

# Détection de fraude dans le domaine de l'assurance

Arbres de décision et prédiction

Léa Darne  
Hugo Munive

Janvier 2024



# Introduction

## Contexte et objectifs

La construction des arbres de prédiction représente un élément essentiel dans le but de prédire au mieux les événements possible de se produire. Le présent rapport s'inscrit dans le cadre du projet de fin d'année du cours d'Analyse de Données, visant à développer un modèle de prédiction du risque de déclarations frauduleuses d'accidents par les assurés. L'objectif principal est d'explorer les ensembles de données, composés de données financières, démographiques et relatives aux accidents, afin de concevoir des classifieurs performants.

## Justification des outils utilisés

La justification des outils choisis repose sur la nécessité d'exploiter des méthodes robustes de classification et de clustering. Ces outils sont cruciaux pour atteindre l'objectif de minimiser les risques financiers en identifiant de manière précise les déclarations frauduleuses. L'utilisation de méthodes de rééquilibrage des classes, de tests de corrélation et de clustering permettra d'optimiser la prédiction du risque.

## Aperçu du projet

L'ensemble du projet se décompose en plusieurs étapes clés. Dans le chapitre 1, nous explorerons les caractéristiques fondamentales des ensembles de données. Le chapitre 2 se concentrera sur la correction des déséquilibres dans les données. Le chapitre 3, dédié aux *Tests de corrélation*, examinera les relations entre les différentes variables. Par la suite, le chapitre 4, *Clustering*, se penchera sur le regroupement des données similaires. Enfin, dans le chapitre 5, *Méthode de rééquilibrage et clustering*, nous combinerons ces approches pour espérer de maximiser la performance du modèle de prédiction.

Chaque sous-partie suivante détaillera les choix, les méthodes et les résultats spécifiques à sa phase correspondante du projet.

Toutes les fonctions auxiliaires, ainsi que les complémentaires du code R seront mis en annexe.

# Sommaire

<b>1</b>	<b>Première approche et visualisation des données</b>	<b>4</b>
1.1	Chargement et exploration des données	4
1.1.1	Chargement et préparation des données	4
1.1.2	Première approche de visualisation des caractéristiques des données	4
1.2	Arbres de décision	6
1.2.1	Pré-traitement des données	6
1.2.2	Création des arbres	6
1.2.3	Visualisation des arbres	6
1.3	Courbes ROC	8
1.3.1	Fonction auxiliaire	8
1.3.2	Graphique des courbes ROC	8
1.4	Aire sous la courbe (AUC)	9
1.4.1	Fonction auxiliaire	9
1.4.2	Interprétation de l'AUC	9
1.5	Mesure d'évaluation	9
1.5.1	Fonction auxiliaire et résultats	9
1.5.2	Interprétation	11
<b>2</b>	<b>Méthode de rééquilibrage des classes</b>	<b>12</b>
2.1	Chargement de données	12
2.2	Principe du rééquilibrage	12
2.3	Application	13
2.4	Arbres prédictifs et étude du modèle	13
2.4.1	Création des arbres et visualisation	13
2.4.2	Courbes ROC	14
2.4.3	AUC	14
2.4.4	Mesure d'évaluation	15
<b>3</b>	<b>Tests de corrélation</b>	<b>17</b>
3.1	Test de significativité	17
3.1.1	Test de Chi2 pour variables catégorielles	17
3.1.2	Test de Pearson, Spearman et Kendall pour variables numériques	17
3.1.3	Test en mosaïque	18
3.2	Fonction AttrEval	19
3.3	Application	19
3.3.1	Chargement et retraitement de données	19
3.3.2	Création des arbres et visualisation	20
3.3.3	Courbes ROC	20
3.3.4	AUC	21
3.3.5	Matrice de confusion	21
<b>4</b>	<b>Clustering</b>	<b>24</b>
4.1	Chargement et prétraitement de données	24
4.2	Matrice de distance et clustering par partitionnement	24
4.3	Application	25
4.3.1	Arbres	25
4.3.2	ROC	25
4.3.3	AUC	25
4.3.4	Matrice de confusion	25

**5 Méthode de variables significatives et clustering** **27**

5.1 Test expérimental . . . . . 27

5.2 Interprétation . . . . . 30

**6 Conclusion** **31**

# Chapitre 1

## Première approche et visualisation des données

Toutes les fonctions auxiliaires utilisées dans ce projet seront documentées dans le code en annexe.

### 1.1 Chargement et exploration des données

Dans cette première partie, nous abordons le chargement des données à partir du fichier CSV fourni, "Data\_Projet1.csv". Nous utilisons des bibliothèques telles que ggplot2 pour la visualisation initiale et inspectons les premières lignes du jeu de données avec la fonction View(). De plus, nous explorons le résumé statistique et les noms des variables pour obtenir une compréhension initiale de la structure des données.

#### 1.1.1 Chargement et préparation des données

Dans toute la suite, les données utilisées dans ce projet seront stockées dans une variable appelée `fraud_data`. Tout d'abord, à l'aide de la commande `str(fraud_data)`, on peut observer qu'il s'agit d'un data frame contenant 1100 observations et 12 variables. Parmi ces 12 variables, qui seront appelées les variables explicatives, il y a la variable `fraudulent`, qui est notre variable cible à prédire.

```
1 fraud_data <- read.csv("Donnees/Data_Projet_1.csv",  
2   header = TRUE, sep = ",", dec = ".", stringsAsFactors = T)
```

En revanche, on est contraint d'enlever les variables du type `id`, vu que leur but est seulement d'organiser les variables pour une meilleure lecture des données, mais ne participent pas à la prédiction de la variable cible. Pour ce faire, on utilise la fonction `subset`

```
1 fraud_data <- subset(fraud_data, select=-c(customer_id, claim_id))
```

#### 1.1.2 Première approche de visualisation des caractéristiques des données

Afin de cerner notre problème de prédiction de la classe cible *fraudulent*, une première idée c'est de regarder, dans un premier temps, la répartition de la classe cible, i.e. le nombre des cas réels *oui* et *non* dans la classe parmi les 1100 observations. Dans un second temps, nous utilisons des graphiques à barres pour examiner la distribution des déclarations frauduleuses en fonction de diverses variables explicatives. On obtient donc les graphiques suivants :



## 1.2 Arbres de décision

### 1.2.1 Pré-traitement des données

Afin de pouvoir créer des arbres de décision, on utilisera la démarche suivante : Créer un ensemble d'apprentissage (EA) et un ensemble de tests (ET) à partir des 1100 observations dans `fraud_data`. Par convention, on choisira une taille pour l'EA de deux tiers de la taille totale de l'échantillon, soit les observations 1 jusqu'à 734. Le reste sera utilisé comme l'ET, afin d'évaluer la performance des arbres de décision.

```
1 fraud_data_EA <- fraud_data[1:734,]  
2 fraud_data_ET <- fraud_data[735:1100,]
```

### 1.2.2 Création des arbres

Nous utilisons différentes méthodes d'apprentissage d'arbres de décision (`rpart`, `C5.0`, `tree`, `ctree`) pour modéliser la relation entre les variables et la déclaration de fraude. Les arbres ainsi générés sont ensuite visualisés pour une meilleure compréhension.

```
1 tree1 <- rpart(fraudulent~., fraud_data_EA, parms = list(split = "gini"))  
2 tree1.1 <- rpart(fraudulent~., fraud_data_EA, parms = list(split = "information"),  
3               control = rpart.control(maxdepth = 9))  
4 tree2 <- C5.0(fraudulent~., fraud_data_EA, param = list(split = "gini"))  
5 tree2.1 <- C5.0(fraudulent~., fraud_data_EA, param = list(split = "information"))  
6 tree3 <- tree(fraudulent~., fraud_data_EA)
```

Comme on peut l'observer, on utilise plusieurs méthodes de création d'arbres de décision. Bien que les 3 méthodes principales sont `rpart`, `C5.0` et `tree`, on utilise différents paramètres pour leur création.

Tout d'abord, on teste l'efficacité de changer le paramètre de création de l'arbre "gini" et "information". Il s'agit de deux méthodes pour interpréter l'information contenue dans le data frame et en déduire des méthodes de sélection et d'apprentissage afin de prédire la classe cible.

### 1.2.3 Visualisation des arbres

Dans toute la suite, on considérera la même méthode pour afficher les arbres de prédiction.

```
1 plot(tree2, type="simple", main="Arbre C5.0")  
2 plot(tree2.1, type="simple", main="Arbre C5.0")  
3 plot(tree3, main="Arbre tree", col="blue", fill="blue")  
4 text(tree3, pretty = 0)
```

Ainsi, on obtient les arbres suivants :



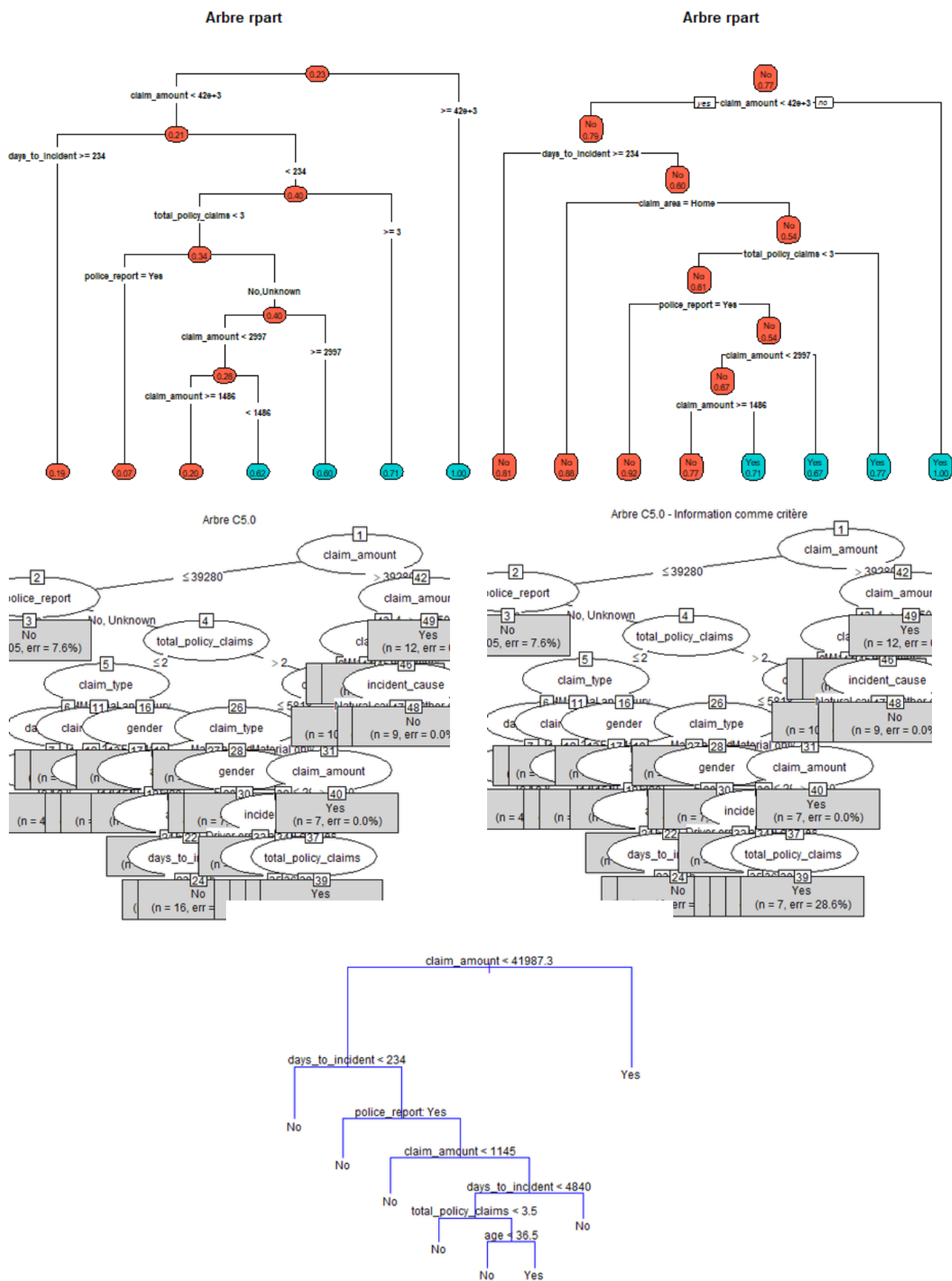


FIGURE 1.2 – Arbres de prédiction

## 1.3 Courbes ROC

Une fois que les arbres de prédiction sont créés, on peut les comparer grâce à leurs courbes ROC (DEF).

### 1.3.1 Fonction auxiliaire

Dans toute la suite, on affiche les courbes ROC à l'aide de cette fonction auxiliaire :

```
1 plot(ROC("rpart"), col = "green")
2 plot(ROC("rpart1"), col = "black", add = TRUE)
3 plot(ROC("C5.0"), col = "red", add = TRUE)
4 plot(ROC("C5.01"), col = "yellow", add = TRUE)
5 plot(ROC("tree"), col = "blue", add = TRUE)
6 legend(0.5, 0.5, legend=c("rpart", "rpart(info)", "C5.0", "C5.0(info)", "tree"),
7 col=c("green", "black", "red", "yellow", "blue"), lty=1:3, cex=0.8)
8 title(main="Courbes ROC")
```

### 1.3.2 Graphique des courbes ROC

Ainsi, on obtient les courbes ROC suivantes :

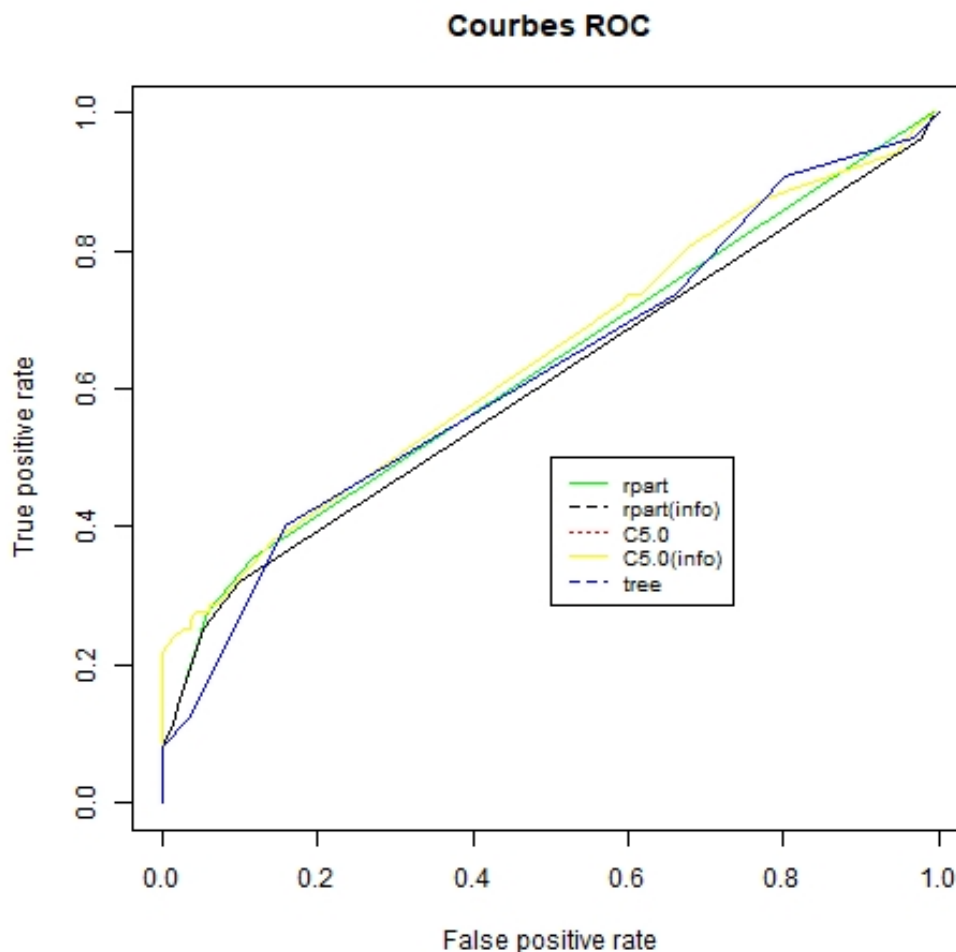


FIGURE 1.3 – Courbes ROC

On peut constater que toutes les courbes restent relativement inchangées. On s'attend donc à des arbres de prédiction peu performants et ce avec n'importe quel arbre choisi dans cette section. Pour déterminer la performance des arbres, on continue encore notre étude d'évaluation avec des tests d'air sous la courbe et des matrices de confusion.

## 1.4 Aire sous la courbe (AUC)

### 1.4.1 Fonction auxiliaire

Dans toute la suite du projet, on utilisera cette fonction qui nous aidera à calculer

```
1 AUC("rpart") #AUC : 0.630783174720883
2 AUC("rpart1") #AUC : 0.608659827792197
3 AUC("C5.0") #AUC : 0.646150043257941
4 AUC("C5.1") # AUC : 0.646150043257941
5 AUC("tree") #AUC : 0.625406830634862
```

### 1.4.2 Interprétation de l'AUC

On constate que toutes les AUC restent en moyenne de 0.631. Ceci peut s'interpréter comme des arbres de prédiction peu performants.

## 1.5 Mesure d'évaluation

### 1.5.1 Fonction auxiliaire et résultats

Dans toute la suite du projet, on utilisera cette fonction qui sert à afficher la matrice de transition associée à l'arbre de prédiction, ainsi que les différentes caractéristiques et propriétés de la matrice.

```
1 MC("rpart")
2 "Confusion Matrix and Statistics"
3
4           Reference
5 Prediction  No  Yes
6           No 263  63
7           Yes  16  24
8
9           Accuracy : 0.7842
10          95% CI : (0.7384, 0.8252)
11 No Information Rate : 0.7623
12 P-Value [Acc > NIR] : 0.1789
13
14           Kappa : 0.2684
15
16 McNemar's Test P-Value : 2.274e-07
17
18           Sensitivity : 0.27586
19           Specificity : 0.94265
20           Pos Pred Value : 0.60000
21           Neg Pred Value : 0.80675
22           Prevalence : 0.23770
23           Detection Rate : 0.06557
24           Detection Prevalence : 0.10929
25           Balanced Accuracy : 0.60926
26
27 'Positive' Class : Yes"
28 MC("C5.0")
29 "Confusion Matrix and Statistics"
30
31           Reference
32 Prediction  No  Yes
33           No 262  62
34           Yes  17  25
35
36           Accuracy : 0.7842
37          95% CI : (0.7384, 0.8252)
38 No Information Rate : 0.7623
39 P-Value [Acc > NIR] : 0.1789
40
```

```

41         Kappa : 0.2754
42
43     McNemar's Test P-Value : 7.407e-07
44
45         Sensitivity : 0.28736
46         Specificity : 0.93907
47         Pos Pred Value : 0.59524
48         Neg Pred Value : 0.80864
49         Prevalence : 0.23770
50         Detection Rate : 0.06831
51         Detection Prevalence : 0.11475
52         Balanced Accuracy : 0.61321
53
54     'Positive' Class : Yes"
55 MC("tree")
56 "Confusion Matrix and Statistics
57
58         Reference
59 Prediction   No  Yes
60      No      269  76
61      Yes      10  11
62
63         Accuracy : 0.765
64         95% CI : (0.7182, 0.8075)
65         No Information Rate : 0.7623
66         P-Value [Acc > NIR] : 0.4798
67
68         Kappa : 0.1226
69
70     McNemar's Test P-Value : 2.398e-12
71
72         Sensitivity : 0.12644
73         Specificity : 0.96416
74         Pos Pred Value : 0.52381
75         Neg Pred Value : 0.77971
76         Prevalence : 0.23770
77         Detection Rate : 0.03005
78         Detection Prevalence : 0.05738
79         Balanced Accuracy : 0.54530
80
81     'Positive' Class : Yes"
82 MC("rpart1")
83 "Confusion Matrix and Statistics
84
85         Reference
86 Prediction   No  Yes
87      No      264  65
88      Yes      15  22
89
90         Accuracy : 0.7814
91         95% CI : (0.7355, 0.8227)
92         No Information Rate : 0.7623
93         P-Value [Acc > NIR] : 0.2135
94
95         Kappa : 0.2482
96
97     McNemar's Test P-Value : 4.293e-08
98
99         Sensitivity : 0.25287
100        Specificity : 0.94624
101        Pos Pred Value : 0.59459
102        Neg Pred Value : 0.80243
103        Prevalence : 0.23770
104        Detection Rate : 0.06011
105        Detection Prevalence : 0.10109
106        Balanced Accuracy : 0.59956

```

```

107      'Positive' Class : Yes"
108 MC("C5.1")
109 "Confusion Matrix and Statistics
110
111      Reference
112 Prediction  No Yes
113      No    262  62
114      Yes    17  25
115
116      Accuracy : 0.7842
117      95% CI : (0.7384, 0.8252)
118      No Information Rate : 0.7623
119      P-Value [Acc > NIR] : 0.1789
120
121      Kappa : 0.2754
122
123      McNemar's Test P-Value : 7.407e-07
124
125      Sensitivity : 0.28736
126      Specificity : 0.93907
127      Pos Pred Value : 0.59524
128      Neg Pred Value : 0.80864
129      Prevalence : 0.23770
130      Detection Rate : 0.06831
131      Detection Prevalence : 0.11475
132      Balanced Accuracy : 0.61321
133
134      'Positive' Class : Yes"
135

```

### 1.5.2 Interprétation

Dans nos résultats, on se situe dans un cadre de détection de fraude. Ainsi, on cherche à maximiser la sensibilité (pour détecter le plus de fraudes possible) et aussi la spécificité (pour éviter de faussement détecter de la fraude). Ainsi, on peut voir que l'arbre C5.0 est le meilleur car il a la meilleure sensibilité et spécificité. De plus, il a la meilleure AUC. On peut aussi voir que les courbes ROC des arbres C5.0 et C5.0(info) sont les plus proches de l'angle supérieur gauche, ce qui est un bon indicateur de la qualité de l'arbre.

De ce fait, on s'intéresse donc à des méthodes et stratégies pour améliorer les performances de nos arbres de prédiction. Dans cette première section, on a pu observer un déséquilibre des classes. En effet, on dispose de 1100 observations avec 846 cas non frauduleux et le reste qui l'est. Ce déséquilibre peut être à l'origine de la mauvaise performance. On va donc procéder avec une stratégie de rééquilibrage des classes.

## Chapitre 2

# Méthode de rééquilibrage des classes

La méthode utilisée dans cette section consiste à rééquilibrer la classe fraudulent par une méthode de sur-échantillonnage. En effet, on cherche à augmenter le nombre de la casse minoritaire en simulant des observations.

### 2.1 Chargement de données

On se place dans le même cadre que dans la première section.

### 2.2 Principe du rééquilibrage

Ainsi, on utilise la fonction `ovun.sample` afin d'utiliser la méthode de sur-échantillonnage

```
1 fraud_data <- resample(fraud_data)
2 qplot(fraudulent, data=fraud_data, fill=fraudulent, geom="bar", main="Fraudulent",
  xlab="Fraudulent", ylab="Nombre de cas")
```

Une fois la librairie activée, on crée la fonction `resample` afin d'automatiser la méthode. De plus, avec la fonction `ovun.sample`, les observations obtenues sont ordonnées. Donc on utilise la fonction `sample` pour avoir nos données mélangées. Une fois nos données prêtes à être étudiées, on observe avec un `plot`<sup>2.1</sup> que, effectivement, nos données ont été rééquilibrées.

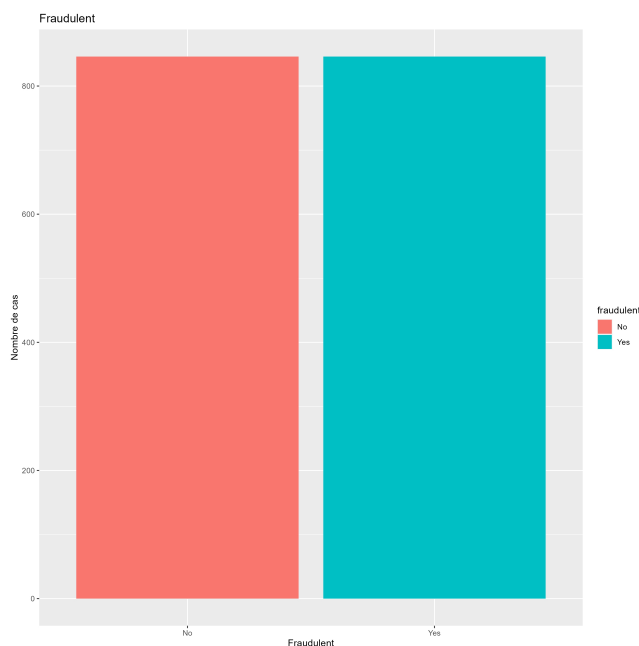


FIGURE 2.1 – Répartition de la classe rééquilibrée

Avec les paramètres correctes, on arrive à équilibrer parfaitement nos classes. Ainsi, on continue de la même manière que dans la section 1. On crée nos arbres de prédiction avec seulement `rpart`, `C5.0` et `tree`. Bien évidemment, on enlève systématiquement les variables ID.

```
1 fraud_data <- subset(fraud_data, select=-c(customer_id, claim_id))
```

## 2.3 Application

Création de l'ensemble de tests et l'ensemble d'apprentissage, on remarque qu'on dispose maintenant de 1692 observations.

```
1 fraud_data_EA <- fraud_data[1:1128,]
2 fraud_data_ET <- fraud_data[1129:1692,]
```

## 2.4 Arbres prédictifs et étude du modèle

On procède de la même façon que dans la section 1

### 2.4.1 Création des arbres et visualisation

```
1 tree1 <- rpart(fraudulent~., fraud_data_EA, parms = list(split = "gini"))
2 tree2 <- C5.0(fraudulent~., fraud_data_EA, param = list(split = "gini"))
3 tree3 <- tree(fraudulent~., fraud_data_EA)
```

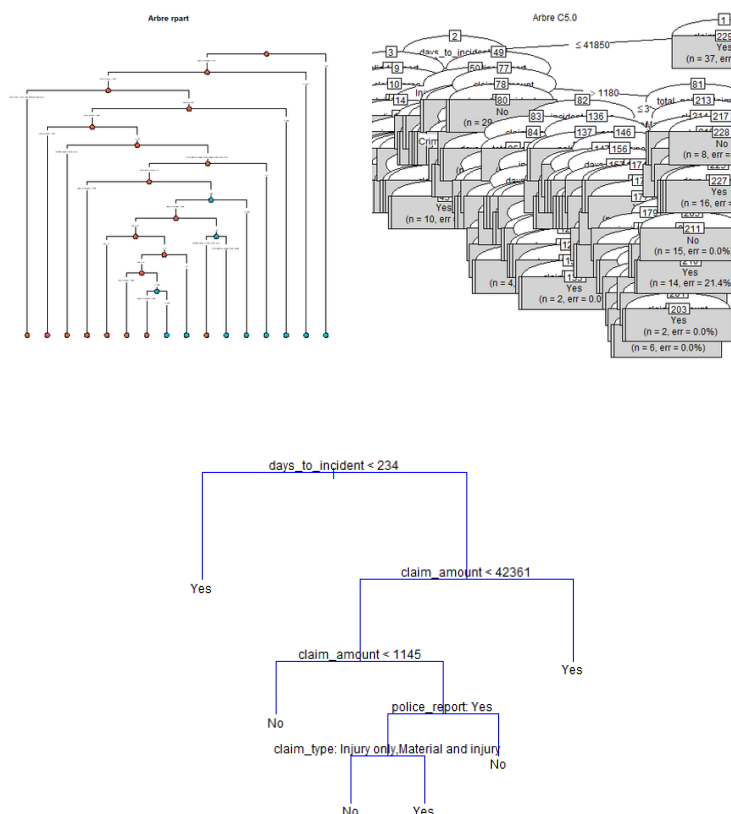


FIGURE 2.2 – Arbres de prédiction

On peut observer que la lecture des arbres est plus compliquée, étant donné que l'on dispose de plus d'observations.

### 2.4.2 Courbes ROC

On utilise la même fonction que dans la section 1, adapté pour seulement 3 arbres.

```
1 plot(ROC("rpart"), col = "green")
2 plot(ROC("C5.0"), col = "red", add = TRUE)
3 plot(ROC("tree"), col = "blue", add = TRUE)
```

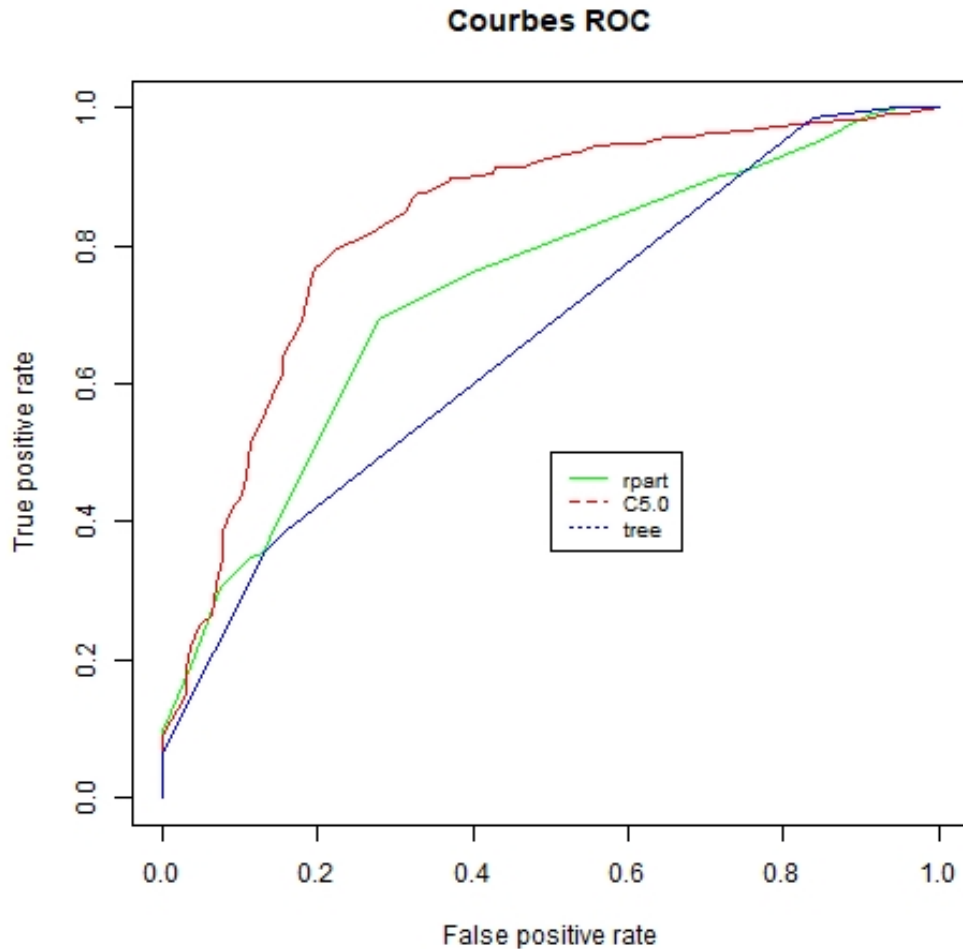


FIGURE 2.3 – Courbes ROC

On constate que l'arbre C5.0 est le plus performant parce que la courbe ROC est la plus proche du coin supérieur gauche, ce qui veut dire que le taux de vrais positifs est le plus élevé et le taux de faux positifs le plus faible. On peut aussi calculer l'aire sous la courbe des arbres pour confirmer cette hypothèse.

### 2.4.3 AUC

On utilise la même fonction que dans la section 1, modifiée pour seulement 3 arbres.

```
1 AUC("rpart") #AUC : 0.733895444756495
2 AUC("C5.0") #AUC : 0.829943453018148
3 AUC("tree") #AUC : 0.665100814830674
```

En effet, l'arbre C5.0 est le plus performant avec un AUC de 0.829943453018148. Si on le compare à l'AUC obtenu dans la section 1 (qui était de 0.646150043257941), on peut remarquer une amélioration, en pourcentage, de 28.4%.



## 2.4.4 Mesure d'évaluation

De même, cette fonction a été modifiée pour seulement 3 arbres

```
1 MC("rpart")
2 "Confusion Matrix and Statistics"
3
4           Reference
5 Prediction  No  Yes
6           No 195  90
7           Yes  76 203
8
9           Accuracy : 0.7057
10            95% CI : (0.6662, 0.743)
11    No Information Rate : 0.5195
12    P-Value [Acc > NIR] : <2e-16
13
14            Kappa : 0.4116
15
16    McNemar's Test P-Value : 0.313
17
18            Sensitivity : 0.6928
19            Specificity : 0.7196
20    Pos Pred Value : 0.7276
21    Neg Pred Value : 0.6842
22    Prevalence : 0.5195
23    Detection Rate : 0.3599
24    Detection Prevalence : 0.4947
25    Balanced Accuracy : 0.7062
26
27    'Positive' Class : Yes"
28 MC("C5.0")
29 "Confusion Matrix and Statistics"
30
31           Reference
32 Prediction  No  Yes
33           No 183  38
34           Yes  88 255
35
36           Accuracy : 0.7766
37            95% CI : (0.7399, 0.8103)
38    No Information Rate : 0.5195
39    P-Value [Acc > NIR] : < 2.2e-16
40
41            Kappa : 0.5494
42
43    McNemar's Test P-Value : 1.27e-05
44
45            Sensitivity : 0.8703
46            Specificity : 0.6753
47    Pos Pred Value : 0.7434
48    Neg Pred Value : 0.8281
49    Prevalence : 0.5195
50    Detection Rate : 0.4521
51    Detection Prevalence : 0.6082
52    Balanced Accuracy : 0.7728
53
54    'Positive' Class : Yes"
55 MC("tree")
56 "Confusion Matrix and Statistics"
57
58           Reference
59 Prediction  No  Yes
60           No 228 180
61           Yes  43 113
62
63           Accuracy : 0.6046
```

```

64          95% CI : (0.5629, 0.6452)
65      No Information Rate : 0.5195
66      P-Value [Acc > NIR] : 2.902e-05
67
68          Kappa : 0.2228
69
70      McNemar's Test P-Value : < 2.2e-16
71
72          Sensitivity : 0.3857
73          Specificity : 0.8413
74          Pos Pred Value : 0.7244
75          Neg Pred Value : 0.5588
76          Prevalence : 0.5195
77          Detection Rate : 0.2004
78          Detection Prevalence : 0.2766
79          Balanced Accuracy : 0.6135
80
81      'Positive' Class : Yes"

```

D'après les matrices de transition obtenues, on remarque que l'arbre C5.0 est le plus performant avec un taux de bonne classification de 77.66%. De plus, on remarque que l'arbre C5.0 est le plus performant en terme de sensibilité avec un taux de 87.03%. On peut critiquer que le nombre de faux négatifs est assez élevé (88) parmi les 338 cas positifs mais le nombre de faux positifs est assez faible (38) parmi les 271 cas négatifs. On peut donc conclure que l'arbre C5.0 est le plus performant pour la détection de fraude. Il nous reste à étudier la performance des arbres en questions de variables explicatives. En effet, on remarque dans le plot des arbres que la lecture est assez compliquée, étant donné le nombre d'observations et de variables explicatives. Ainsi, on peut se demander si toutes les variables explicatives ont un rôle pertinent dans la détection de fraude.

## Chapitre 3

# Tests de corrélation

Dans cette section, on vise à déterminer la significativité des variables explicatives par rapport à la variable de classe. Pour ce faire, on utilisera des tests de corrélation, ainsi que les graphiques en mosaïques. Le résultat attendu sera en fonction du test effectué et de la p-value obtenue.

### 3.1 Test de significativité

Pour ces premiers tests, on charge les données dans une variable, on rééquilibre les données et on enlève les deux variables ID.

#### 3.1.1 Test de Chi2 pour variables catégorielles

On utilise la fonction `chisq.test` qui nous permet de visualiser la p-value pour un test de corrélation de chi2. Si le test est réalisé pour niveau de 91%, on en conclue que si la p-value est inférieure à 0.05, alors cette variable est significative et donc corrélée avec la variable de classe. Ici, on a créé des fonctions pour afficher le résultat directement.

```
1 "pvchi2_table()  
2 police_report      claim_type incident_cause      gender      claim_area  
3 2.474865e-10 1.775015e-08 9.940190e-02 1.585209e-01 2.931374e-01  
4 > pvfish_table()  
5 police_report      claim_type incident_cause      gender      claim_area  
6 1.441902e-10 1.358592e-08 9.969912e-02 1.585007e-01 2.931189e-01"
```

Ainsi, on peut en déduire que les variables qui ne sont pas corrélées sont `gender`, et `claim area`.

#### 3.1.2 Test de Pearson, Spearman et Kendall pour variables numériques

De même, on fait un test de corrélation pour les variables catégorielles. On crée une fonction pour afficher le résultat directement. La seule différence c'est que pour pouvoir évaluer les variables numériques entre-elles, la variable de classe doit elle-même être numérique. On va donc changer les valeurs de la variable `fraudulent` et poser "Yes" = 1 et "Non" = 0. Afin de ne pas changer la variable où sont stockées nos données initiales, on crée une variable temporaire `fraud data copie`.

```
1 psk("age")  
2 '[1] "p-value Pearson = 4.61337284270551e-07"  
3 [1] "p-value Spearman = 2.48938561558996e-07"  
4 [1] "p-value Kendall = 2.73961583160717e-07" '  
5 psk("days_to_incident")  
6 '"[1] "p-value Pearson = 1.63154239370855e-06"  
7 [1] "p-value Spearman = 4.29831904665197e-11"  
8 [1] "p-value Kendall = 5.59674069164966e-11" '  
9 psk("claim_amount")  
10 '[1] "p-value Pearson = 0.455960214966064"  
11 [1] "p-value Spearman = 0.0703230077916243"  
12 [1] "p-value Kendall = 0.0703407348935491" '  
13 psk("total_policy_claims")  
14 '[1] "p-value Pearson = 0.00980647494827477"  
15 [1] "p-value Spearman = 0.0372912761703214"  
16 [1] "p-value Kendall = 0.0373492251937505" '
```

Avec une p-value  $< 0.05$ , on rejette l'hypothèse nulle, il y a donc une corrélation entre la variable et la variable cible

### 3.1.3 Test en mosaïque

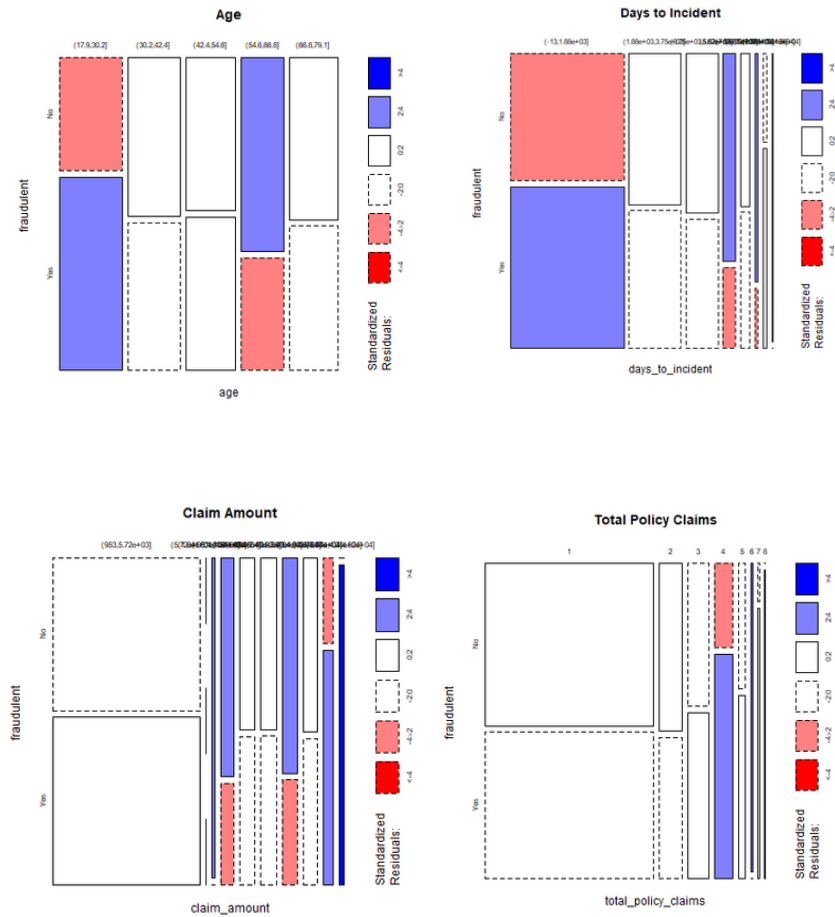


FIGURE 3.1 – Graphique en mosaïque

Ainsi, on a pu étudier la corrélation entre les variables et la variable classe étudiée. On peut remarquer certaines incohérences au niveau de la p-value calculée. De ce fait, notre critère se basera dans le graphique en mosaïque. On remarquera que si on al'apparition des rectangles rouges ou bleu, on gardera alors cette variable car ceci veut dire que, en plus d'être significative, elle joue un rôle important dans la caractérisation de la variable de classe.

## 3.2 Fonction AttrEval

On a une fonction prédefinie qui nous indique déjà le niveau de significativité par rapport à l'ensemble global de l'information dans nos données. Avec la fonction AttrEval on peut déterminer ce niveau.

```
1 Eval("GainRatio")
2 "total_policy_claims      police_report      days_to_incident      age
3      0.104655533          0.019078534          0.009190180          0.008606188
4      claim_type          claim_area          incident_cause          gender
5      0.007791333          0.003102022          0.001044547          0.000363893
6      claim_amount
7      -1.000000000"
8 Eval("Gini")
9 "claim_amount      police_report      age      days_to_incident
10      0.0303624371      0.0171140163      0.0136374392      0.0120551733
11      claim_type      total_policy_claims      incident_cause      claim_area
12      0.0070925484      0.0055454957      0.0016318665      0.0009763576
13      gender
14      0.0002522038"
15 Eval("Accuracy")
16 "claim_amount      age      police_report      days_to_incident
17      0.07210402      0.06560284      0.05910165      0.05614657
18      claim_type      total_policy_claims      incident_cause      claim_area
19      0.04314421      0.03546099      0.02127660      0.01300236
20      gender
21      0.01122931"
22 Eval("Relief")
23 "      age      days_to_incident      incident_cause      police_report
24      0.12529551      0.10283688      0.10224586      0.09929078
25 total_policy_claims      claim_amount      claim_type      gender
26      0.09574468      0.07978723      0.03250591      0.02836879
27      claim_area
28      0.01536643"
29 Eval("MDL")
30 "      claim_amount      police_report      age      days_to_incident
31      0.0382797154      0.0211697834      0.0113500490      0.0067407198
32 total_policy_claims      claim_type      claim_area      gender
33      0.0060765507      0.0059883394      -0.0005329839      -0.0020233076
34      incident_cause
35      -0.0060326374"
```

## 3.3 Application

### 3.3.1 Chargement et retraitement de données

Les données utilisées dans toute la suite seront rééquilibrées.

```
1 fraud_data <- read.csv("Donnees/Data_Projet_1.csv",
2                       header = TRUE, sep = ",", dec = ".", stringsAsFactors = T)
3 resample <- function(fraud_data){
4   resample <- ovun.sample(fraudulent ~ ., data = fraud_data, method = "over", N =
5     1.539 * length(fraud_data$fraudulent))
6   indices <- sample(nrow(resample$data))
7   resample <- resample$data[indices,]
8   return(resample)
9 }
10 fraud_data_sample <- resample(fraud_data)
11 "2/3 de 1692 pour fraud_data_EA et 1/3 pour fraud_data_ET"
12 fraud_data_EA <- fraud_data_sample[1:1128,]
13 fraud_data_ET <- fraud_data_sample[1129:1692,]
14 "On enlève les variables non significatives"
15 fraud_data_EA <- subset(fraud_data_EA, select=-c(customer_id, claim_id, claim_area,
16   gender, incident_cause))
```

### 3.3.2 Création des arbres et visualisation

```

1 tree1 <- rpart(fraudulent~., fraud_data_EA, parms = list(split = "gini"))
2 tree2 <- C5.0(fraudulent~., fraud_data_EA, param = list(split = "gini"))
3 tree3 <- tree(fraudulent~., fraud_data_EA)
4 rpart.plot(tree1, type=4, extra=7, box.col=c("tomato", "darkturquoise")[
  tree1$frame$yval], main="Arbre rpart")
5 plot(tree2, type="simple", main="Arbre C5.0")
6 plot(tree3, main="Arbre tree", col="blue")
7 text(tree3, pretty = 0)

```

Ainsi, on peut visualiser les arbres

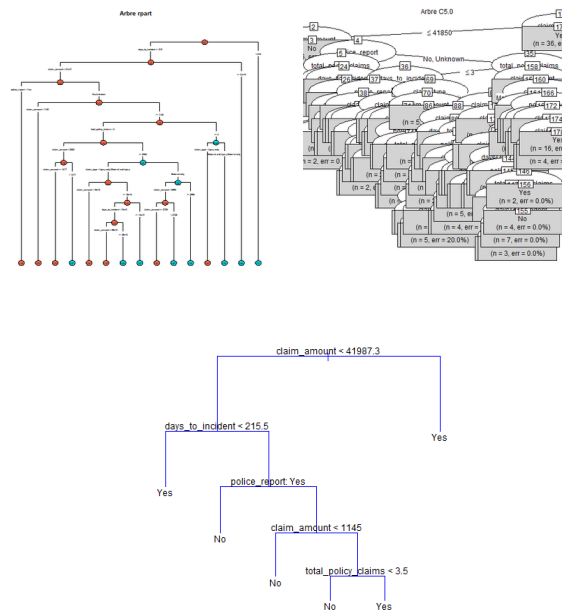


FIGURE 3.2 – Arbres de prédiction

### 3.3.3 Courbes ROC

On utilise la même fonction que dans la section 2.

```

1 plot(ROC("rpart"), col = "green")
2 plot(ROC("C5.0"), col = "red", add = TRUE)
3 plot(ROC("tree"), col = "blue", add = TRUE)
4 legend(0.5, 0.5, legend=c("rpart", "C5.0", "tree"), col=c("green", "red", "blue"),
  lty=1:3, cex=0.8)
5 title(main="Courbes ROC")

```

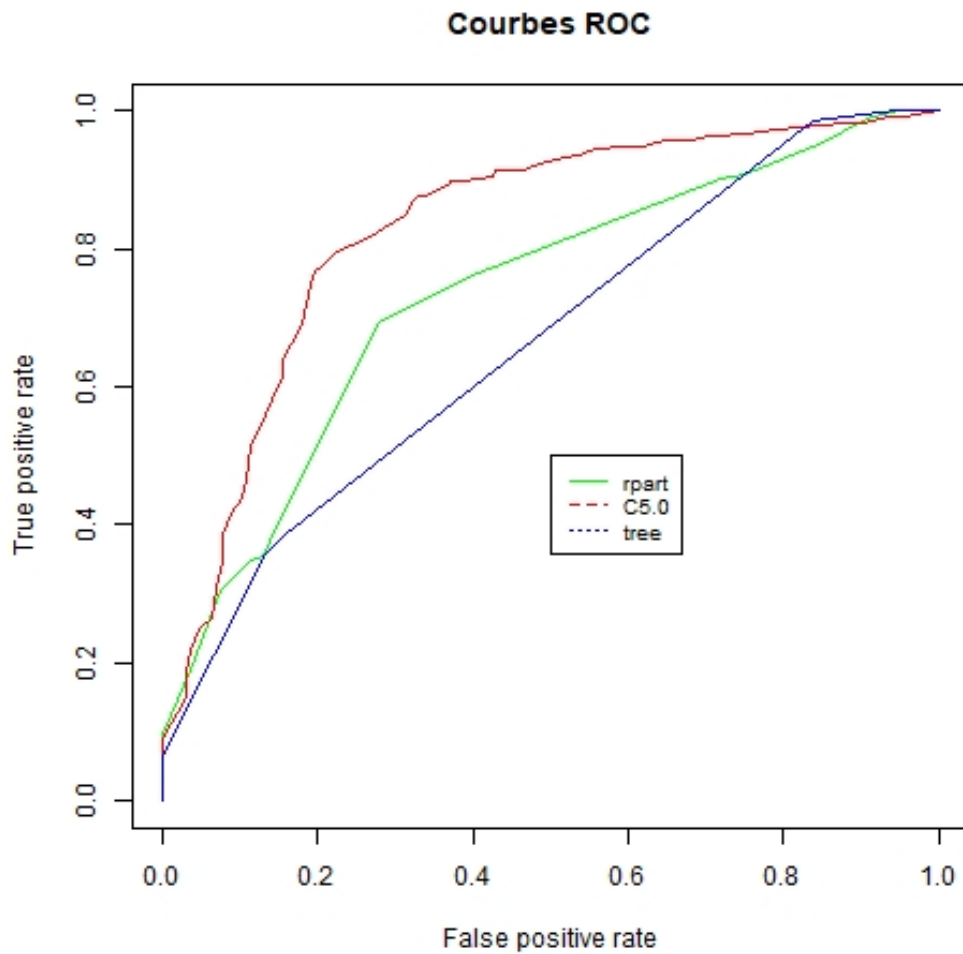


FIGURE 3.3 – Courbes ROC

### 3.3.4 AUC

On utilise la même fonction que dans la section 2.

```
1 AUC("rpart") #AUC : 0.746698943661972
2 AUC("C5.0") #AUC : 0.818228118712274
3 AUC("tree") #AUC : 0.679564889336016
```

### 3.3.5 Matrice de confusion

On utilise la même fonction que dans la section 2.

```
1 MC("rpart")
2 "Confusion Matrix and Statistics"
3
4      Reference
5 Prediction No Yes
6      No    209 104
7      Yes    71 180
8
9      Accuracy : 0.6897
10     95% CI : (0.6497, 0.7277)
11     No Information Rate : 0.5035
12     P-Value [Acc > NIR] : < 2e-16
13
14     Kappa : 0.3799
```

```

16  McNemar's Test P-Value : 0.01556
17
18      Sensitivity : 0.6338
19      Specificity : 0.7464
20      Pos Pred Value : 0.7171
21      Neg Pred Value : 0.6677
22      Prevalence : 0.5035
23      Detection Rate : 0.3191
24      Detection Prevalence : 0.4450
25      Balanced Accuracy : 0.6901
26
27      'Positive' Class : Yes"
28  MC("C5.0")
29  "Confusion Matrix and Statistics
30
31      Reference
32  Prediction  No  Yes
33      No      223  95
34      Yes      57 189
35
36      Accuracy : 0.7305
37      95% CI : (0.6918, 0.7667)
38      No Information Rate : 0.5035
39      P-Value [Acc > NIR] : < 2e-16
40
41      Kappa : 0.4615
42
43  McNemar's Test P-Value : 0.00269
44
45      Sensitivity : 0.6655
46      Specificity : 0.7964
47      Pos Pred Value : 0.7683
48      Neg Pred Value : 0.7013
49      Prevalence : 0.5035
50      Detection Rate : 0.3351
51      Detection Prevalence : 0.4362
52      Balanced Accuracy : 0.7310
53
54      'Positive' Class : Yes"
55  MC("tree")
56  "Confusion Matrix and Statistics
57
58      Reference
59  Prediction  No  Yes
60      No      235 167
61      Yes      45 117
62
63      Accuracy : 0.6241
64      95% CI : (0.5827, 0.6642)
65      No Information Rate : 0.5035
66      P-Value [Acc > NIR] : 5.493e-09
67
68      Kappa : 0.2505
69
70  McNemar's Test P-Value : < 2.2e-16
71
72      Sensitivity : 0.4120
73      Specificity : 0.8393
74      Pos Pred Value : 0.7222
75      Neg Pred Value : 0.5846
76      Prevalence : 0.5035
77      Detection Rate : 0.2074
78      Detection Prevalence : 0.2872
79      Balanced Accuracy : 0.6256
80
81      'Positive' Class : Yes"

```



Ainsi, en enlevant certaines variables non significatives, on arrive avec l'arbre C5.0 à obtenir moins de résultats faux négatifs, améliorant le fait de n'arrive pas à prédire les cas qui sont réellement positifs. C'est donc le résultat cherché. On pourrait donc se demander, en plus de la significativité individuelle des variables explicatives, leur comportement en tant que variables par groupes. Et donc, on procédera par une étude de clustering des variables.

# Chapitre 4

## Clustering

### 4.1 Chargement et prétraitement de données

On charge les données et on modifie la variable total policy claims pour la transformer en factor.

### 4.2 Matrice de distance et clustering par partitionnement

```
1 dmatrix <- daisy(fraud_data)
2 km4 <- kmeans(dmatrix, 5)
3 table(km4$cluster, fraud_data$fraudulent)
4 "      No Yes
5  1 391  38
6  2  91   0
7  3 135  52
8  4 124  39
9  5 105 125"
```

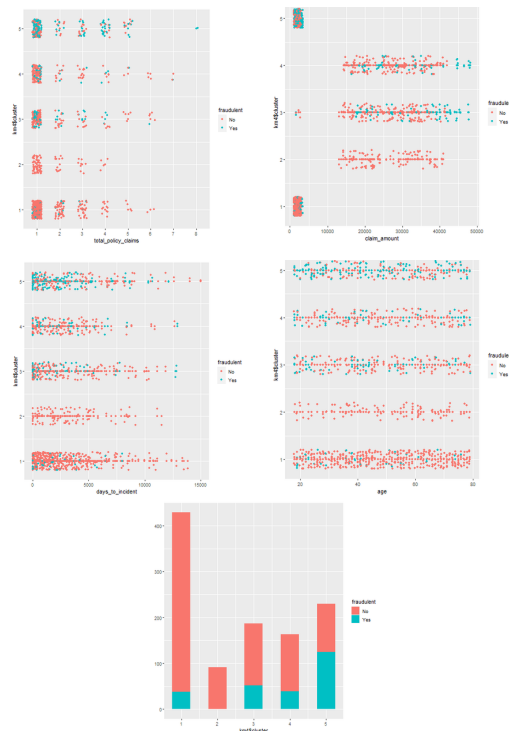


FIGURE 4.1 – Clustering

## 4.3 Application

L'idée dans cette section c'est de créer une variable de cluster. C'est-à-dire, en utilisant la méthode de k-means, avec un paramètre de 6, on peut créer 6 groupes avec lesquels on va associer la variable de classe. de Cette façon, on va pouvoir ajouter une variable qui est encore plus représentative que les autres variables explicatives. De plus, à l'aide de cette nouvelle variable, les arbres de prédiction gagneront en performance dans tous les aspects.

### 4.3.1 Arbres

```
1 tree1 <- rpart(fraudulent~., fraud_data_EA)
2 tree2 <- C5.0(fraudulent~., fraud_data_EA)
3 tree3 <- tree(fraudulent~., fraud_data_EA)
```

### 4.3.2 ROC

On utilise la même fonction que dans la section 2.

```
1 plot(ROC("rpart"), col = "green")
2 plot(ROC("C5.0"), col = "red", add = TRUE)
3 plot(ROC("tree"), col = "blue", add = TRUE)
```

### 4.3.3 AUC

On utilise la même fonction que dans la section 2.

```
1 AUC("rpart")
2 AUC("C5.0")
3 AUC("tree")
```

### 4.3.4 Matrice de confusion

```
1 MC("rpart")
2 "Confusion Matrix and Statistics"
3
4           Reference
5 Prediction  No  Yes
6      No    244  57
7      Yes    43 220
8
9           Accuracy : 0.8227
10           95% CI : (0.7886, 0.8533)
11      No Information Rate : 0.5089
12      P-Value [Acc > NIR] : <2e-16
13
14           Kappa : 0.645
15
16      McNemar's Test P-Value : 0.1936
17
18           Sensitivity : 0.7942
19           Specificity : 0.8502
20      Pos Pred Value : 0.8365
21      Neg Pred Value : 0.8106
22      Prevalence : 0.4911
23      Detection Rate : 0.3901
24      Detection Prevalence : 0.4663
25      Balanced Accuracy : 0.8222
26
27      'Positive' Class : Yes"
28 MC("C5.0")
29 "Confusion Matrix and Statistics"
30
```

```

31         Reference
32 Prediction  No Yes
33         No  257  28
34         Yes  30  249
35
36             Accuracy : 0.8972
37             95% CI : (0.8691, 0.921)
38         No Information Rate : 0.5089
39         P-Value [Acc > NIR] : <2e-16
40
41             Kappa : 0.7943
42
43         McNemar's Test P-Value : 0.8955
44
45             Sensitivity : 0.8989
46             Specificity : 0.8955
47             Pos Pred Value : 0.8925
48             Neg Pred Value : 0.9018
49             Prevalence : 0.4911
50             Detection Rate : 0.4415
51             Detection Prevalence : 0.4947
52             Balanced Accuracy : 0.8972
53
54         'Positive' Class : Yes"
55 MC("tree")
56 "Confusion Matrix and Statistics
57
58         Reference
59 Prediction  No Yes
60         No  242  55
61         Yes  45  222
62
63             Accuracy : 0.8227
64             95% CI : (0.7886, 0.8533)
65         No Information Rate : 0.5089
66         P-Value [Acc > NIR] : <2e-16
67
68             Kappa : 0.6451
69
70         McNemar's Test P-Value : 0.3681
71
72             Sensitivity : 0.8014
73             Specificity : 0.8432
74             Pos Pred Value : 0.8315
75             Neg Pred Value : 0.8148
76             Prevalence : 0.4911
77             Detection Rate : 0.3936
78             Detection Prevalence : 0.4734
79             Balanced Accuracy : 0.8223
80
81         'Positive' Class : Yes"

```

Dans tout notre projet, on a cherché à minimiser le taux de faux négatifs, aussi appelé sensibilité ( $TP/(TP+FN)$ ). Ainsi, on applique ce modèle aux données où on doit prédire la classe fraudulent. La méthode de clustering semble être la plus efficace pour donner à l'arbre C5.0 la capacité de prédire la classe.

## Chapitre 5

# Méthode de variables significatives et clustering

### 5.1 Test expérimental

Dans cette dernière partie, on utilise en même temps la méthode de clustering et la méthode de variables significatives. On charge les données, on les rééquilibre, on enlève les variables non significatives et on en extrait des clusters. Une fois les clusters formés, on ajoute cette nouvelle variable comme dans le chapitre 3 et on en déduit les arbres de prédiction.

```
1 dmatrix <- daisy(fraud_data)
2 kmi <- kmeans(dmatrix, 6)
3 fraud_data$cluster <- kmi$cluster
4 tree1 <- rpart(fraudulent~., fraud_data_EA)
5 tree2 <- C5.0(fraudulent~., fraud_data_EA)
6 tree3 <- tree(fraudulent~., fraud_data_EA)
7 # ROC
8 plot(ROC("rpart"), col = "green")
9 plot(ROC("C5.0"), col = "red", add = TRUE)
10 plot(ROC("tree"), col = "blue", add = TRUE)
11 AUC("rpart")
12 "Aire de l'arbre rpart est = 0.941854907682246"
13 AUC("C5.0")
14 "Aire de l'arbre C5.0 est = 0.965714141973135"
15 AUC("tree")
16 "Aire de l'arbre tree est = 0.908732454595764"
```

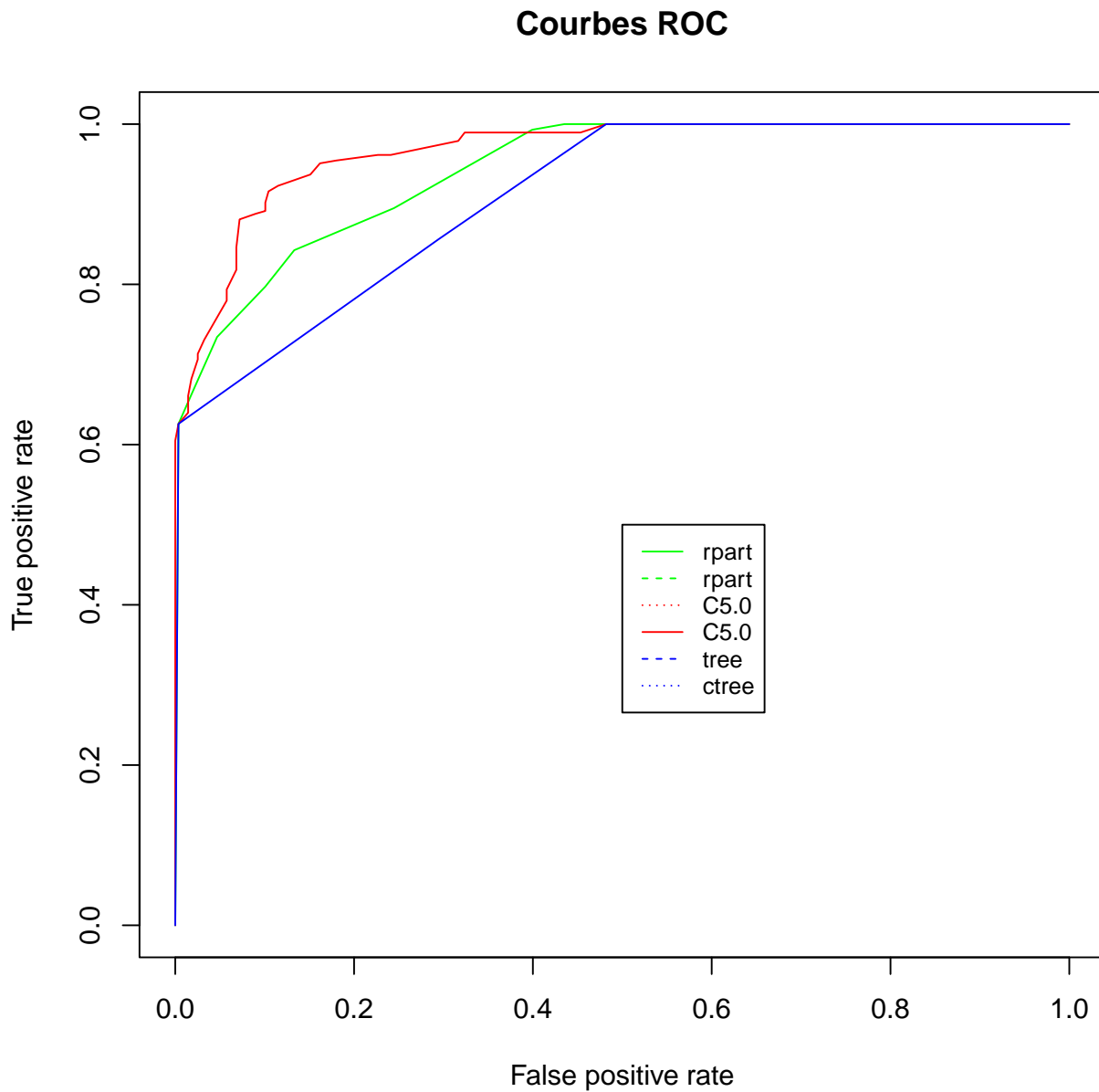


FIGURE 5.1 – Courbes ROC

```

1  "RPART
2  Confusion Matrix and Statistics
3
4      Reference
5  Prediction  No  Yes
6      No     241  45
7      Yes     37  241
8
9      Accuracy : 0.8546
10     95% CI : (0.8228, 0.8827)
11     No Information Rate : 0.5071
12     P-Value [Acc > NIR] : <2e-16
13
14     Kappa : 0.7093
15
16     McNemar's Test P-Value : 0.4395
17
18     Sensitivity : 0.8427
19     Specificity : 0.8669

```

```

20      Pos Pred Value : 0.8669
21      Neg Pred Value : 0.8427
22      Prevalence : 0.5071
23      Detection Rate : 0.4273
24      Detection Prevalence : 0.4929
25      Balanced Accuracy : 0.8548
26
27      'Positive' Class : Yes
28
29 C5.0
30 Confusion Matrix and Statistics
31
32      Reference
33 Prediction   No  Yes
34      No      246  22
35      Yes      32 264
36
37      Accuracy : 0.9043
38      95% CI : (0.8769, 0.9273)
39      No Information Rate : 0.5071
40      P-Value [Acc > NIR] : <2e-16
41
42      Kappa : 0.8084
43
44      McNemar's Test P-Value : 0.2207
45
46      Sensitivity : 0.9231
47      Specificity : 0.8849
48      Pos Pred Value : 0.8919
49      Neg Pred Value : 0.9179
50      Prevalence : 0.5071
51      Detection Rate : 0.4681
52      Detection Prevalence : 0.5248
53      Balanced Accuracy : 0.9040
54
55      'Positive' Class : Yes
56
57 TREE
58 Confusion Matrix and Statistics
59
60      Reference
61 Prediction   No  Yes
62      No      277 107
63      Yes       1 179
64
65      Accuracy : 0.8085
66      95% CI : (0.7736, 0.8402)
67      No Information Rate : 0.5071
68      P-Value [Acc > NIR] : < 2.2e-16
69
70      Kappa : 0.619
71
72      McNemar's Test P-Value : < 2.2e-16
73
74      Sensitivity : 0.6259
75      Specificity : 0.9964
76      Pos Pred Value : 0.9944
77      Neg Pred Value : 0.7214
78      Prevalence : 0.5071
79      Detection Rate : 0.3174
80      Detection Prevalence : 0.3191
81      Balanced Accuracy : 0.8111
82
83      'Positive' Class : Yes"

```

## 5.2 Interprétation

Dans un premier, on peut observer des changements au niveau des matrices de confusion. En effet, bien que le côté aléatoire de la fonction `ovun.sample` joue un rôle, on remarque que le nombre de faux négatifs est presque nul. Par contre, le nombre de faux positifs a augmenté, et la précision a diminuée. On remarque aussi pour l'arbre C5.0, qui était supposé être le plus performant, une légère augmentation dans le nombre des faux positifs et négatifs. Ainsi, on peut émettre l'hypothèse de que, même si on a fait un test de corrélation entre les variables explicatives et la variable de classe, on n'a pas étudié l'impact que les couples ou les groupes qui forment ces variables ont envers la variable de classe. On préférera donc garder l'arbre de prédiction avec la méthode de clustering sans enlever les variables non significatives.



## Chapitre 6

# Conclusion

En conclusion, on peut remarquer que la construction des arbres de prédiction a été faite de façon intuitive. Peu à peu on se rapprochait d'un modèle qui augmentait les performances avec des nouvelles méthodes telles que le rééquilibrage de classes, la suppression des variables non significatives et la construction des clusters. Pour pouvoir choisir un bon modèle, on a dû se concentrer sur les méthodes de classification telle que la matrice de confusion. On regarde le taux de succès et en particulier, le taux de faux négatifs. Dans le cadre de la détection de fraude, on a comme priorité la détection des possibles cas frauduleux et minimiser ainsi les coûts éventuels d'une fraude.

Dans le fichier annexe Data result.csv, on a utilisé notre arbre de prédiction décrit dans la section 4, l'arbre C5.0 afin d'utiliser les données disponibles dans le fichier Data projet 1 New et prédire la classe.

Cependant il faut remarquer que l'utilisation de l'arbre de prédiction C5.0 construite par notre méthode reste aléatoire à chaque génération de rééquilibrage de la classe. Ainsi, si on veut améliorer encore la prédiction, il est donc recommandable d'essayer en boucle jusqu'à trouver un arbre de prédiction, qui semble être toujours C5.0, pour lequel nos attentes de réduction de faux négatifs et une bonne précision sont accomplies.