

# DEVOIR\_S2 DATA\_MINING

SABAYE Fried-Junior, Caleb KASHALA, Mohamad EL KAWASS, Hicham AZOUD

## Table des matières

<b>1</b>	<b>Contexte et problématique</b>	<b>2</b>
1.1	Méthodologie . . . . .	2
<b>2</b>	<b>Arbres</b>	<b>3</b>
2.1	Arbre complet . . . . .	3
2.2	Arbre élagué . . . . .	5
<b>3</b>	<b>Bagging et Random forest</b>	<b>7</b>
3.1	Random forest . . . . .	7
3.2	Comparaison bagging et Random forest . . . . .	9
3.3	Application des modèles à l'échantillon Test . . . . .	9
<b>4</b>	<b>Modèle randomForest automatisé</b>	<b>10</b>
4.1	Application à l'échantillon test . . . . .	10
<b>5</b>	<b>Boosting</b>	<b>11</b>
5.1	Application à l'échantillon test . . . . .	12
<b>6</b>	<b>Scoring</b>	<b>13</b>
6.1	Variable à expliquer . . . . .	13
6.2	Analyse exploratoire . . . . .	13
6.3	Construction des modèles . . . . .	14
6.4	Validation des modèles : Indicateurs de qualité et de robustesse . . . . .	16
6.5	Application à l'échantillon test . . . . .	18
<b>7</b>	<b>Choix du meilleur modèle</b>	<b>21</b>
<b>8</b>	<b>Discussion</b>	<b>21</b>

# 1 Contexte et problématique

Cette étude fait suite à l'analyse factorielle effectuée au **semestre précédent**.

La base de données renseigne sur diverses variables et paramètres concernant des accidents de vélo.

Les dix premières lignes et colonnes de la base de données figurent dans le tableau suivant :

	Bike_Age	Bike_Alc_D	Bike_Dir	Bike_Injur	Bike_Pos	Bike_Race	Bike_Sex	Developmen	Drvr_Alc_D	Drvr_Injur
1	6	No	Not Applicable	Non	Driveway / Alley	Black	Female	Residential	No	O: No Injury
2	51	No	With Traffic	Non	Travel Lane	Black	Male	Commercial	No	O: No Injury
3	10	No	With Traffic	Oui	Travel Lane	Black	Male	Residential	No	O: No Injury
6	52	No	With Traffic	Oui	Travel Lane	White	Male	Commercial	No	O: No Injury
9	6	No	Facing Traffic	Oui	Travel Lane	White	Male	Residential	No	O: No Injury
12	30	No	With Traffic	Oui	Travel Lane	Black	Male	Commercial	No	O: No Injury
13	17	No	With Traffic	Oui	Travel Lane	White	Male	Residential	No	O: No Injury
14	20	No	With Traffic	Oui	Travel Lane	White	Male	Residential	No	O: No Injury
15	14	No	Facing Traffic	Non	Travel Lane	White	Male	Residential	No	B: Evident Injury
17	19	No	With Traffic	Oui	Travel Lane	Black	Male	Residential	No	O: No Injury

Le but de cette étude est de prédire la variable : *Bike\_Injur*. Pour cela nous allons utiliser diverses méthodes d'apprentissage supervisé.

TABLE 1 – Effectif de la variable à prédire

Bike_Injur	Nombre
Non	2725
Oui	2991

## 1.1 Méthodologie

Nous allons séparer la base de données en deux échantillons : Apprentissage et Test.

La base de données sera séparée comme suit : 2/3 des données constitueront l'*échantillon d'Apprentissage* qui sera utilisé pour construire les modèles ; les 1/3 restant qui constitueront l'*échantillon Test* seront utilisés pour tester nos modèles.

Pour chacune des méthodes d'apprentissage supervisé qui sera utilisée : on commencera par construire le(s) modèle(s) en fonction des différents paramètres propres à chacun d'eux et ce, sur l'échantillon d'apprentissage. Ensuite, dans le but de déterminer le meilleur paramétrage de chaque méthode : nous appliquerons ce(s) modèle(s) ainsi construit(s) sur l'échantillon test, puis nous comparerons les diverses mesures de performance. Ce qui nous permettra de choisir, si il le faut, le meilleur modèle de chaque méthode. Enfin, nous comparerons les résultats des meilleurs modèles de chaque méthode afin de déterminer celui permettant de mieux prédire.

## Quelques précisions

1) Après la transformation, la sélection des variables et l'affectation des effectifs réalisée **au premier semestre**, les variables restantes pouvant être utilisées sont : *Bike\_Age, Bike\_Alc\_D, Bike\_Dir, Bike\_Injur, Bike\_Pos, Bike\_Race, Bike\_Sex, Crash\_Hour, Crash\_Loc, Crash\_Time, Crash\_Type, Crash\_Ty\_1, Developmen, DrvrAge\_Gr, Drvr\_Age, Drvr\_Alc\_D, Drvr\_Injur, Drvr\_Race, Drvr\_Sex, Hit\_Run, Light\_Cond, Num\_Lanes, Num\_Units, Rd\_Charact, Rd\_Class, Rd\_Conditi, Region, Rural\_Urba et Workzone\_I*.

2) En apprentissage automatique supervisé, une matrice de confusion est une matrice qui mesure la qualité d'un système de classification. Elles s'interprètent comme suit :

TABLE 2 – Effectif de la variable à prédire

Classe réelle	Classe estimée par le classificateur	
	Non	Oui
Non	Vrais négatifs	Faux négatifs
Oui	Faux positifs	Vrais positifs

Un vrai positif (VP) est un résultat où le modèle prédit correctement la classe positive. De façon analogue, un vrai négatif (VN) est un résultat où le modèle prédit correctement la classe négative.

Un faux positif (FP) est un résultat où le modèle prédit incorrectement la classe positive et un faux négatif (FN) est un résultat où le modèle prédit incorrectement la classe négative.

3) À partir de la matrice de confusion on peut dériver tout un tas de critères de performance. Parmi lesquels :

- **La sensibilité** : le taux de vrais positifs, c'est à dire la proportion de blessés que l'on a correctement identifiés. C'est la capacité de notre modèle à détecter toutes les blessures.

$$\text{Sensibilité} = \frac{VP}{VP + FN}$$

- **La spécificité** : le taux de vrais négatifs, autrement dit la capacité à détecter toutes les situations où il n'y a pas de blessures. C'est une mesure complémentaire de la sensibilité.

$$\text{Spécificité} = \frac{VN}{FP + VN}$$

- **La précision** : la proportion de prédictions correctes parmi les blessures que l'on a prédites positives. C'est la capacité de notre modèle à prédire qu'un individu a été blessé que si il a été réellement blessé.

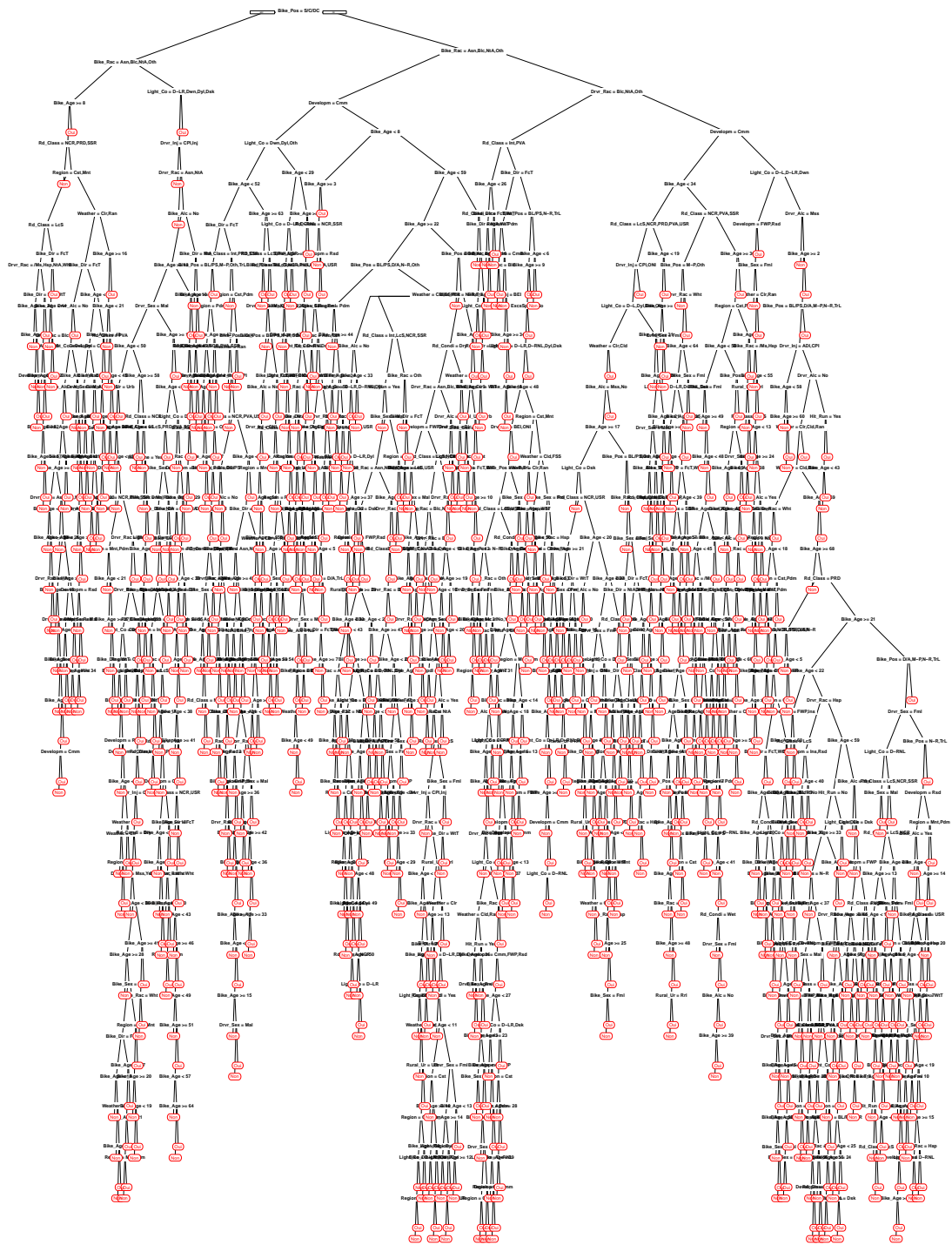
$$\text{Précision} = \frac{VP}{VP + FP}$$

## 2 Arbres

Nous allons commencer par réaliser un modèle d'arbre complet sur *l'échantillon d'apprentissage* et ce, pour avoir une vue d'ensemble.

### 2.1 Arbre complet

## Arbre complet



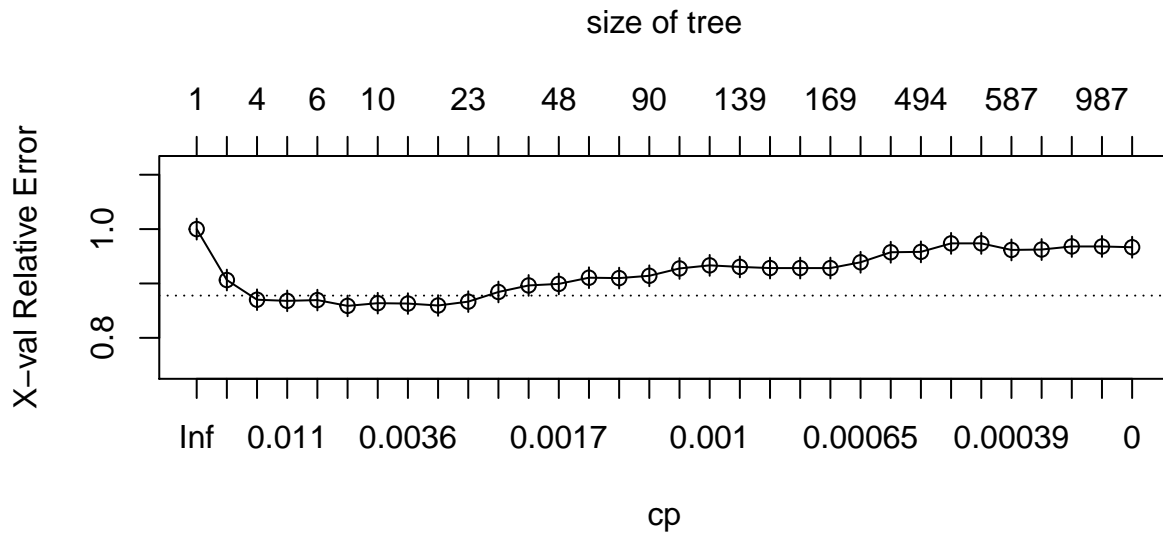
Nous savons qu'un arbre complet n'est généralement pas utilisé pour les prédictions car il ne commet aucune erreur sur l'échantillon sur lequel il est construit, puisqu'il en épouse toutes les caractéristiques. Par conséquent il est difficilement généralisable.

Nous devons donc élaguer notre arbre selon un certain niveau de complexité afin de palier à ce problème de surapprentissage.

## 2.2 Arbre élagué

Cherchons le niveau de complexité qui minimise l'erreur estimée.

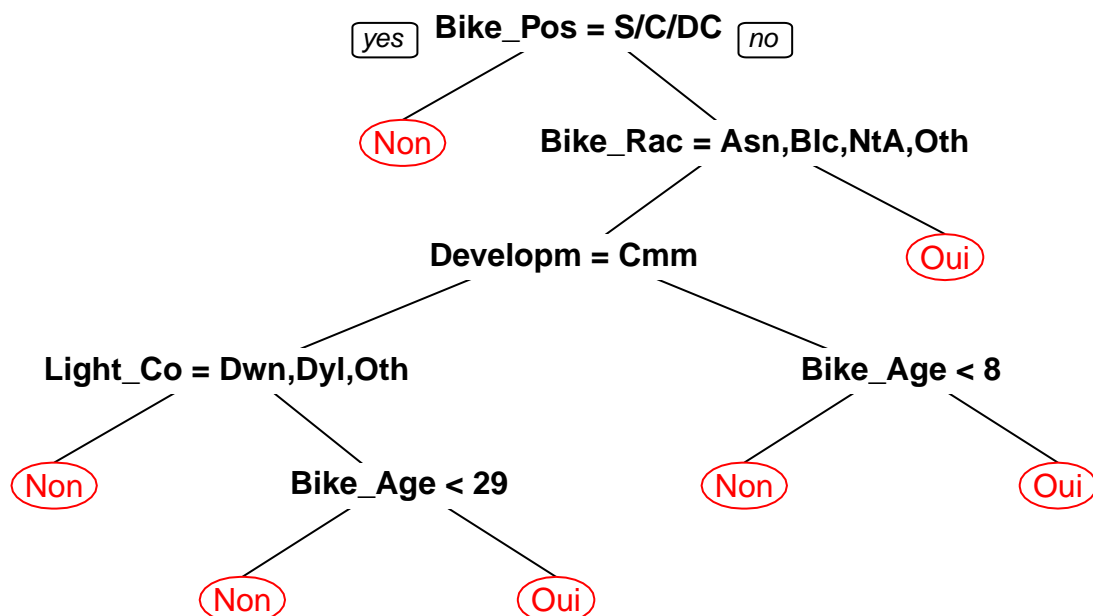
Nous allons réaliser un graphique qui nous montre le taux d'erreur en fonction de la complexité.



Ce graphique nous montre que la complexité qui permet de minimiser l'erreur estimée est de 0.001420455 avec une erreur égale à 0.88 environ.

En appliquant ce niveau de complexité on obtient l'arbre élagué suivant :

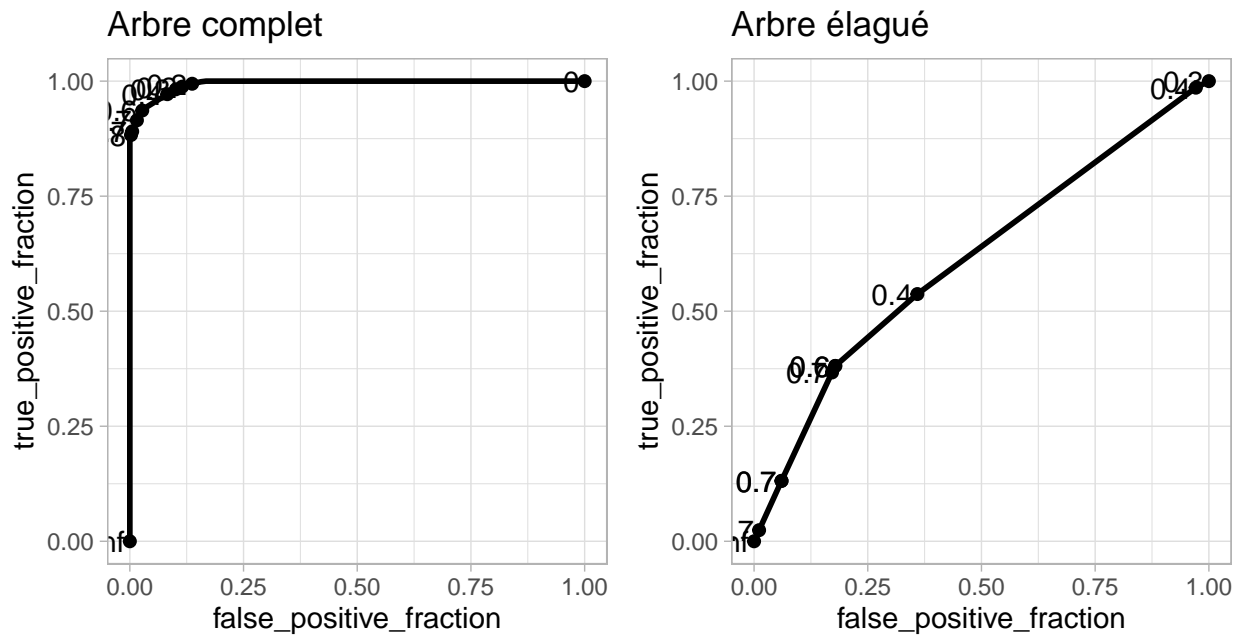
### Arbre élagué



Pour se convaincre de l'utilisation d'un arbre élagué par rapport à l'arbre complet on peut représenter les courbes ROC.

Une courbe ROC (receiver operating characteristic) est une représentation graphique de la relation qu'il existe entre la sensibilité (qui mesure sa capacité à donner un résultat positif lorsqu'une hypothèse est vérifiée) et la

spécificité (qui mesure la capacité d'un test à donner un résultat négatif lorsque l'hypothèse est vérifiée) d'un test pour chaque valeur seuil considérée.



Comme prévu, la courbe ROC de l'arbre complet est parfaite à cause du surapprentissage. La courbe ROC de l'arbre élagué est près de la bissectrice mais suffisamment au dessus pour conclure que l'utilisation de l'arbre élagué est préférable à celui d'un classificateur aléatoire .

Nous poursuivrons la comparaison en appliquant ces deux modèles à l'échantillon d'apprentissage

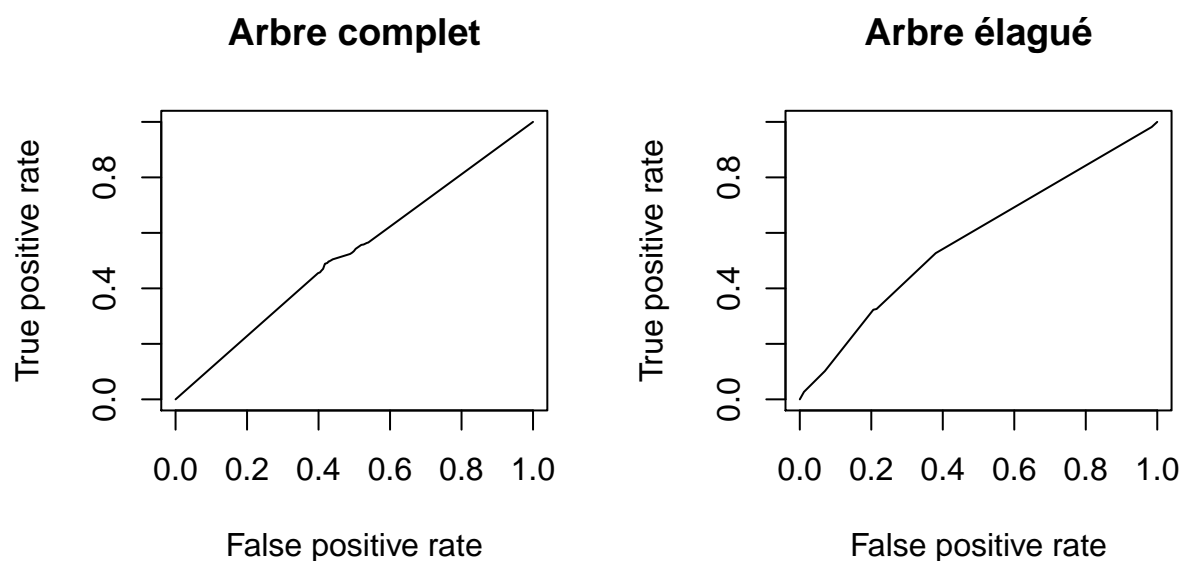
## Application à l'échantillon test

En appliquant ces modèles à *l'échantillon Test* on obtient la matrice de confusion suivante :

TABLE 3 – Matrice de confusion Arbres : test

Prédiction	Référence			
	Arbre complet		Arbre élagué	
	Non	Oui	Non	Oui
Non	351	330	230	165
Oui	355	438	476	603

## Courbe ROC



Les mesures de performance de nos arbres appliqués à l'échantillon test sont renseignés dans le tableau suivant :

TABLE 4 – mesures de performance

	Sensitivity	Specificity	Precision	AUC
Arbre complet	0.4971671	0.5703125	0.5154185	0.5240138
Arbre élagué	0.3257790	0.7851562	0.5822785	0.5773375

**Conclusion:** Il semblerait, dans le cas particulier de notre base de données, que l'arbre complet et l'arbre élagué aient des performances qui se valent.

## 3 Bagging et Random forest

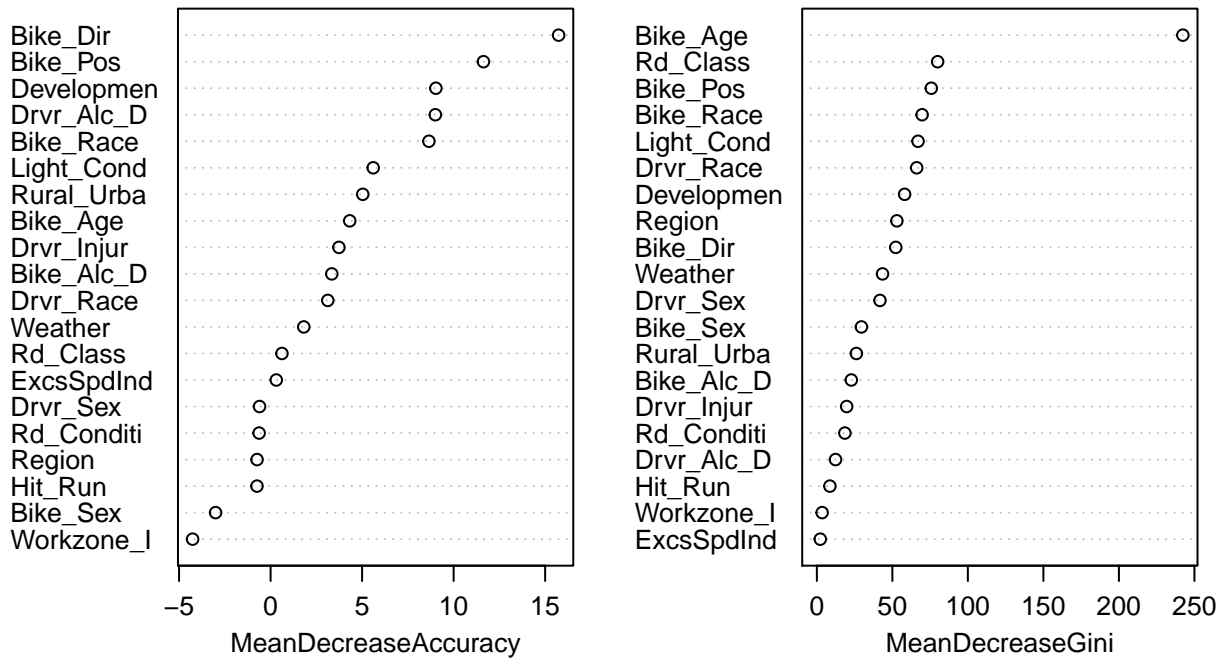
Dans cette partie nous allons comparer la méthode *RandomForest* classique avec celle du *Bagging*. Nous allons d'abord construire ces modèles sur les données d'apprentissage, comparer les résultats et enfin nous les appliquerons sur nos données test avant de conclure.

### 3.1 Random forest

L'algorithme Random Forest, est l'un des plus couramment utilisé, il s'agit d'un type spécial de bagging appliqué aux arbres de décision. Cet algorithme de forêt aléatoire de Breiman(basé sur le code Fortran original de Breiman et Cutler) combine les concepts de sous-espaces aléatoires et de bagging. L'algorithme des forêts d'arbres décisionnels effectue un apprentissage sur de multiples arbres de décision entraînés sur des sous-ensembles de données légèrement différents.

Commençons par déterminer l'importance des variables selon la méthode de *Mean Decrease Accuracy*.

## Random Forest



Ce graphique montre l'importance qu'a chaque variable dans la construction de notre forêt. On peut noter que les variables : *Bike\_Race*, *Bike\_Dir*, *Bike\_Pos* et *Rd\_Class* sont les plus importantes.

En construisant la forêt aléatoire sur nos données d'apprentissage on obtient :

TABLE 5 – Matrice de confusion Random Forest : apprentissage

Prédiction	Référence		class.error
	Non	Oui	
Non	694	715	0.5074521
Oui	542	998	0.3519481

Cette méthode permet d'obtenir une erreur OOB ( Out Of Bag ) de 42.96%.



## ## Bagging

*Bagging* signifie bootstrap aggregating. C'est un méta-algorithme d'ensemble d'apprentissage automatique conçu pour améliorer la stabilité et la précision des algorithmes d'apprentissage automatique utilisés dans la classification statistique et la régression.

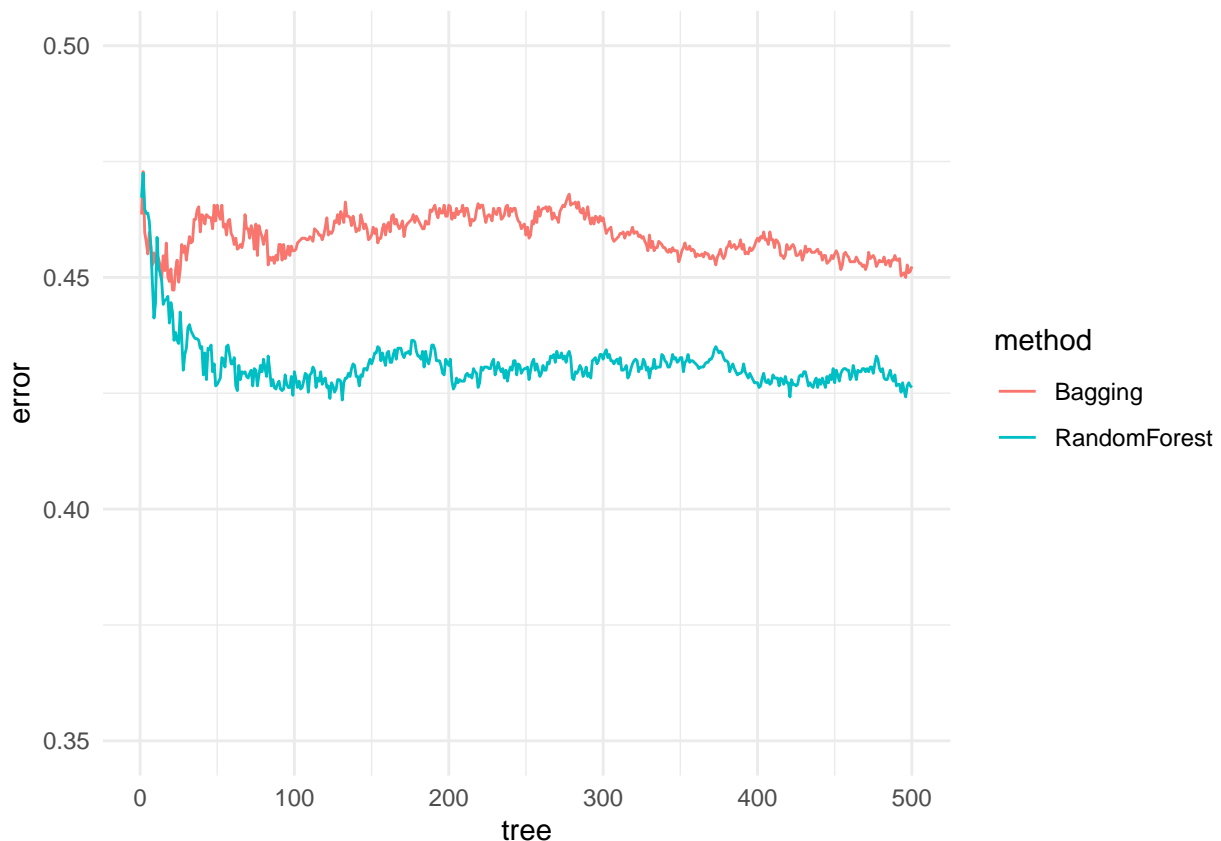
TABLE 6 – Matrice de confusion bagging : apprentissage

Prédiction	Référence		class.error
	Non	Oui	
Non	689	720	0.5110007
Oui	614	926	0.3987013

En utilisant le Bagging : on obtient une erreur OOB ( Out Of Bag ) égale à 44.6% sur nos données d'entraînement.

### 3.2 Comparaison bagging et Random forest

Comparons les niveaux d'erreurs des deux modèles. Pour 500 arbres, on a :



En comparant les erreurs sur les données d'entraînement des 2 modèles pour 500 arbres, on remarque que l'erreur provenant du *Random Forest* est plus faible que l'erreur du bagging quelque soit le nombre d'arbre considéré.

### 3.3 Application des modèles à l'échantillon Test

Les mesures de performances de ce modèle sont les suivantes :

**Conclusion :** Le Random forest permet un taux de précision supérieur à celui du Random Forest. De plus

TABLE 7 – Matrice de confusion test

Prédiction	Référence			
	Random Forest		Bagging	
	Non	Oui	Non	Oui
Non	354	279	344	283
Oui	352	489	362	485

TABLE 8 – Mesures de performance

	Sensitivity	Specificity	Precision	AUC
Random Forest	0.5014164	0.6367188	0.5592417	0.6005546
Bagging	0.4872521	0.6315104	0.5486443	0.5879653

## 4 Modèle randomForest automatisé

En laissant le logiciel déterminer les paramètres optimaux pour le modèle Random Forest, on obtient les paramètres suivants : 100 noeuds, une variable et 250 arbres.

En appliquant ces paramètres et en construisant le modèle sur notre échantillon d'apprentissage, on obtient la matrice de confusion :

TABLE 9 – Matrice de confusion Random Forest automatisé : apprentissage

Prédiction	Référence		class.error
	Non	Oui	
Non	587	822	0.5833925
Oui	359	1181	0.2331169

### 4.1 Application à l'échantillon test

En appliquant ce modèle à l'échantillon test on obtient :

TABLE 10 – Matrice de confusion Random Forest automatisé : test

Prédiction	Référence	
	Non	Oui
Non	276	182
Oui	430	586

Les mesures de performances de ce modèle sont les suivantes :

TABLE 11 – Mesures de performance

	Sensitivity	Specificity	Precision	AUC
Random Forest automatisé	0.3909348	0.7630208	0.6026201	0.6027346

**Conclusion :** Les performances de ce modèle semblent bien meilleures que tous les précédents.

## 5 Boosting

Le *Boosting* génère une séquence de modèles de classification, chaque modèle de classification successif dans la séquence permettant de mieux prévoir la classification des observations qui était mal classée par les modèles de classification précédents. Lors du déploiement, les prévisions issues des différents modèles de classification pourront alors être combinées afin d'obtenir la meilleure prévision ou classification. Le Boosting, est similaire à la méthode bagging. En revanche, les étapes se produisent de manière séquentielles et non pas simultanées.

En construisant le boosting, on obtient la matrice de confusion suivante :

TABLE 12 – Matrice de confusion apprentissage : Boosting

Référence	Prédiction	
	Non	Oui
Non	1392	17
Oui	23	1517

L'erreur Out of bag est de 0.019.

On constate, comme on pouvait s'y attendre, un problème de surapprentissage. On décide donc d'appliquer différentes pénalisations à notre modèle et ce afin de trouver le modèle optimal.

Les pénalisations généralement utilisées sur le Boosting sont : 0.1, 0.01 et 0.001.

**Pénalisation = 0.1**

TABLE 13 – Matrice de confusion apprentissage : Boosting p=0.1

Référence	Prédiction	
	Non	Oui
Non	689	720
Oui	446	1094

L'erreur Out of bag est de 0.398.

**Pénalisation = 0.01**

TABLE 14 – Matrice de confusion apprentissage : Boosting p=0.01

Référence	Prédiction	
	Non	Oui
Non	462	947
Oui	252	1288

L'erreur Out of bag est de 0.408.

**Pénalisation = 0.001**

L'erreur Out of bag est de 0.409.

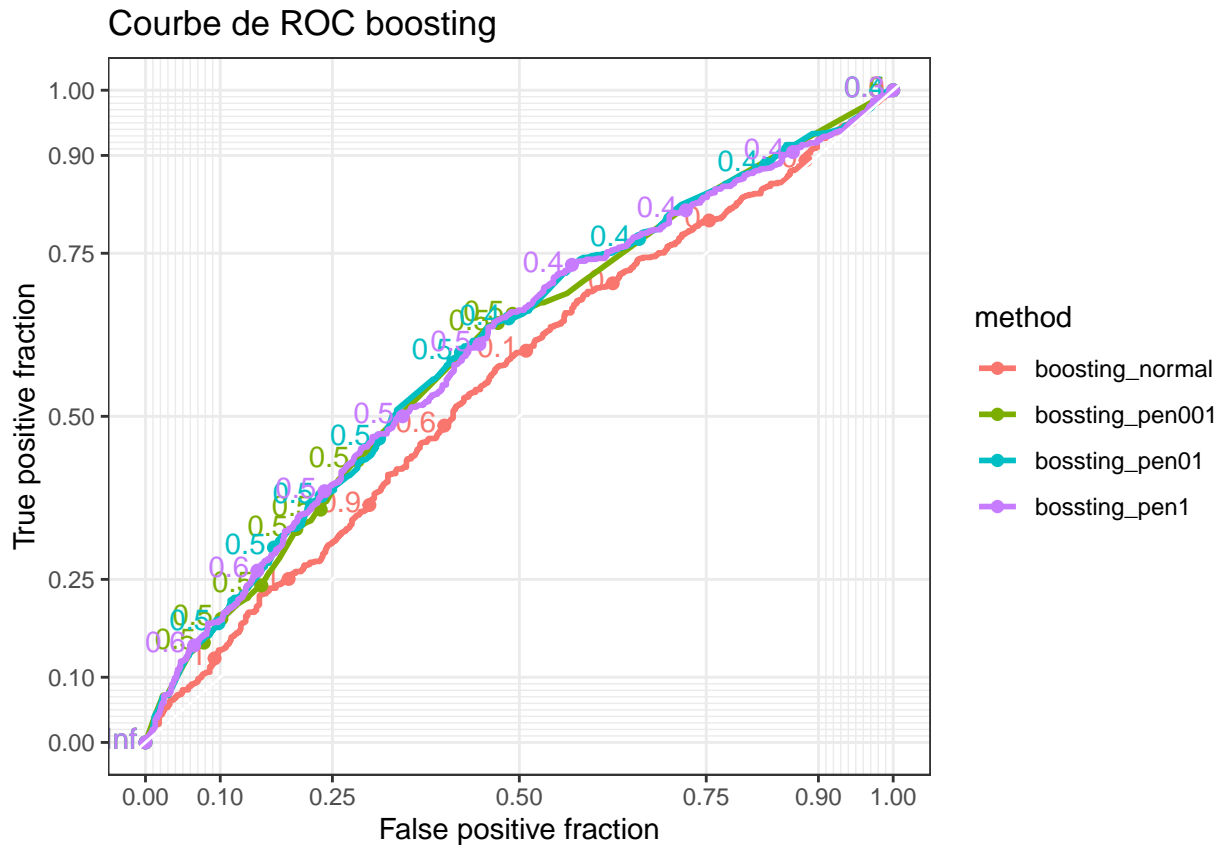
Afin de déterminer le modèle boosting optimal, on va appliquer chacun d'entre eux à l'échantillon test. Le meilleur modèle sera celui permettant la meilleure prévision.

TABLE 15 – Matrice de confusion apprentissage : Boosting p=0.001

Référence	Prédiction	
	Non	Oui
Non	336	1073
Oui	174	1366

### 5.1 Application à l'échantillon test

Pour déterminer le meilleur modèle Boosting on va comparer les courbes ROC ainsi que les mesures de performance.



Le modèle ayant l'AUC la plus grande est le modèle boosting pénalisé à 0.01.

On pourrait avant de conclure, comparer également les mesures de performance des tests :

TABLE 16 – Comparaisons des boosting

	Sensibilité	Spécificité	Précision	AUC
Boosting	0.5127479	0.5859375	0.5323529	0.5558144
Boosting p=0.1	0.4603399	0.6979167	0.5834829	0.6054494
Boosting p=0.01	0.2932011	0.8281250	0.6106195	0.6077031
Boosting p=0.001	0.2209632	0.8632812	0.5977011	0.6031265

**Conclusion :** Les AUC étant très similaires, pour le choix du meilleur modèle Boosting on se base sur la Précision : c'est donc le Boosting pénalisé à 0.01 qui est choisi.

## 6 Scoring

Comme toute bonne démarche de modélisation, la construction d'un bon **score** se fait par une succession d'étapes : nous commencerons par vérifier la liaison entre les descripteurs, ensuite nous construirons les modèles sur l'*échantillon d'apprentissage* et enfin nous les appliquerons sur l'*échantillon test*. Nous pourrons ensuite comparer les mesures de performance afin de déterminer le meilleur modèle.

Nous comparerons plusieurs modèles et retiendrons le modèle le plus adéquat selon l'objectif de l'étude.

### 6.1 Variable à expliquer

Revenons sur nos données initiales. La variable **Bike\_Injur** est séparée en plusieurs modalités, comme suit :

A: Disabling Injury	B: Evident Injury	C: Possible Injury	Injury	K: Killed	O: No Injury
291	2405	2199	172	123	526

Nous allons affecter à ces modalités les valeurs 0 et 1.

La valeur 1 : à celles qui concernent les blessures avérées et graves (Killed, Disabling Injury, Evident Injury et Injury) et la valeur 0 à celles qui concernent l'absence, évidente ou non, de blessure (No Injury et Possible Injury).

On obtient donc :

	Bike_Injur	
	0	1
Effectif	2725	2991

On souhaite déterminer la probabilité qu'un cycliste soit blessé ou non compte tenu des paramètres de son accident.

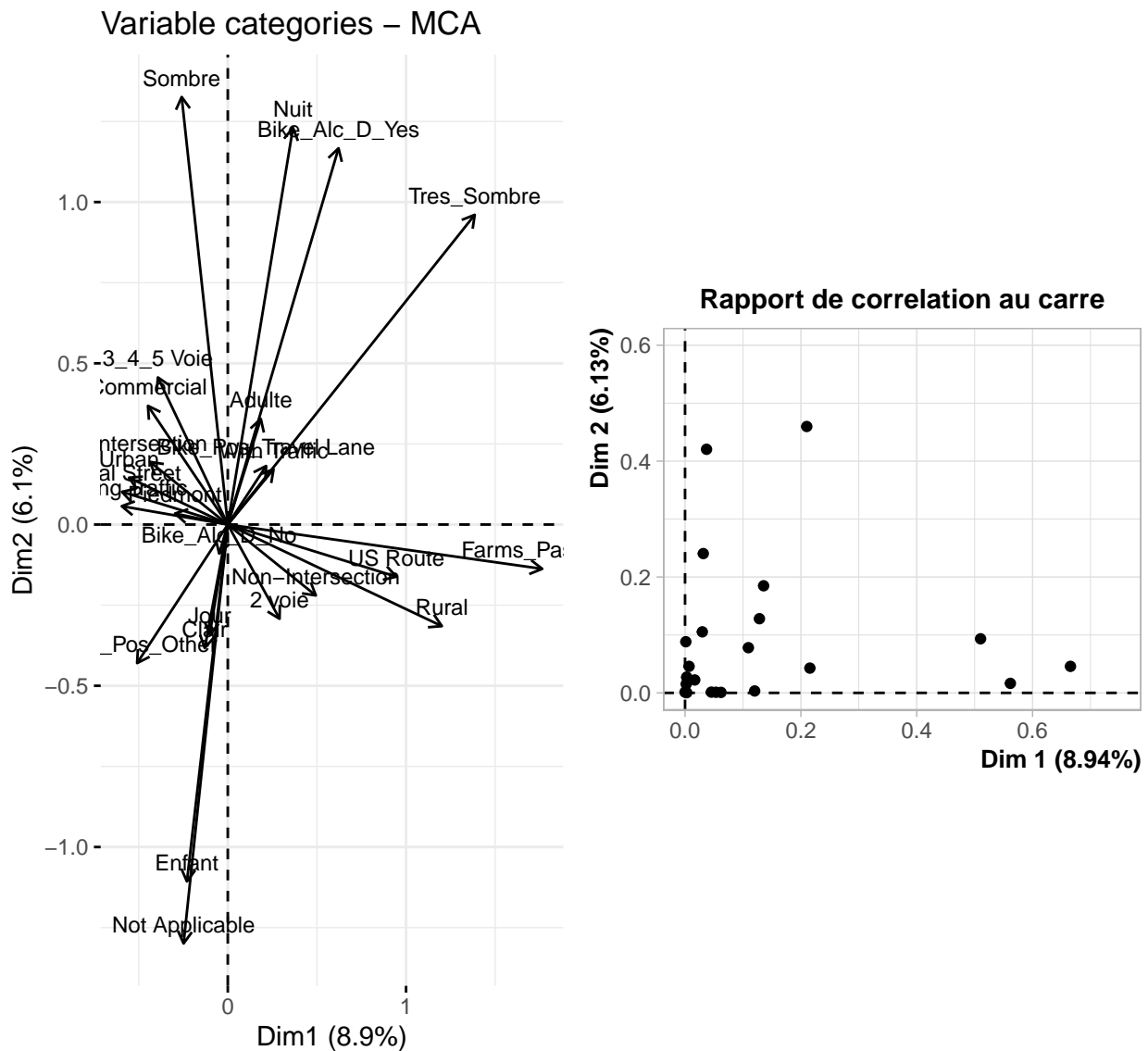
Pour cela, nous utiliserons des **Régression logistique**.

Nous allons effectuer une sélection de variables et ce pour plusieurs raisons.

D'abord parce que certaines des variables de notre base de données sont inutilisables en l'état et d'autre part parce qu'un modèle avec peu de variables sera plus facilement généralisable en terme de robustesse : *Principe du rasoir d'Occam*.

### 6.2 Analyse exploratoire

Construisons une ACM afin d'essayer d'identifier les corrélations éventuelles entre les variables explicatives.



Plusieurs variables semblent très corrélées, ce qui pourrait entraîner des problèmes de colinéarité lors de la régression.

### 6.3 Construction des modèles

Le **modèle général** que nous allons construire est un modèle naïf. Il prend en compte toutes les variables, sans aucune spécification particulière :

$$\begin{aligned}
 \text{Bike\_Injur}_i = & \beta_0 + \beta_1 \text{Bike\_Age}_i + \beta_2 \text{Bike\_Alc\_D}_i + \beta_3 \text{Bike\_Dir}_i + \beta_4 \text{Bike\_Pos}_i + \beta_5 \text{Bike\_Race}_i \\
 & + \beta_6 \text{Bike\_Sex}_i + \beta_7 \text{Crash\_Hour}_i + \beta_8 \text{Crash\_Loc}_i + \beta_9 \text{Developmen}_i + \beta_{10} \text{Drvr\_Alc\_D}_i + \beta_{11} \\
 & \text{Drvr\_Injur}_i + \beta_{12} \text{Drvr\_Race}_i + \beta_{13} \text{Drvr\_Sex}_i + \beta_{14} \text{Hit\_Run}_i + \beta_{15} \text{Light\_Cond}_i + \beta_{16} \\
 & \text{Num\_Lanes}_i + \beta_{17} \text{Rd\_Class}_i + \beta_{18} \text{Rd\_Condit}_i + \beta_{19} \text{Region}_i + \beta_{20} \text{Rural\_Urba}_i + \beta_{21} \text{Workzone\_I}_i \\
 & + \varepsilon_i
 \end{aligned}$$

Nous savons que ce modèle sera très mauvais étant donné le nombre de variable utilisée. Pour obtenir un modèle pertinent nous devons effectuer une sélection des variables et ce pour plusieurs raisons, la principale étant que : un modèle avec peu de variables sera plus facilement généralisable en terme de robustesse *Principe du rasoir d'Occam*.

Le second modèle, le **modèle AIC** : sera obtenu en faisant une sélection automatique de variables, sur le critère d'Akaike (AIC). Ce dernier s'écrit comme suit :  $AIC = 2k - 2 \ln(L)$  ; où k est le nombre de paramètres à estimer du modèle et L est le maximum de la fonction de vraisemblance du modèle. Si l'on considère un ensemble de modèles candidats, le modèle choisi sera celui qui aura la plus faible valeur d'AIC.

Le dernier modèle, le **modèle AIC modifié** sera issue du second. On supprimera toutes les variables non significatives du second modèle.

Les différents résultats des régressions obtenus sont renseignés dans le tableau ci-dessous :

TABLE 17 – Résultats

	<i>Dependent variable:</i>		
	Bike_Injur		
	(1)	(2)	(3)
Bike_AgeEnfant			0.066 (0.128)
Bike_AgeJeune			-0.073 (0.088)
Bike_Alc_DYes			0.431*** (0.163)
Bike_DirNot Applicable	0.676*** (0.124)	0.709*** (0.123)	0.550*** (0.161)
Bike_DirWith Traffic	0.456*** (0.082)	0.486*** (0.081)	0.414*** (0.102)
Bike_PosTravel Lane	0.301*** (0.080)	0.316*** (0.080)	0.320*** (0.100)
Bike_RaceOther	0.470*** (0.119)	0.473*** (0.119)	0.421*** (0.148)
Bike_RaceWhite	0.495*** (0.071)	0.480*** (0.070)	0.477*** (0.095)
Bike_SexMale			0.036 (0.111)
Crash_HourNuit			0.138 (0.096)
Crash_LocNon-Intersection			0.123 (0.086)
DevelopmenFarms_Pastures	0.195 (0.121)		0.007 (0.149)
DevelopmenInstitutional	0.263 (0.195)		0.004 (0.245)
DevelopmenResidential	0.206*** (0.071)		0.122 (0.091)
Drvr_Alc_DNo	0.408 (0.675)	0.393 (0.670)	0.123 (0.719)
Drvr_Alc_DYes	1.300* (0.722)	1.288* (0.717)	1.070 (0.787)
Drvr_Injur1			-0.063 (0.259)
Drvr_RaceBlack			-0.318 (0.607)
Drvr_RaceOther			-0.265 (0.621)
Drvr_RaceWhite			-0.294 (0.605)
Drvr_SexMale			0.095 (0.080)
Rd_ConditiWet			-0.004 (0.155)
RegionMountains			0.267* (0.161)
RegionPiedmont			0.160* (0.090)
Rural_UrbaUrban	-0.185** (0.083)	-0.242*** (0.070)	-0.288*** (0.108)
Workzone_IYes	0.062 (0.473)		0.265 (0.508)
Constant	-1.218* (0.684)	-1.066 (0.677)	-0.742 (0.960)
Observations	4,191	4,191	2,794
Log Likelihood	-2,796.461	-2,801.077	-1,851.478
Akaike Inf. Crit.	5,618.923	5,620.154	3,756.957
Note: * p<0.1; ** p<0.05; *** p<0.01			

Modèle AIC:

$$Bike\_Injur_i = \beta_0 + \beta_1 Bike\_Alc\_D_i + \beta_2 Bike\_Dir_i + \beta_3 Bike\_Pos_i + \beta_4 Bike\_Race_i + \beta_9 Developmen_i + \beta_5 Drvr\_Alc\_D_i + \beta_6 Rural\_Urba_i + \beta_7 Workzone\_I_i + \varepsilon_i$$

Modèle AIC modifié:

$$Bike\_Injur_i = \beta_0 + \beta_1 Bike\_Alc\_D_i + \beta_2 Bike\_Dir_i + \beta_3 Bike\_Pos_i + \beta_4 Bike\_Race_i + \beta_5 Drvr\_Alc\_D_i + \beta_6 Rural\_Urba_i + \varepsilon_i$$

## 6.4 Validation des modèles : Indicateurs de qualité et de robustesse

On exclu le modèle général car la plupart des coefficients ne sont pas significatifs.



On s'intéressera donc exclusivement aux modèles AIC (Régression 1) et AIC modifié (Régression 2).

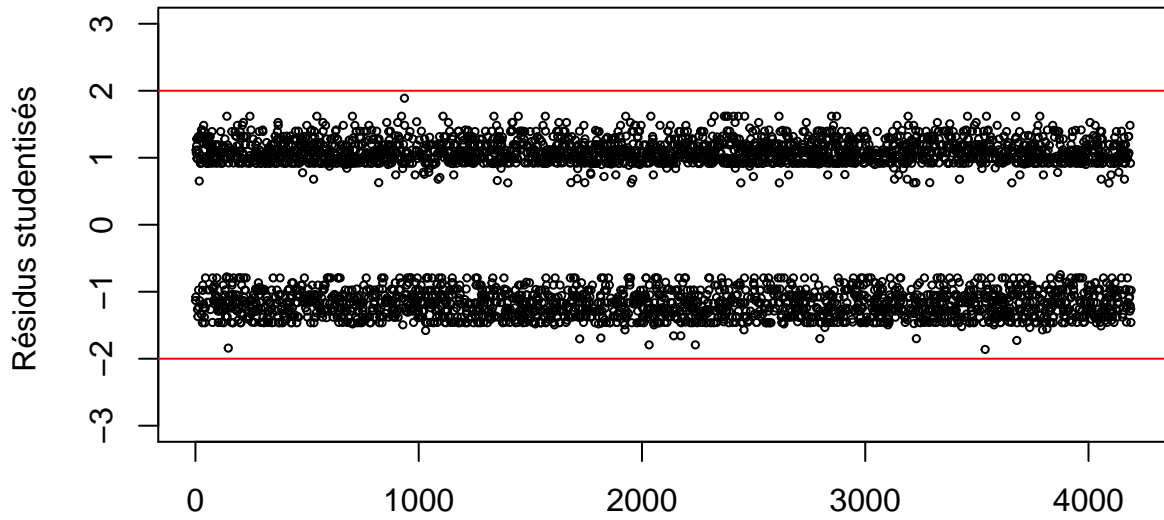
TABLE 18 – Résultats des modèles 2 et 3

	<i>Dependent variable:</i>	
	Bike_Injur	
	(1)	(2)
Bike_DirNot Applicable	0.676*** (0.124)	0.709*** (0.123)
Bike_DirWith Traffic	0.456*** (0.082)	0.486*** (0.081)
Bike_PosTravel Lane	0.301*** (0.080)	0.316*** (0.080)
Bike_RaceOther	0.470*** (0.119)	0.473*** (0.119)
Bike_RaceWhite	0.495*** (0.071)	0.480*** (0.070)
DevelopmenFarms_Pastures	0.195 (0.121)	
DevelopmenInstitutional	0.263 (0.195)	
DevelopmenResidential	0.206*** (0.071)	
Drvr_Alc_DNo	0.408 (0.675)	0.393 (0.670)
Drvr_Alc_DYes	1.300* (0.722)	1.288* (0.717)
Rural_UrbaUrban	-0.185** (0.083)	-0.242*** (0.070)
Workzone_IYes	0.062 (0.473)	
Constant	-1.218* (0.684)	-1.066 (0.677)
Observations	4,191	4,191
Log Likelihood	-2,796.461	-2,801.077
Akaike Inf. Crit.	5,618.923	5,620.154
<i>Note:</i> * p<0.1; ** p<0.05; *** p<0.01		

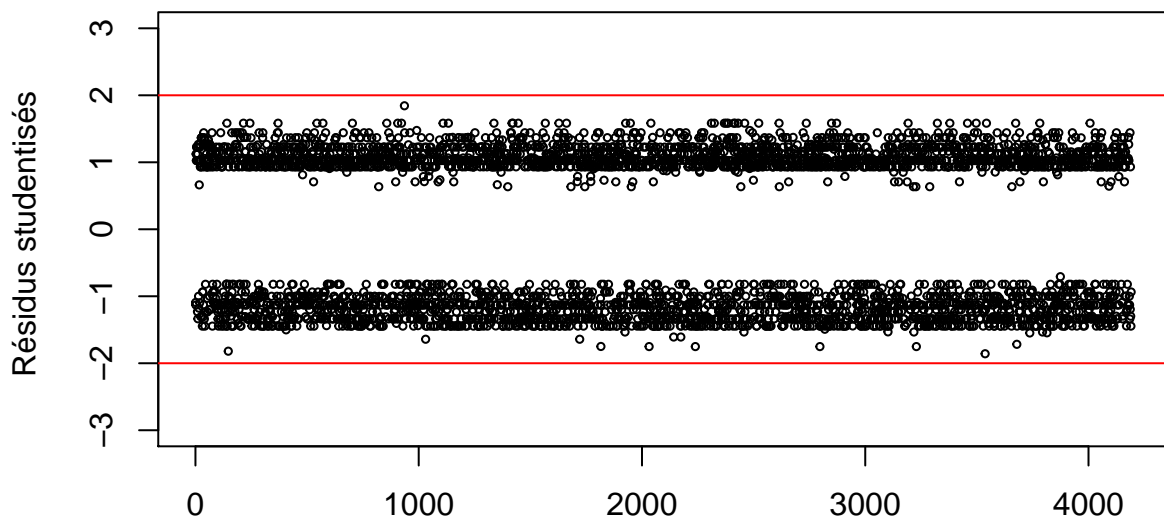
## Résidus de déviations

Pour les régressions logistiques, on s'intéresse la plupart du temps aux résidus de déviance. Ils prennent généralement des valeurs qui oscillent entre -2 et 2.

### Modèle 1



### Modèle 2



Il semblerait qu'il n'y ait pas de valeurs aberrantes.

Les deux modèles sont donc utilisables dans l'état.

## 6.5 Application à l'échantillon test

En appliquant les deux modèles à l'échantillon *Test* on a :

TABLE 19 – Matrice de confusion Scoring : test

Prédiction	Référence			
	Modèle AIC		Modèle AIC modifié	
	0	1	0	1
0	352	242	334	223
1	322	481	340	500

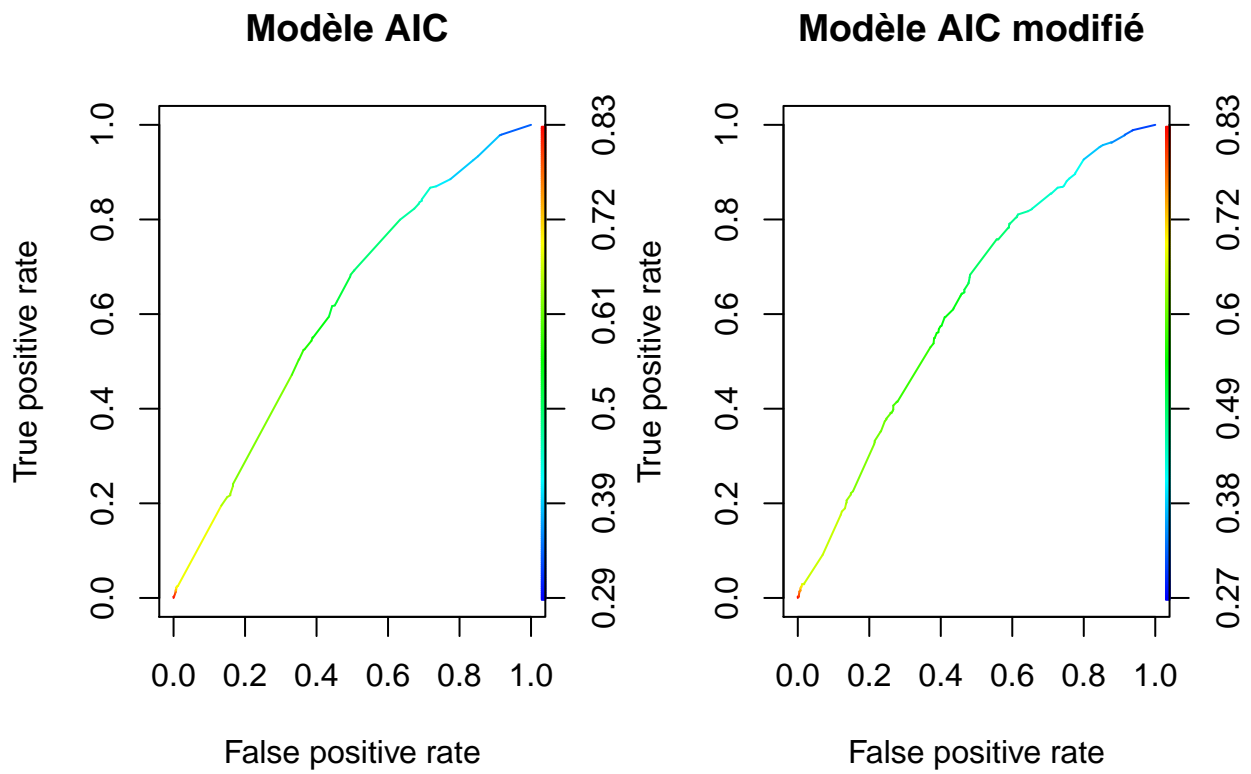
Qu'en est-il du taux d'erreur ?

TABLE 20 – Taux d'erreur test

Modèle AIC	Modèle AIC modifié
0.4037223	0.4030064

Le taux d'erreur est relativement le même pour les deux modèles. Il est inférieur à 0.5, ce qui est suffisant pour conclure que les modèles sont pertinents.

## Courbes ROC



## Aire sous les courbes

TABLE 21 – Aire sous les courbes

Modèle AIC	Modèle AIC modifié
0.6203463	0.6107239

L'aire sous les courbes est relativement le même lui aussi.

## Mesure de performance :

TABLE 22 – Mesure de performance : Scoring

	Sensibilité	Spécificité	Précision
Modèle AIC	0.5222552	0.6652835	0.5925926
Modèle AIC modifié	0.4955490	0.6915629	0.5996409

**Conclusion :** Les deux modèles sont très similaires et ont des mesures de performances qui se valent mais on décide de privilégier celui étant le plus précis : le modèle AUC modifié.

## 7 Choix du meilleur modèle

TABLE 23 – Comparaison des modèles

	Sensitivity	Specificity	Precision	AUC
Arbre complet	0.4971671	0.5703125	0.5154185	0.5240138
Arbre élagué	0.3257790	0.7851562	0.5822785	0.5773375
Random Forest	0.5014164	0.6367188	0.5592417	0.6005546
Bagging	0.4872521	0.6315104	0.5486443	0.5879653
Random Forest automatisé	0.3909348	0.7630208	0.6026201	0.6027346
Boosting p=0.01	0.2932011	0.8281250	0.6106195	0.6077031
Scoring	0.4955490	0.6915629	0.5996409	0.6107239

## 8 Discussion

La sélection des mesures de performance les plus pertinentes se fait en fonction de la problématique à traiter. La notre pourrait être soit de déterminer la probabilité que l'accident ait causé une blessure soit de déterminer la probabilité que l'accident n'ait pas causé de blessure. Il n'y a pas grand intérêt à déterminer la probabilité que l'accident ait causé une blessure pour les hôpitaux par exemple car ces derniers envoient systématiquement une ambulance. En revanche il serait très intéressant pour eux de déterminer la probabilité que l'accident n'ait pas causé de blessures et ce afin de gérer le flux des ambulances ou simplement de faire un choix prioritaire parmi deux situations par exemple. En ce sens la sensibilité ne donne pas une mesure très pertinente. La spécificité en revanche, qui mesure le taux de vrais négatifs (dans notre cas qu'il n'y ait pas de blessure) semble plus pertinente.

Aussi, La précision est de fait une mesure très intéressante. Dans le cadre de notre étude : c'est la capacité de nos modèles à ne prédire non à une blessure si l'accident n'a effectivement pas entraîné une blessure. Pour évaluer un compromis entre sensibilité et précision, on peut calculer la "F-mesure", qui est leur moyenne harmonique.

Le calcul du F-mesure ou encore F-score est le suivant:

$$F_{mesure} = \frac{2 \times Précision \times sensibilité}{précision + sensibilité}$$

Ces considérations ainsi prises en compte nous permettent de conclure que le modèle RandomForest automatisé (qui pour rappel est un modèle Random Forest ayant les paramètres suivant : 250 arbres, 50 nœud et 1 ) et le Boosting pénalisé à 0.01 sont les modèles permettant d'obtenir les meilleurs résultats pour la problématique : Quelle est la probabilité que l'accident n'est pas engendré de blessures.

TABLE 24 – Comparaison des F-mesure des modèles

	F-mesure
Arbre complet	0.5061283
Arbre élagué	0.4178020
Random Forest	0.5287528
Bagging	0.5161290
Random Forest automatisé	0.4742268
Boosting p=0.01	0.3961722
Scoring	0.5426483