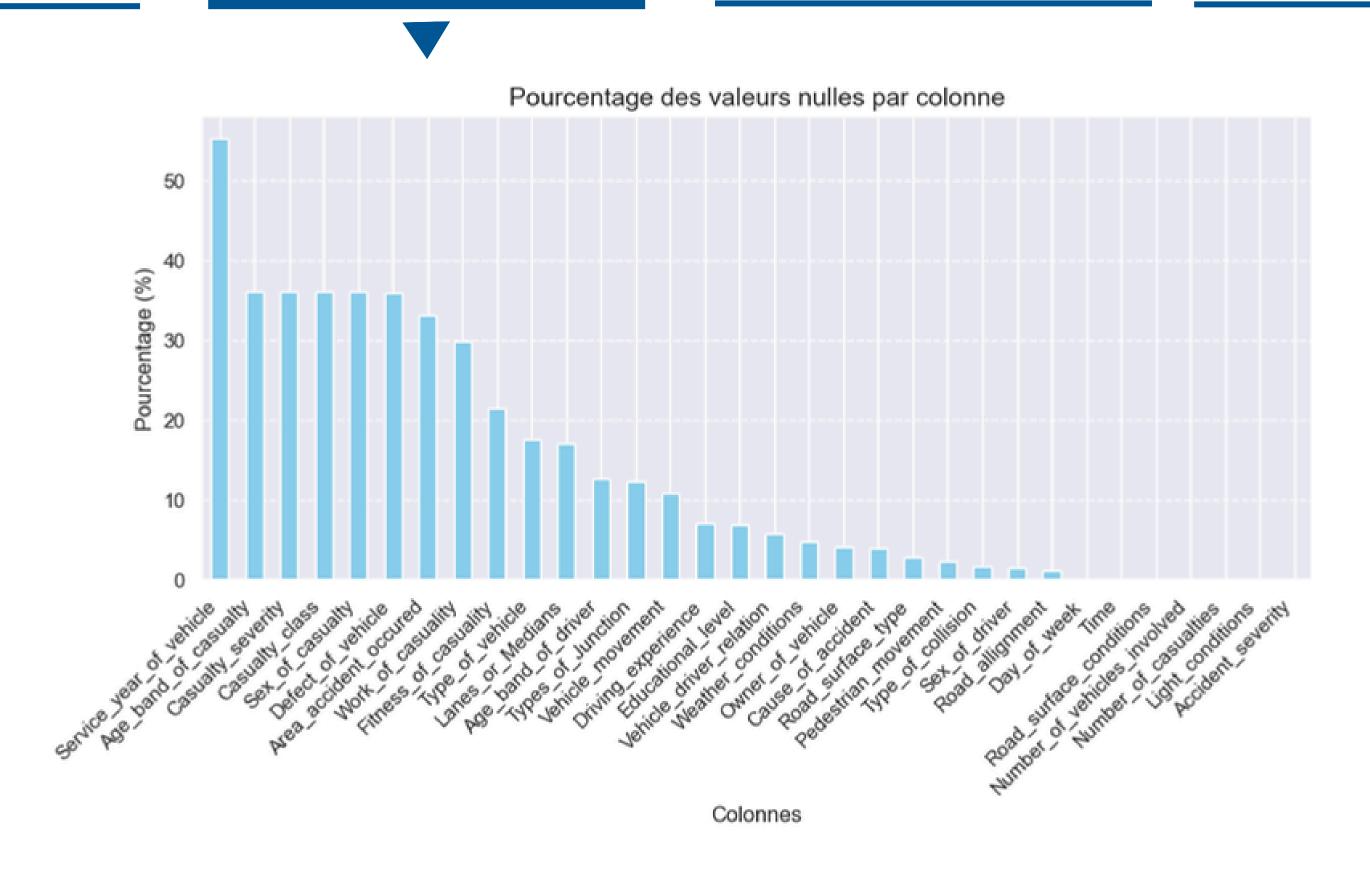
Données — Définition du format standard — Application des transformations

Standardisation

Valeurs manquantes

Encodage

Matrice de corrélation



Standardisation

Valeurs manquantes

Encodage

Matrice de corrélation



Colonne Time

Colonne Ordinale

Colonne Non Ordinale





Extraction des composants : heures, minutes, secondes



(heures * 3600) + (minutes * 60) + secondes]



Temps en secondes

	Time
0	61320
1	61320
2	61320
3	3960
4	3960
12311	58500
12312	64800
12313	50100
12314	50100
12315	50100
10016	

Standardisation

Valeurs manquantes

Encodage

Matrice de corrélation

Colonne Time

Colonnes Ordinales

Colonnes Non Ordinales

Colonne "age band driver"	Nouvelle Valeur
Under 18	1
18-30	2
31-50	3
Over 51	4

Standardisation

Valeurs manquantes

Encodage

Matrice de corrélation

Colonne Time

Colonne Ordinale

Colonne Non Ordinale

LabelEncoder

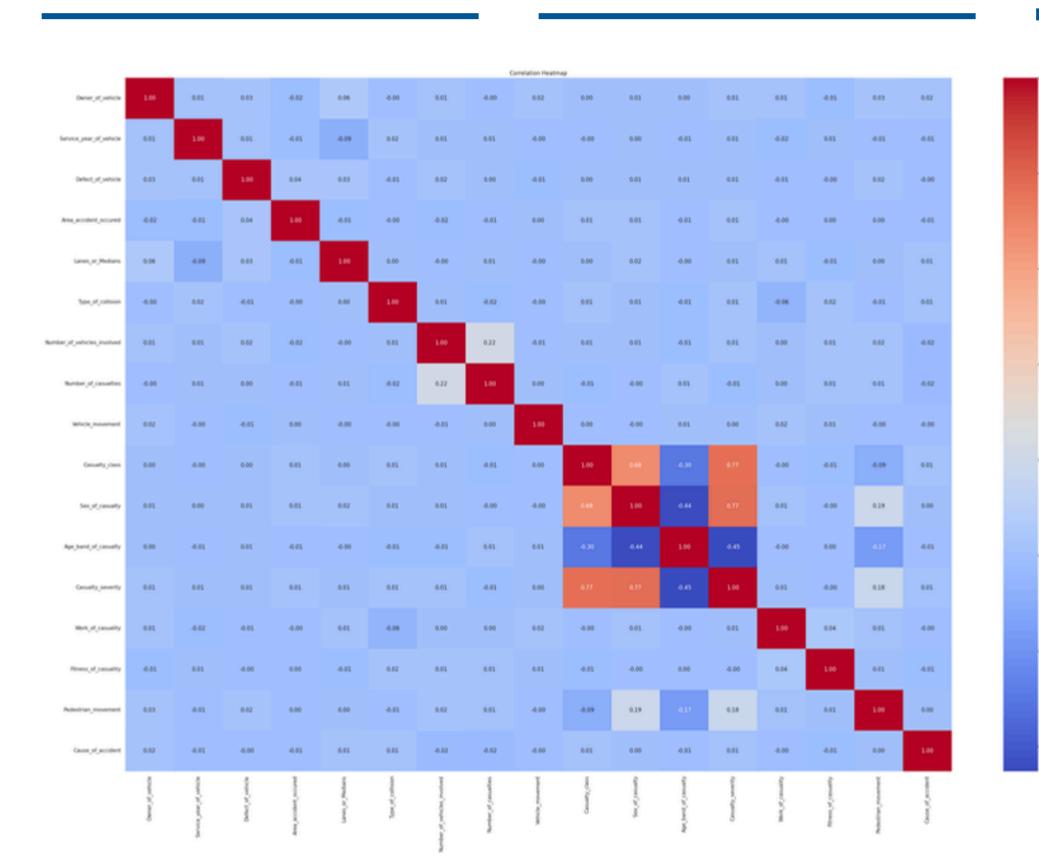
Colonne "Owner_of_vehicle"	Nouvelle Valeur
Owner	0
Governmental	1
Organization	2

Standardisation

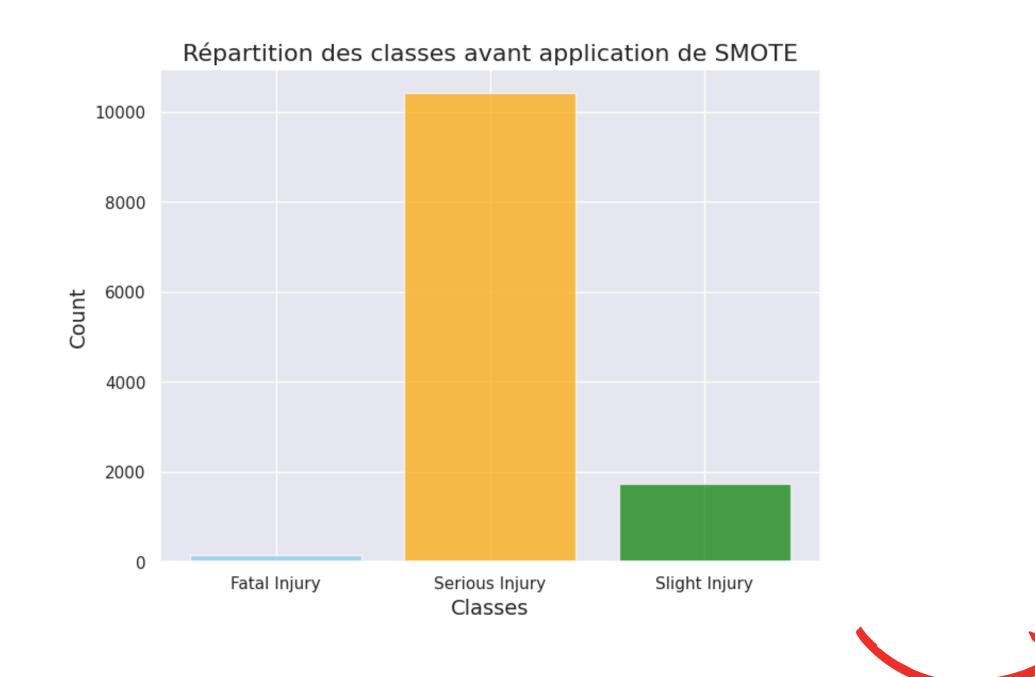
Valeurs manquantes

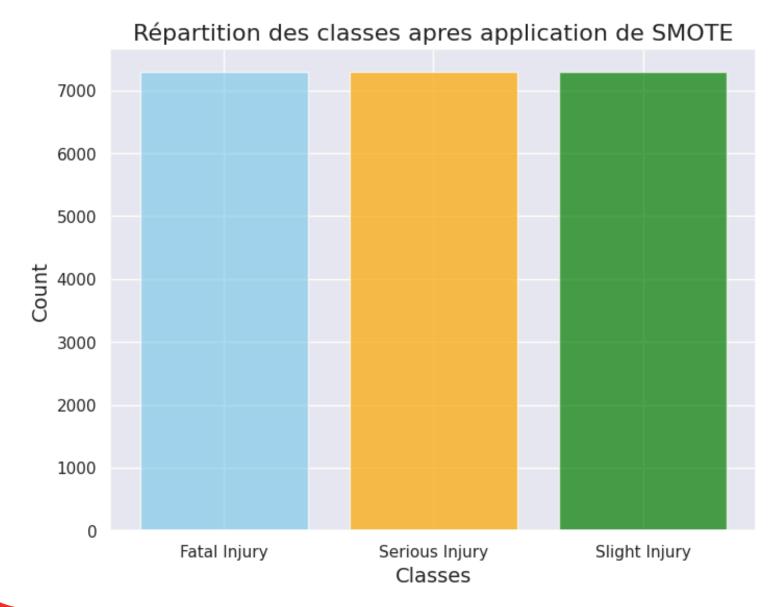
Encodage

Matrice de corrélation



Problème des classes déséquilibrées





04 AutoML et Modélisation

04 Exploration AutoML

PICARE

prétraitement de données

Données:

Nombre d'observations : 12 316

Nombre de caractéristiques : 25

Données avec valeurs manquantes : 8.2%

Division des données : Entraînement (8621) / Test (3695)

Prétraitement:

Imputation des valeurs manquantes :

Numériques : Moyenne

Catégorielles : Mode

123	Session id	0
Accident_severity	Target	1
Multiclass	Target type	2
(12316, 25)	Original data shape	3
(12316, 25)	Transformed data shape	4
(8621, 25)	Transformed train set shape	5
(3695, 25)	Transformed test set shape	6
24	Numeric features	7
8.2%	Rows with missing values	8
True	Preprocess	9
simple	Imputation type	10
mean	Numeric imputation	11
mode	Categorical imputation	12
StratifiedKFold	Fold Generator	13
10	Fold Number	14
-1	CPU Jobs	15
False	Use GPU	16
False	Log Experiment	17
clf-default-name	Experiment Name	18
4ee6	USI	19

04 Exploration AutoML

Comparaison des modèles

1- Meilleur modèle : LightGBM

Accuracy: **85.41%**

Precison: 82.94%

Recall: 85.41%

2- Deuxième modèle : Random Forest

Accuracy: **85.22**%

Precision: 82.77%

3- Troisième modèle : XGBoost

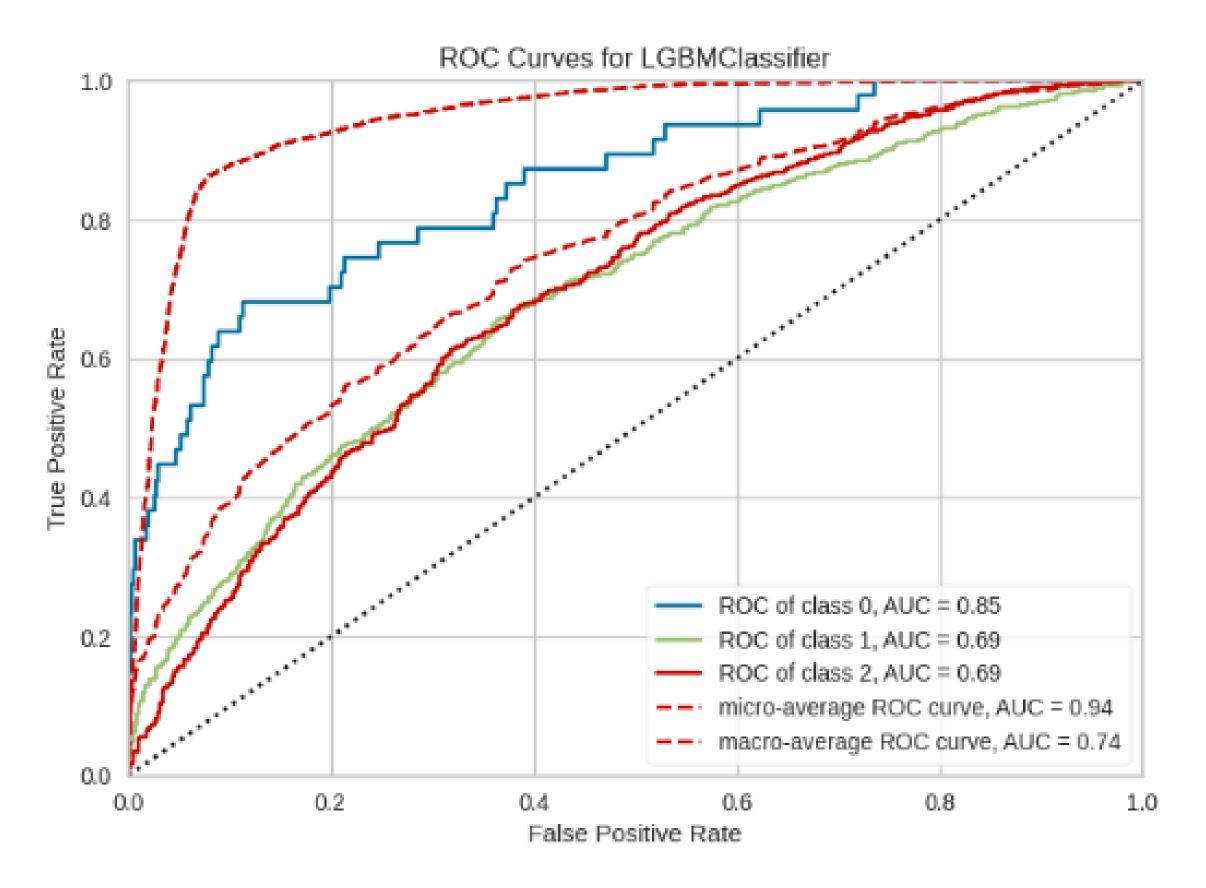
Accuracy: **85.05**%

Precision: 81.34%

Model	Accuracy	AUC	Recall	Prec.	F1	Карра	MCC	TT (Sec)
Light Gradient Boosting Machine	0.8541	0.6918	0.8541	0.8294	0.8018	0.1298	0.2183	5.0490
Random Forest Classifier	0.8522	0.6694	0.8522	0.8277	0.7940	0.0907	0.1860	1.3350
Extra Trees Classifier	0.8505	0.6524	0.8505	0.8229	0.7904	0.0729	0.1603	1.2580
Extreme Gradient Boosting	0.8505	0.6982	0.8505	0.8134	0.8042	0.1461	0.2096	0.5820
Gradient Boosting Classifier	0.8496	0.0000	0.8496	0.8167	0.7906	0.0772	0.1558	2.9910
Logistic Regression	0.8456	0.0000	0.8456	0.7151	0.7749	0.0000	0.0000	3.1570
Ridge Classifier	0.8456	0.0000	0.8456	0.7151	0.7749	0.0000	0.0000	0.0440
Dummy Classifier	0.8456	0.5000	0.8456	0.7151	0.7749	0.0000	0.0000	0.0370
Ada Boost Classifier	0.8455	0.0000	0.8455	0.7578	0.7774	0.0168	0.0566	0.3440
Linear Discriminant Analysis	0.8438	0.0000	0.8438	0.7148	0.7739	-0.0020	-0.0085	0.0500
Naive Bayes	0.8402	0.5821	0.8402	0.7145	0.7723	-0.0046	-0.0128	0.0460
K Neighbors Classifier	0.8291	0.6541	0.8291	0.7805	0.7945	0.1290	0.1488	0.1450
SVM - Linear Kernel	0.7625	0.0000	0.7625	0.6437	0.6975	0.0000	0.0000	0.3270
Decision Tree Classifier	0.7609	0.5841	0.7609	0.7770	0.7684	0.1563	0.1571	0.0890
Quadratic Discriminant Analysis	0.3682	0.0000	0.3682	0.5345	0.3986	0.0174	0.0216	0.0490
	Light Gradient Boosting Machine Random Forest Classifier Extra Trees Classifier Extreme Gradient Boosting Gradient Boosting Classifier Logistic Regression Ridge Classifier Dummy Classifier Ada Boost Classifier Linear Discriminant Analysis Naive Bayes K Neighbors Classifier SVM - Linear Kernel Decision Tree Classifier	Light Gradient Boosting Machine Random Forest Classifier Extra Trees Classifier Extreme Gradient Boosting Gradient Boosting Classifier Logistic Regression Classifier Dummy Classifier Ada Boost Classifier Discriminant Analysis Naive Bayes K Neighbors Classifier Decision Tree Classifier 0.8505 0.8505 0.8496 0.8496 0.8456 0.8456 0.8456 0.8456 0.8455 Linear Discriminant Analysis 0.8438 Naive Bayes 0.8402 K Neighbors Classifier 0.7625 Decision Tree Classifier 0.7609	Light Gradient Boosting Machine 0.8541 0.6918 Random Forest Classifier 0.8522 0.6694 Extra Trees Classifier 0.8505 0.6524 Extreme Gradient Boosting 0.8505 0.6982 Gradient Boosting Classifier 0.8496 0.0000 Logistic Regression 0.8456 0.0000 Ridge Classifier 0.8456 0.0000 Dummy Classifier 0.8456 0.5000 Ada Boost Classifier 0.8455 0.0000 Linear Discriminant Analysis 0.8438 0.0000 Naive Bayes 0.8402 0.5821 K Neighbors Classifier 0.8291 0.6541 SVM - Linear Kernel 0.7625 0.0000 Decision Tree Classifier 0.7609 0.5841	Light Gradient Boosting Machine 0.8541 0.6918 0.8541 Random Forest Classifier 0.8522 0.6694 0.8522 Extra Trees Classifier 0.8505 0.6524 0.8505 Extreme Gradient Boosting 0.8505 0.6982 0.8505 Gradient Boosting Classifier 0.8496 0.0000 0.8496 Logistic Regression 0.8456 0.0000 0.8456 Ridge Classifier 0.8456 0.0000 0.8456 Dummy Classifier 0.8456 0.5000 0.8456 Ada Boost Classifier 0.8455 0.0000 0.8455 Linear Discriminant Analysis 0.8438 0.0000 0.8438 Naive Bayes 0.8402 0.5821 0.8402 K Neighbors Classifier 0.8291 0.6541 0.8291 SVM - Linear Kernel 0.7625 0.0000 0.7625 Decision Tree Classifier 0.7609 0.5841 0.7609	Light Gradient Boosting Machine 0.8541 0.6918 0.8541 0.8294 Random Forest Classifier 0.8522 0.6694 0.8522 0.8277 Extra Trees Classifier 0.8505 0.6524 0.8505 0.8229 Extreme Gradient Boosting 0.8505 0.6982 0.8505 0.8134 Gradient Boosting Classifier 0.8496 0.0000 0.8496 0.8167 Logistic Regression 0.8456 0.0000 0.8456 0.7151 Ridge Classifier 0.8456 0.0000 0.8456 0.7151 Dummy Classifier 0.8456 0.5000 0.8456 0.7151 Ada Boost Classifier 0.8455 0.0000 0.8456 0.7578 Linear Discriminant Analysis 0.8438 0.0000 0.8438 0.7148 Naive Bayes 0.8402 0.5821 0.8402 0.7145 K Neighbors Classifier 0.8291 0.6541 0.8291 0.7805 SVM - Linear Kernel 0.7609 0.5841 0.7609 0.7770	Light Gradient Boosting Machine 0.8541 0.6918 0.8541 0.8294 0.8018 Random Forest Classifier 0.8522 0.6694 0.8522 0.8277 0.7940 Extra Trees Classifier 0.8505 0.6524 0.8505 0.8229 0.7904 Extreme Gradient Boosting 0.8505 0.6982 0.8505 0.8134 0.8042 Gradient Boosting Classifier 0.8496 0.0000 0.8496 0.8167 0.7906 Logistic Regression 0.8456 0.0000 0.8456 0.7151 0.7749 Ridge Classifier 0.8456 0.0000 0.8456 0.7151 0.7749 Dummy Classifier 0.8456 0.5000 0.8456 0.7151 0.7749 Ada Boost Classifier 0.8455 0.0000 0.8455 0.7578 0.7774 Linear Discriminant Analysis 0.8438 0.0000 0.8438 0.7148 0.7723 K Neighbors Classifier 0.8291 0.6541 0.8291 0.7805 0.6437 0.6975 SVM - Linear Kernel </td <td>Light Gradient Boosting Machine 0.8541 0.6918 0.8541 0.8294 0.8018 0.1298 Random Forest Classifier 0.8522 0.6694 0.8522 0.8277 0.7940 0.0907 Extra Trees Classifier 0.8505 0.6524 0.8505 0.8229 0.7904 0.0729 Extreme Gradient Boosting 0.8505 0.6982 0.8505 0.8134 0.8042 0.1461 Gradient Boosting Classifier 0.8496 0.0000 0.8496 0.8167 0.7906 0.0772 Logistic Regression 0.8456 0.0000 0.8456 0.7151 0.7749 0.0000 Ridge Classifier 0.8456 0.5000 0.8456 0.7151 0.7749 0.0000 Dummy Classifier 0.8456 0.5000 0.8456 0.7151 0.7749 0.0000 Ada Boost Classifier 0.8455 0.0000 0.8455 0.7578 0.7774 0.0168 Linear Discriminant Analysis 0.8402 0.5821 0.8402 0.7145 0.7723 -0.0046</td> <td>Light Gradient Boosting Machine 0.8541 0.6918 0.8541 0.8294 0.8018 0.1298 0.2183 Random Forest Classifier 0.8522 0.6694 0.8522 0.8277 0.7940 0.0907 0.1860 Extra Trees Classifier 0.8505 0.6524 0.8505 0.8229 0.7904 0.0729 0.1603 Extreme Gradient Boosting 0.8505 0.6982 0.8505 0.8134 0.8042 0.1461 0.2096 Gradient Boosting Classifier 0.8496 0.0000 0.8496 0.8167 0.7906 0.0772 0.1558 Logistic Regression 0.8456 0.0000 0.8456 0.7151 0.7749 0.0000 0.0000 Ridge Classifier 0.8456 0.5000 0.8456 0.7151 0.7749 0.0000 0.0000 Dummy Classifier 0.8456 0.5000 0.8456 0.7151 0.7749 0.0000 0.0000 Ada Boost Classifier 0.8455 0.0000 0.8456 0.7151 0.7749 0.0168 0.0566 </td>	Light Gradient Boosting Machine 0.8541 0.6918 0.8541 0.8294 0.8018 0.1298 Random Forest Classifier 0.8522 0.6694 0.8522 0.8277 0.7940 0.0907 Extra Trees Classifier 0.8505 0.6524 0.8505 0.8229 0.7904 0.0729 Extreme Gradient Boosting 0.8505 0.6982 0.8505 0.8134 0.8042 0.1461 Gradient Boosting Classifier 0.8496 0.0000 0.8496 0.8167 0.7906 0.0772 Logistic Regression 0.8456 0.0000 0.8456 0.7151 0.7749 0.0000 Ridge Classifier 0.8456 0.5000 0.8456 0.7151 0.7749 0.0000 Dummy Classifier 0.8456 0.5000 0.8456 0.7151 0.7749 0.0000 Ada Boost Classifier 0.8455 0.0000 0.8455 0.7578 0.7774 0.0168 Linear Discriminant Analysis 0.8402 0.5821 0.8402 0.7145 0.7723 -0.0046	Light Gradient Boosting Machine 0.8541 0.6918 0.8541 0.8294 0.8018 0.1298 0.2183 Random Forest Classifier 0.8522 0.6694 0.8522 0.8277 0.7940 0.0907 0.1860 Extra Trees Classifier 0.8505 0.6524 0.8505 0.8229 0.7904 0.0729 0.1603 Extreme Gradient Boosting 0.8505 0.6982 0.8505 0.8134 0.8042 0.1461 0.2096 Gradient Boosting Classifier 0.8496 0.0000 0.8496 0.8167 0.7906 0.0772 0.1558 Logistic Regression 0.8456 0.0000 0.8456 0.7151 0.7749 0.0000 0.0000 Ridge Classifier 0.8456 0.5000 0.8456 0.7151 0.7749 0.0000 0.0000 Dummy Classifier 0.8456 0.5000 0.8456 0.7151 0.7749 0.0000 0.0000 Ada Boost Classifier 0.8455 0.0000 0.8456 0.7151 0.7749 0.0168 0.0566

04 Exploration AutoML

Courbe ROC AUC



Processus de modélisation



Test des modèles

Random Forest Xgboost LightGBM 2

Optimisation des hyperparamètres

> GridSearch Cross Validation

3

Evaluation des modèles

Courbe ROC AUC

Matrice de confusion



Choix du modèle

Meilleure Accuracy

Choix des hyperparamètres

Random Forest:

- n_estimator : 500
- max_depth : 20
- min_samples_split: 10

LightGBM et Xgboost :

- n_estimators : 500
- max_depth : 20
- learning_rate : 0.1
- subsample: 1.0

04 Modélisation

Test des modèles

Evaluation

Modèle	F1-Score	Accuracy	Recall	Precision
Random Forest	0,79	0,82	0,82	0,78
Xgboost	0,80	0,83	0,83	0.79
LightGBM	0,80	0.83	0,83	0,79

	precision	recall	f1-score
0 1	0.50 0.38	0.11 0.17	0.18 0.24
2	0.86	0.95	0.91
	precision	recall	f1-score
0	0.91	0.21	0.34
1	0.36	0.18	0.24

0.95

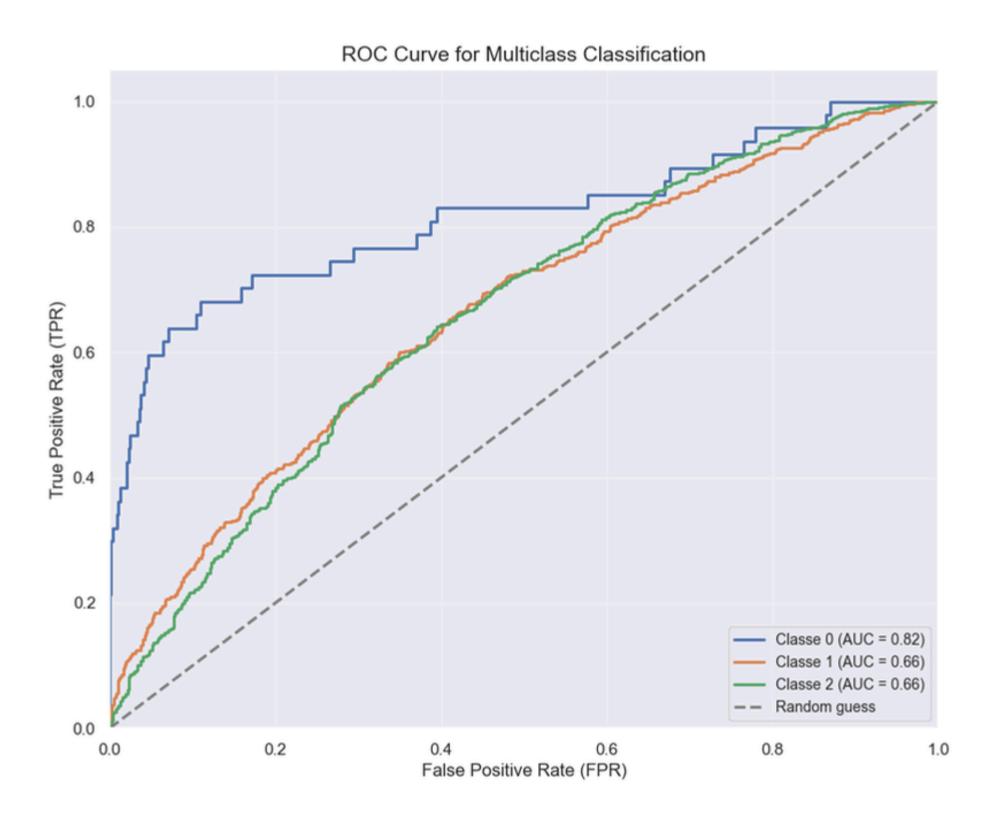
0.90

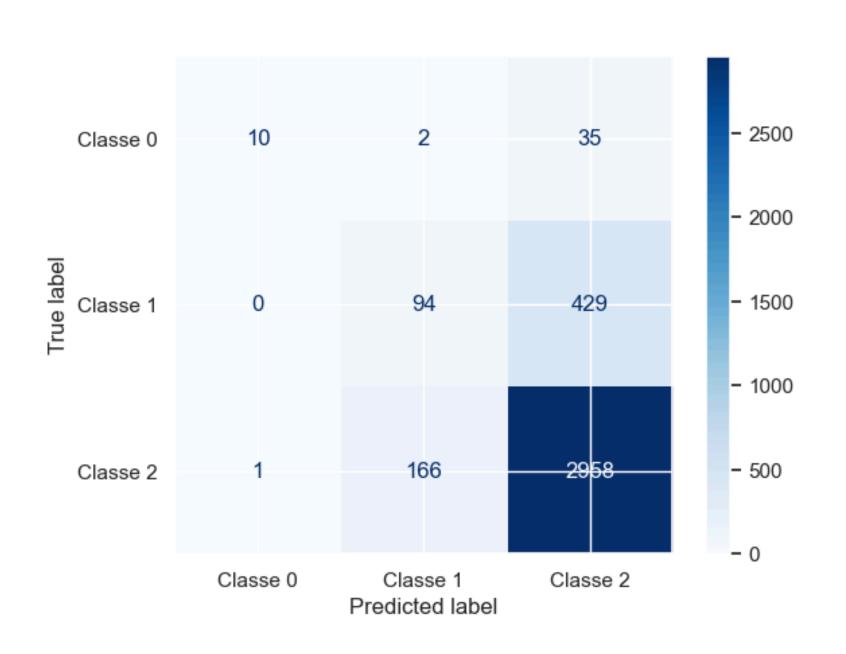
0.86

04 Modélisation

Test des modèles

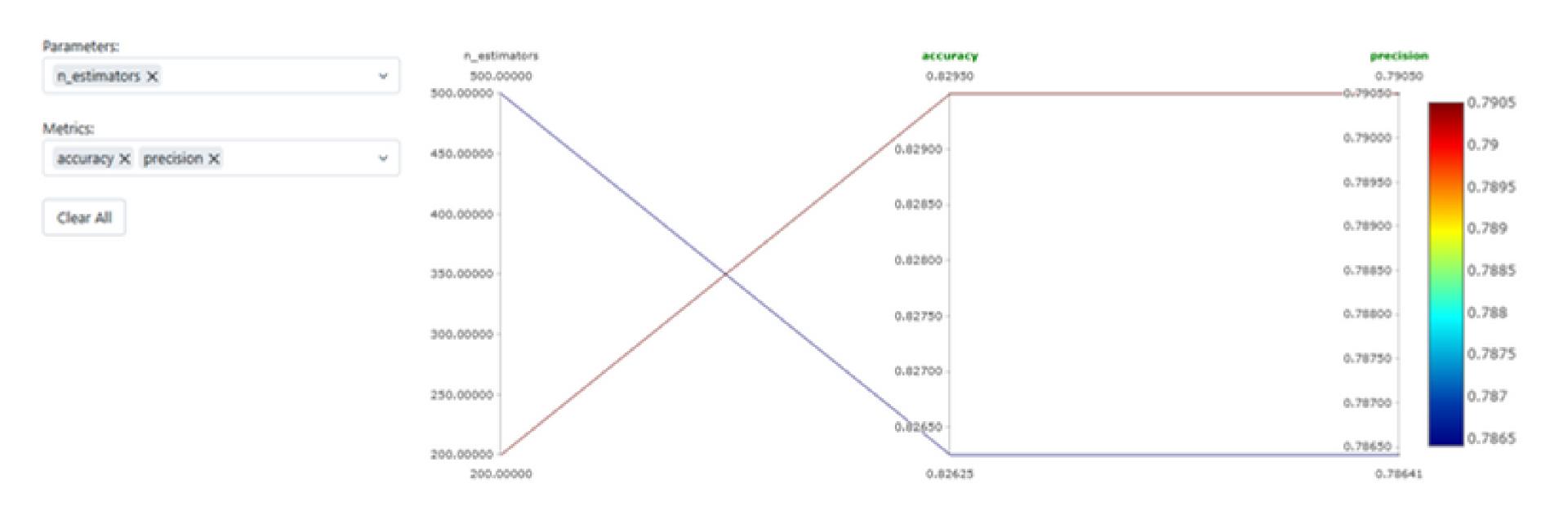




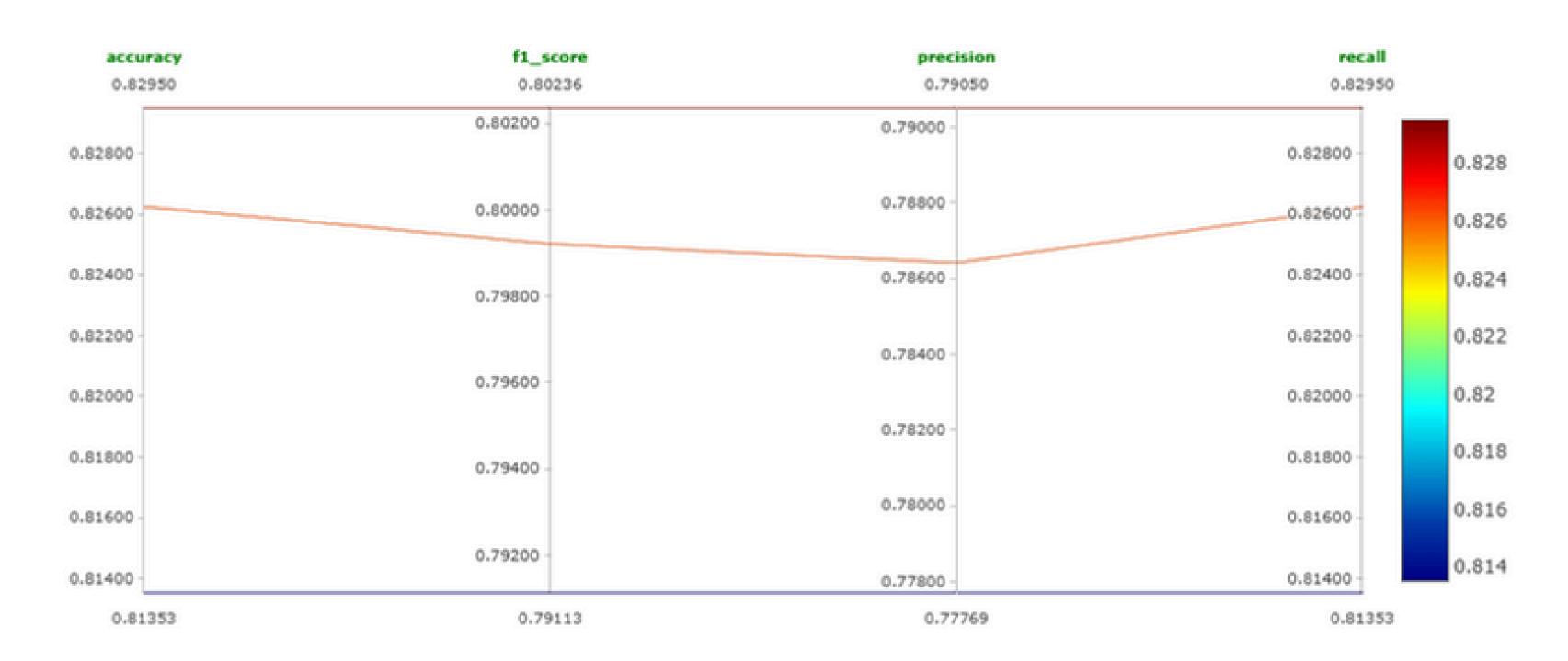




Comparaison des runs



Comparaison des runs



Déploiement du modèle en API REST locale

- MLflow models serve
- Bibliothèque requests

Résultat:

```
else:
    print(f"{response.status_code}: {response.text}")

L'accident est classé en : Serious Injury
```

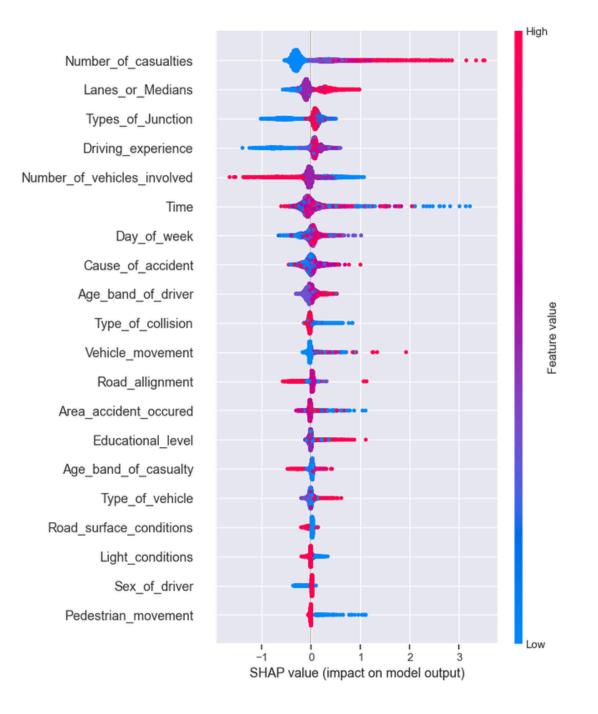
```
data = {
    'Time': 61320,
    'Day of week': 'Monday',
    'Age band of driver': '18-30'
    'Educational level': 'Above high school',
    'Driving experience': '2-5yr',
    'Light conditions': 'Daylight',
    'Sex of driver': 'Male',
    'Vehicle_driver_relation': 'Employee',
    'Type of vehicle': 'Automobile',
    'Owner_of_vehicle': 'Owner',
    'Service year of vehicle': 'Above 10yr',
    'Defect_of_vehicle': 'No defect',
    'Area_accident_occured': 'Residential areas',
    'Lanes or Medians': 'Undivided Two way',
    'Road_allignment': 'Flat terrain',
    'Types of Junction': 'No junction',
    'Road surface type': 'Asphalt roads',
    'Road surface conditions': 'Dry',
    'Weather conditions': 'Normal',
    'Type_of_collision': 'Vehicle with vehicle collision',
    'Number_of_vehicles_involved': 2,
    'Number_of_casualties': 2,
    'Vehicle movement': 'Going straight',
    'Casualty class': 'Driver or rider',
    'Sex of casualty': 'Male',
    'Age band of casualty': '31-50',
    'Casualty_severity': 'Slight Injury',
    'Work of casuality': 'Driver',
    'Fitness of casuality': 'Normal',
    'Pedestrian_movement': 'Not a Pedestrian',
    'Cause of accident': 'Moving Backward'
# Encodage des données
encoded data = encode data(data)
# Préparation des données pour L'API
request_data = {
    "dataframe records": [encoded data]
url = "http://127.0.0.1:1234/invocations"
headers = {"Content-Type": "application/json"}
```

O6 Interprétabilité du Modèle

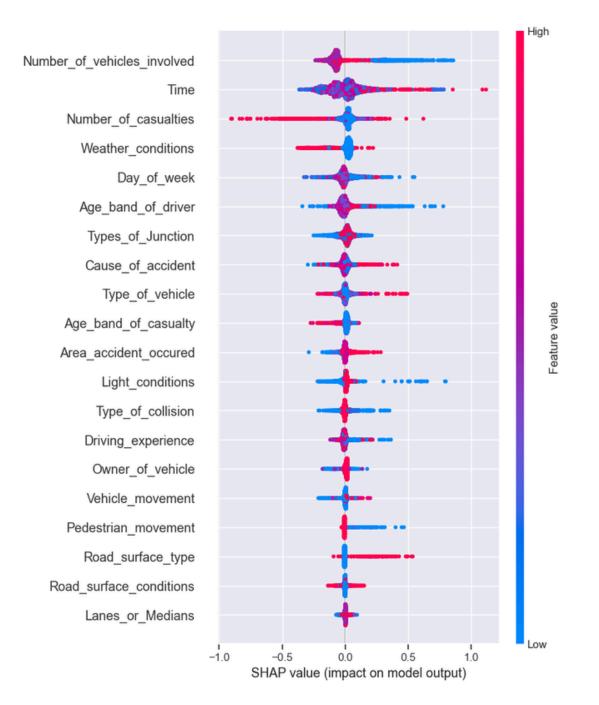
06 Interprétabilité du Modèle

Explication du modèle avec SHAP

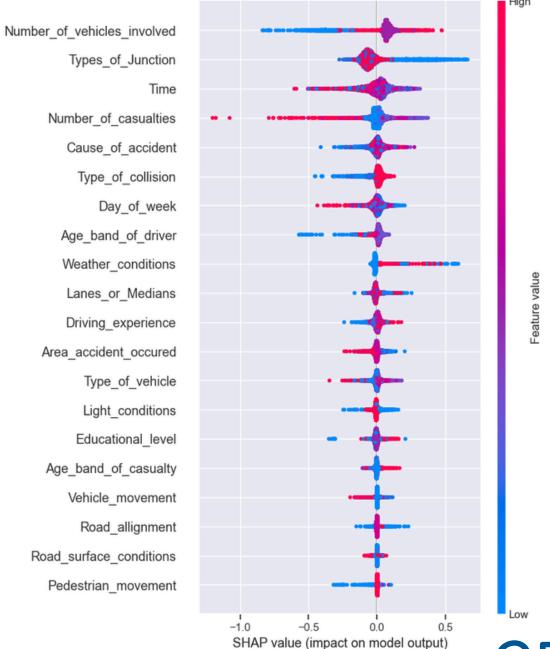
Classe 0 : Fatal Injury



Classe 1 : Serious Injury



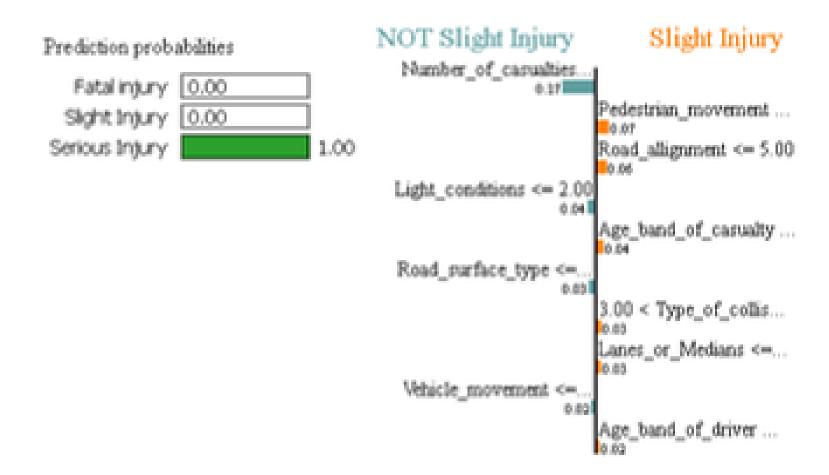
Classe 2 : Slight Injury



25

06 Interprétabilité du Modèle

Explication du modèle avec LIME



Feature	Value
Number_of_casualties	4.00
Pedestrian_movement	5.00
Road_allignment	5.00
Light_conditions	2.00
Age_band_of_casualty	0.00
Road_surface_type	0.00
Type_of_collision	6.00
Lanes_or_Medians	2.00
Vehicle_movement	2.00
Age_band_of_driver	2.00

O7 Conclusion

Merci de votre attention