

# Projet d'analyse de données : qu'est ce qui fait un bon étudiant ?

Rudio et Léo-Paul

2023-05-09

On a beaucoup d'idées reçues par rapport aux facteurs qui feraient d'un étudiant un bon étudiant, notamment sur l'alcool. Le jeu de données que nous avons étudié permet alors de confronter ces idées à des données concrètes.

## 1 Présentation du jeu de données.

Le jeu de données, nommé "Student alcohol consumption", est constitué d'informations sur la vie d'étudiants dans un lycée du Portugal. Ces informations vont de leur résultats universitaires ou de leur vie familiale à leur consommation d'alcool. Le jeu a été construit à partir d'une enquête menée auprès d'étudiants en mathématiques et en portugais.

L'objectif serait alors d'analyser le jeu de données afin de comprendre les facteurs qui impactent la réussite scolaire de ces étudiants. L'intérêt du jeu est la grande variété de facteurs proposée qui permet de couvrir un maximum d'hypothèses, notamment celles sur la consommation d'alcool proposée directement par le nom du jeu de données.

Voici les variables présentes dans ce jeu de données ;

- **school** - école (binaire: 'GP' - Gabriel Pereira ou 'MS' - Mousinho da Silveira)
- **sex** - sexe (binaire: 'F' - female ou 'M' - male)
- **age** - age (numérique: de 15 à 22)
- **address** - adresse (binaire: 'U' - urbain or 'R' - rural)
- **famsize** - taille de la famille (binaire : 'LE3' - inférieur ou égal à 3 or 'GT3' - supérieur à 3)
- **Pstatus** - parents qui habitent ensemble ? (binary: 'T' - living together or 'A' - apart)
- **Medu** - niveau d'études de la mère (numeric: 0 - none, 1 - primary education (4th grade), 2 - 5th to 9th grade, 3 - secondary education or 4 - higher education)
- **Fedu** - niveau d'études du père (numeric: 0 - none, 1 - primary education (4th grade), 2 - 5th to 9th grade, 3 - secondary education or 4 - higher education)
- **Mjob** - travail de la mère (nominal: 'teacher', 'health' care related, civil 'services' (e.g. administrative or police), 'at\_home' or 'other')
- **Fjob** - travail du père (nominal: 'teacher', 'health' care related, civil 'services' (e.g. administrative or police), 'at\_home' or 'other')
- **reason** - raison derrière le choix d'école (nominal: 'home', school 'reputation', 'course' preference ou 'other')
- **guardian** - représentant légal (nominal: 'mother', 'father' ou 'other')
- **traveltime** - temps de trajet (numérique en min)
- **studytime** - temps d'étude hebdomadaire (numérique en h)
- **failures** - nombre d'échecs (numeric: n if  $1 \leq n < 3$ , else 4)
- **schoolsup** - aide scolaire supplémentaire (binaire : yes ou no)
- **famsup** - support familial (binaire : yes ou no)
- **paid** - cours supplémentaires (binaire : yes ou no)
- **activities** - activités extra-scolaires (binaire : yes ou no)
- **nursery** - est allé à la crèche (binaire : yes ou no)

- **higher** - volonté de poursuite d'études (binaire : yes ou no)
- **internet** - accès à internet (binaire : yes ou no)
- **romantic** - en couple ? (binaire : yes ou no)
- **famrel** - états des relation familiale (numérique: de 1 - très mauvais à 5 - excellent)
- **freetime** - temps libre après les cours (numérique: de 1 - très mauvais à 5 - excellent)
- **goout** - sortie entre amis (numérique: de 1 - très bas à 5 - très élevé)
- **Dalc** - consommation d'alcool en semaine (numérique: de 1 - très basse à 5 - très élevée)
- **Walc** - consommation d'alcool le week-end (numérique: de 1 - très basse à 5 - très élevée)
- **health** - état de santé (numeric: from 1 - very bad to 5 - very good)
- **absences** - nombre d'absences (numérique: de 0 à 93)
- **G1** - note du 1<sup>er</sup> (numérique: de 0 à 20)
- **G2** - note du 2<sup>ème</sup> (numérique: de 0 à 20)
- **G3** - note du 3<sup>ème</sup> (numérique: de 0 à 20)

Au cours de ce projet, nous nous concentrons sur la moyenne des étudiants qui est à calculer et représente la note des élèves sur l'année. En addition, nous nous intéressons aussi à la réussite scolaire des élèves sur l'année qui va directement découler de leur moyenne. Il s'agirait donc ici d'étudier un problème de classification supervisé sur la réussite. Le but final serait alors d'avoir une meilleure compréhension des facteurs qui impacteraient la réussite scolaire et de les confronter à nos propres expériences en tant qu'étudiants.

Dans un premier temps, nous avons étudié chaque variable notamment leur corrélation avec la moyenne et la réussite. Ensuite, nous avons utilisé des méthodes et comparé des méthodes de Machine Learning dans le but de prédire la réussite des élèves.

Voici les bibliothèques à installer : ggplot2, FactoMineR, pROC, MASS, randomForest, gbm, gridExtra, dplyr, klaR, rpart, rpart.plot, corrplot, GGally, glmnet, e1071

## 2 Les données

### 2.1 Chargement des données

Le dataset est composé de 2 fichiers csv représentant les élèves de portugais et de maths. Il faut donc concaténer les deux jeux de données pour obtenir le jeu final. On peut noter qu'il y a 382 élèves qui suivent les deux cours.

### 2.2 Nettoyage et vérification des données

Afin d'adapter le jeu de données à notre étude, nous l'avons modifié. Nous avons notamment modifié en amont les variables *traveltime* et *studytime* afin de les rendre numérique.

On transforme les variables qualitatives en factor et on vérifie que le jeu ne contient pas de NaN.

Pour préparer concrètement les données, nous avons calculé la moyenne pour chaque élève (variable *Moy*), et nous avons rajouté une variable pour la réussite scolaire (variable *RS*). On garde 3 modalités différentes pour *RS* :

1. "admission" pour des Moyennes supérieures à 10
2. "redoublement" pour des Moyennes entre 8.50 et 10
3. "exclusion" pour des Moyennes inférieures à 8.50

On a choisi cette séparation étant donné qu'elle est plus intéressante à étudier qu'une simple variable binaire (on a essayé). Cette répartition a été calculée sur celle appliquée en France pour le lycée.

## 3 Exploration des données : analyse des variables

Cette partie consiste à appliquer d'abord des méthodes de statistiques descriptives afin de mieux comprendre le jeu de données et d'analyser les variables qui nous semblent intéressantes. La suite de l'analyse consiste

alors à vérifier la corrélation des variables avec la moyenne et la réussite scolaire. On présente dans cette partie les résultats sur les variables les plus intéressantes, le reste des résultats est disponible en annexe.

### 3.1 Les variables qualitatives

Pour chaque variable, nous étudions sa distribution avec soit un diagramme en bâton soit un diagramme circulaire. On effectue ensuite une ANOVA1 entre la variable et la moyenne pour en connaître l'impact à partir, notamment du test de Fisher. Un test de  $\chi^2$  est alors effectué pour vérifier la corrélation entre la variable et la réussite scolaire. S'il y a corrélation, on effectue alors une Analyse Factorielle des Correspondances (AFC).

#### 3.1.1 Les sorties

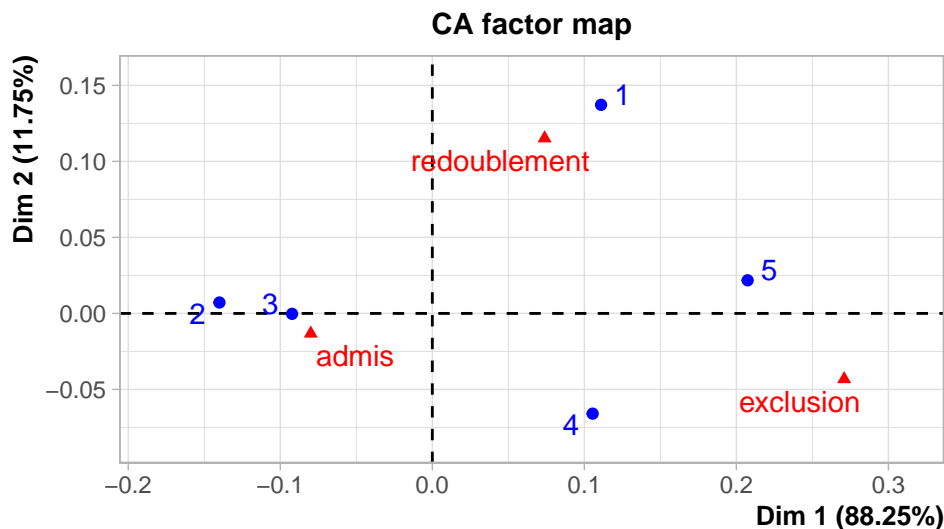
On remarque que les élèves maintiennent leur vie sociale. La grosse majorité sont intermédiaires en termes de sorties ce qui est quand même rassurant. Il y a quand même plus de personnes qui sortent vraiment beaucoup que de personnes qui ne sortent pas. Le test de Fisher indique les sorties sont très corrélées au notes et le test de Chi2 montre que la réussite scolaire est aussi corrélée aux sorties. Ainsi, on retrouve des résultats qui semblent cohérents et représentatifs de la vie étudiante.

Etant donné, la corrélation entre RS et goout, on peut effectuer une AFC pour préciser. On peut remarquer que ceux qui sortent peu-moyennement auront tendance à être admis alors que ce qui ne sortent pas (retrait/exclusion social) vont plutôt redoubler et les autres vont avoir tendances à se faire exclure. On obtient donc des résultats qui semblent plutôt pertinents.



```
##
## Call:
## lm(formula = Moy ~ goout, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.5887  -1.8876  -0.0015   2.1124   7.6652
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  10.5493    0.3770  27.980 < 2e-16 ***
## goout2       1.3727    0.4276   3.210  0.00137 **
## goout3       1.0049    0.4151   2.421  0.01564 *
## goout4       0.4522    0.4320   1.047  0.29548
## goout5      -0.1853    0.4517  -0.410  0.68178
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.177 on 1039 degrees of freedom
## Multiple R-squared:  0.02957,    Adjusted R-squared:  0.02583
## F-statistic: 7.915 on 4 and 1039 DF,  p-value: 2.766e-06
##
## Pearson's Chi-squared test
##
## data:  df$goout and df$RS
## X-squared = 20.537, df = 8, p-value = 0.008485
```

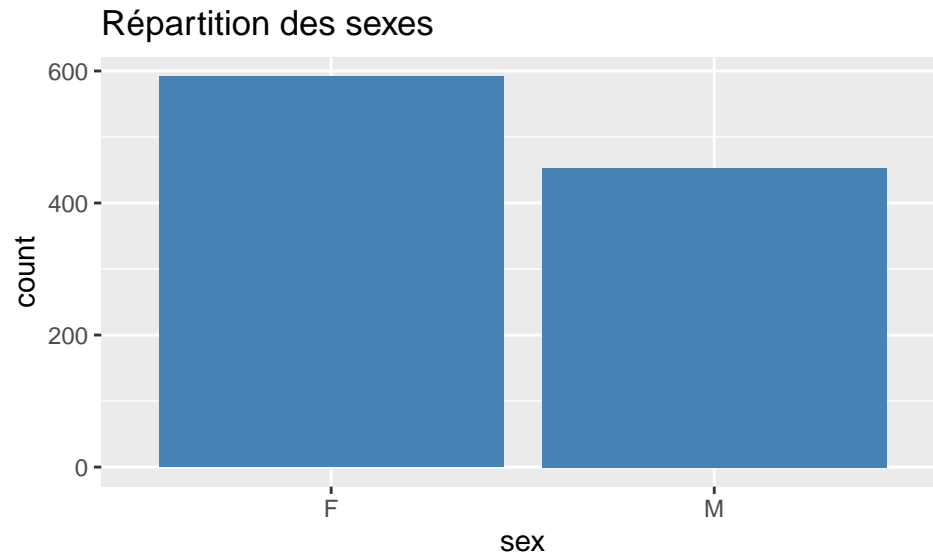


Avec un test du  $\chi^2$  on observe que les variables goout et RS sont corrélées (p-valeur petite devant 5%). Nous allons donc réaliser une AFC dessus. Egalement la p-valeur associée au test de fisher (sortie de anova) sur les variables Moy et goout montre que ces grandeurs sont aussi corrélées.

L'AFC nous montre ici que les étudiants qui sortent raisonnablement sont ceux qui réussissent le plus. En effet, ceux qui sortent le plus consacrent moins de temps à leur études ce qui peut expliquer ce résultat. Egalement les étudiants qui ne sortent quasiment pas échouent aussi beaucoup. Ce manque de sortie peut denoter d'un défaut de sociabilisation ou des problèmes de santé qui impact gravement la réussite de l'élève. Le diagramme en baton nous permet de voir que la majorité des étudiants sortent de manière modéré (modalité 3). L'AFC nous montre que cela n'est pas un frein à leur réussite.

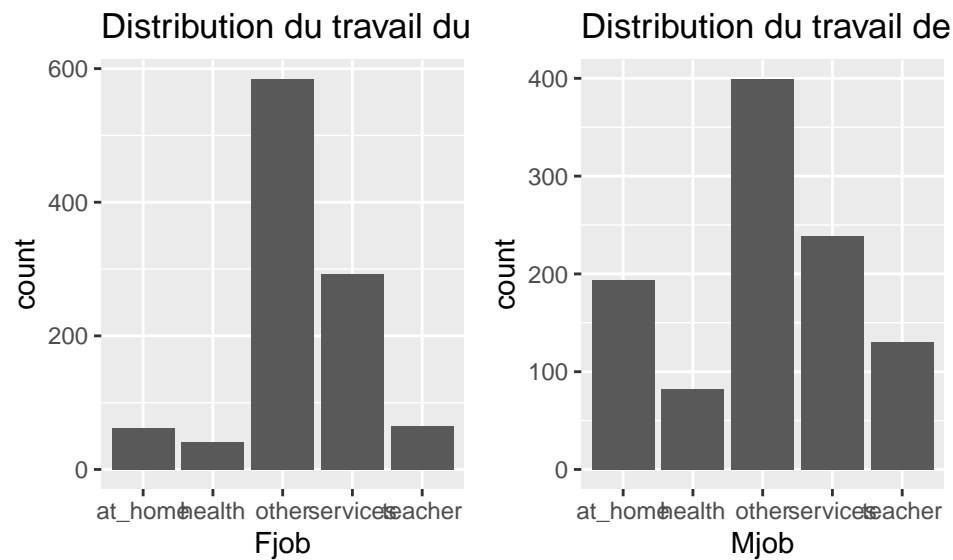
### 3.1.2 Le sexe des étudiants

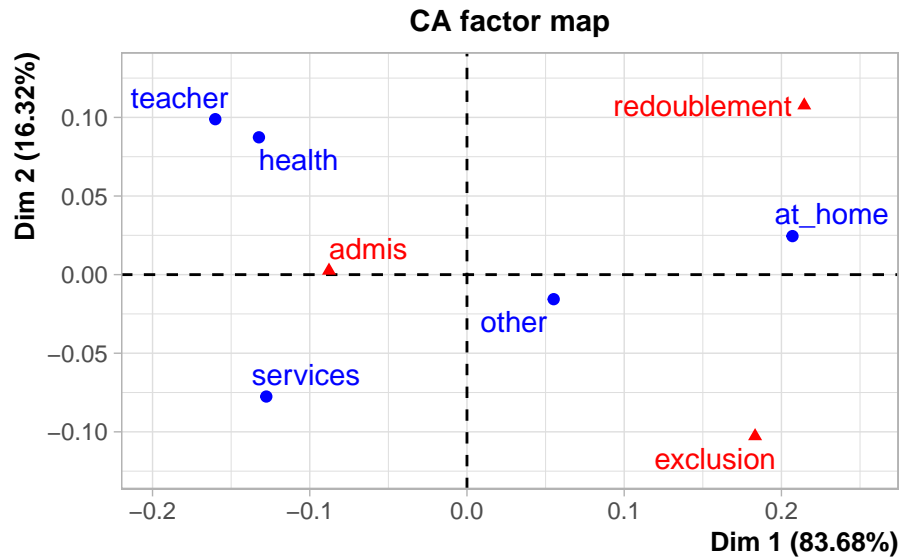
D'après le diagramme, le dataset est plutôt équilibré en terme d'hommes et de femmes, il y a même plus de femmes que d'hommes dans ce lycée. On étudie ensuite le lien entre le sexe et les notes en effectuant une ANOVA1. D'après le test de Fisher, p-value > 5% donc il n'y a pas d'effet du sexe sur les notes. D'après le test d'indépendances de Chi2 avec l'admission, le sexe des élèves n'a pas de lien avec leur réussite scolaire.



### 3.1.3 Travail des parents

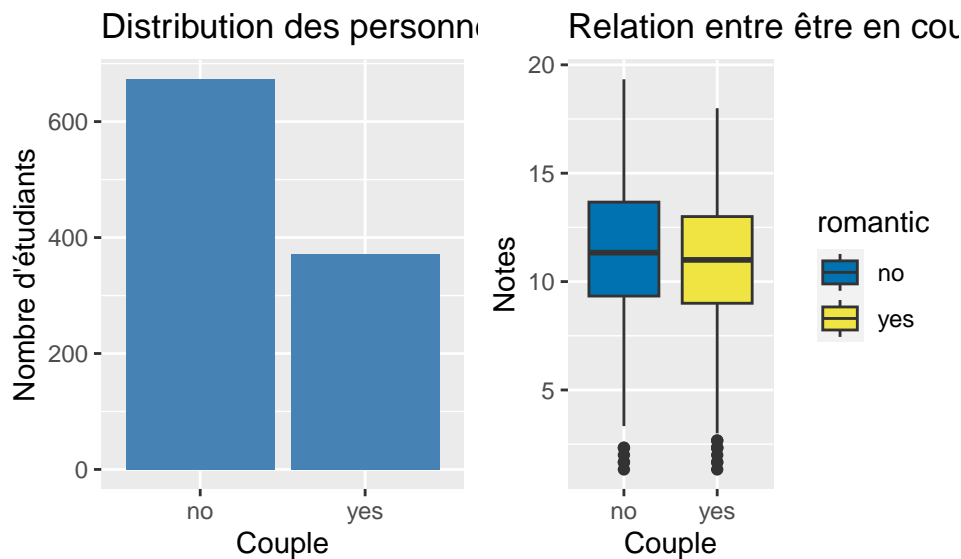
Dans les deux cas, others et services sont les catégories qui dominent. Une différence notable est la que la proportion de femme au-foyer est bien plus élevée que celle des hommes. D'après le test de Fisher, le travail de la mère a un impact sur les notes, contrairement à celui du père. Les résultats des test de Chis2 suivent les résultats des test de Fisher : le travail de la mère et la réussite scolaire sont bien corrélés mais celui du père n'a pas d'impact.





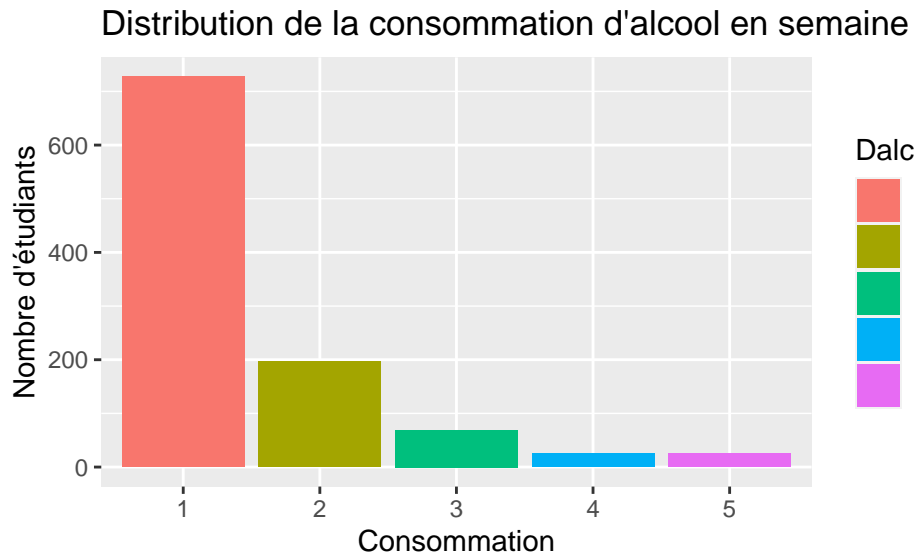
### 3.1.4 Les relations

Il y a environ deux fois plus de jeunes célibataires que de jeunes en couple. On peut penser qu'être en couple réduit le temps passé à étudier et rajoute des distractions, donc il devrait avoir un impact négatif sur les notes. D'après le test de Fisher, la p-value est fortement inférieure à 5%, donc on rejette  $H_0$ : il y a bien un lien entre situation romantique et notes, ce qui rejoint bien l'idée de départ. Il serait donc intéressant d'étudier la distribution des notes selon la situation romantique. D'après les boxplots, les différences sont assez minimes, même si on peut apercevoir que les notes des célibataires sont légèrement meilleures. Cependant, la présence de relation amoureuse n'a pas d'impact sur la réussite scolaire. Ainsi, être en couple fait baisser la moyenne mais n'est pas un facteur d'échec.

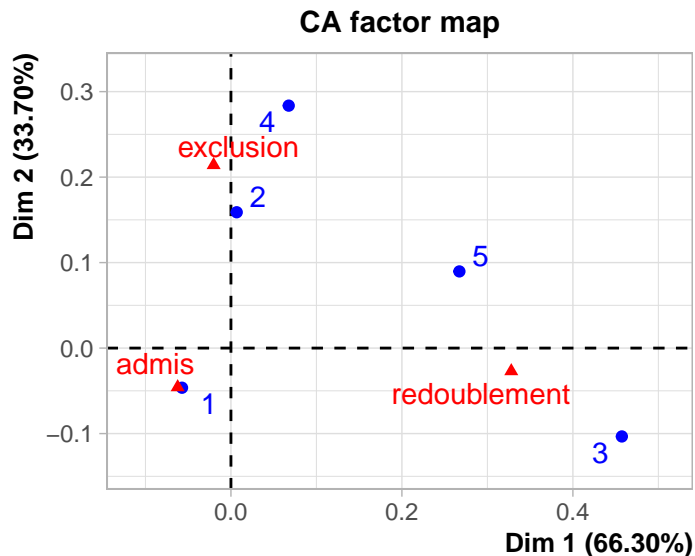


### 3.1.5 La consommation d'alcool

On s'intéresse enfin à la feature "principale" de ce jeu de données, la consommation d'alcool des étudiants.



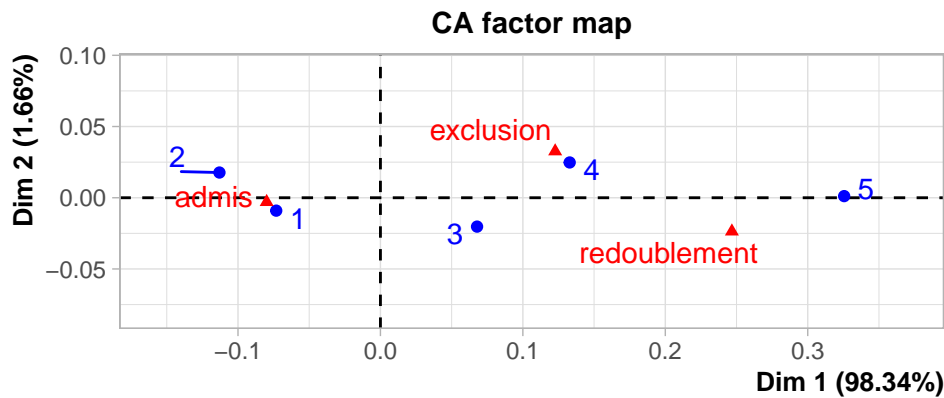
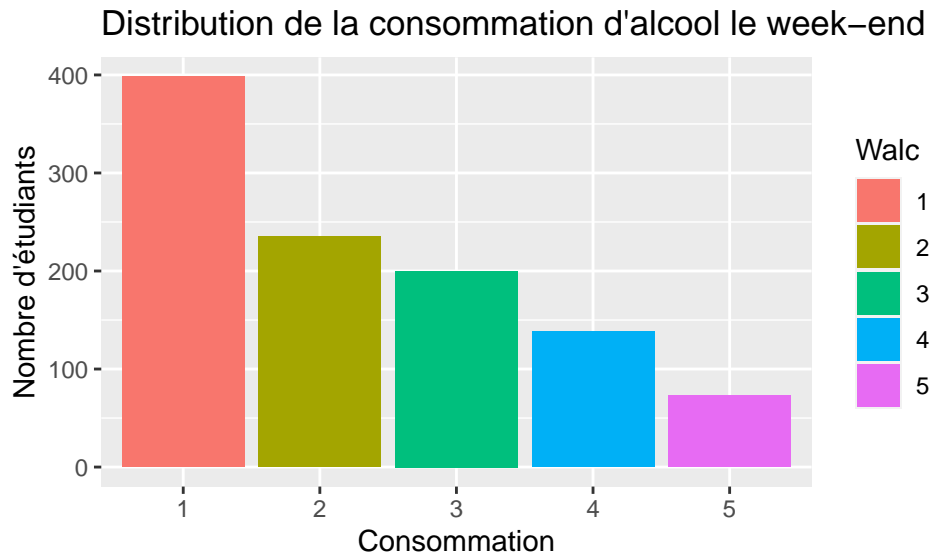
```
## Warning in chisq.test(df$Dalc, df$RS): Chi-squared approximation may be
## incorrect
```



Avec le test du  $\chi^2$  on voit que les variables Dalc et RS sont corrélées. On va donc réaliser une AFC dessus. De même avec la p-valeur du test de Fisher sur les variables Moy et Dalc, on voit que ces variables sont aussi corrélées (et c'est logique au vu du test du  $\chi^2$ ).

On voit très clairement avec l'AFC que les étudiants qui consomment le plus d'alcool sont ceux qui réussissent le moins. En effet, une forte consommation d'alcool témoigne d'un grand nombre de sortie ou bien d'un grave problème de santé (alcoolisme). Ceux qui réussissent le plus sont ceux qui consomment le moins d'alcool.

Avec le diagramme en bâton, on voit que la majorité des étudiants ne consomme quasiment pas d'alcool en semaine. L'AFC montre que cela n'as pas du tout été un frein pour leur réussite.



Tout d'abord on obtiens une p-valeur plus petite que 5% avec le test du  $\chi^2$  ce qui montre que les variables Walc (consommation alcool le week end) et RS (réussite scolaire) sont corrélées. Nous allons réaliser une AFC dessus afin de mieux les expliquer. De même avec le test de Fisher (réalisé à l'aide de l'anova) réalisé sur les variables Walc et Moy montre qu'elles sont corrélées.

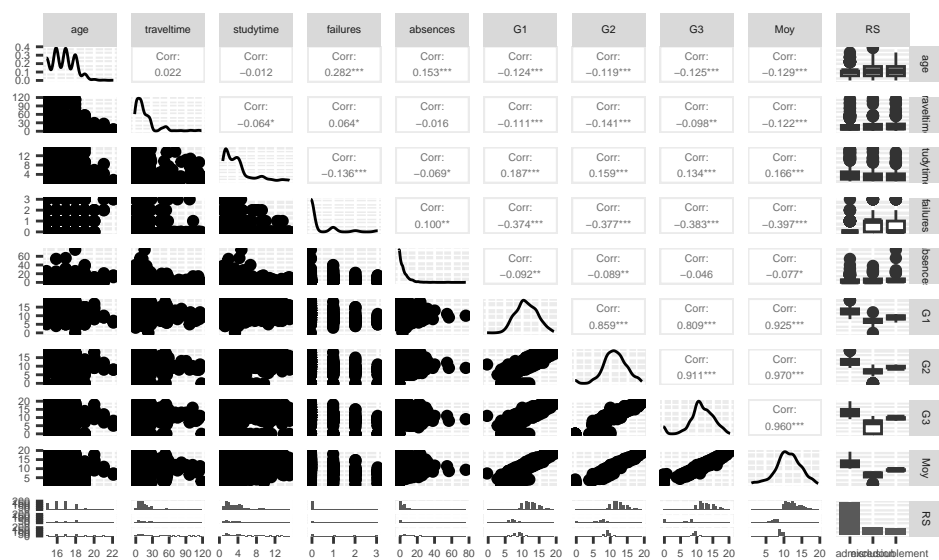
De même on obtiens le même résultat avec la consommation d'alcool le week end (ceux qui consomment le moins réussissent le plus), un peu plus nuancé cependant. En effet, on voit à travers les différents diagrammes en batons que globalement il y a plus d'étudiants qui consomment de l'alcool le week end qu'en semaine. On voit donc grâce aux deux AFC que les étudiants qui consomment plutôt de l'alcool le week end réussissent mieux que les étudiants qui consomment de l'alcool la semaine et le week end. Ainsi la variable avec la modalité 2 témoigne bien du fait que consommé de l'alcool en semaine est bien plus néfaste qu'en consommé en week-end (dans un contexte de soirée).

### 3.2 Les variables quantitatives

Dans cette partie, on s'intéresse aux variables quantitatives du jeu de données. De même que pour les variables qualitatives, on cherche à identifier les facteurs qui impactent la moyenne ainsi que la réussite scolaire.



### 3.2.1 Statistiques descriptives et corrélation



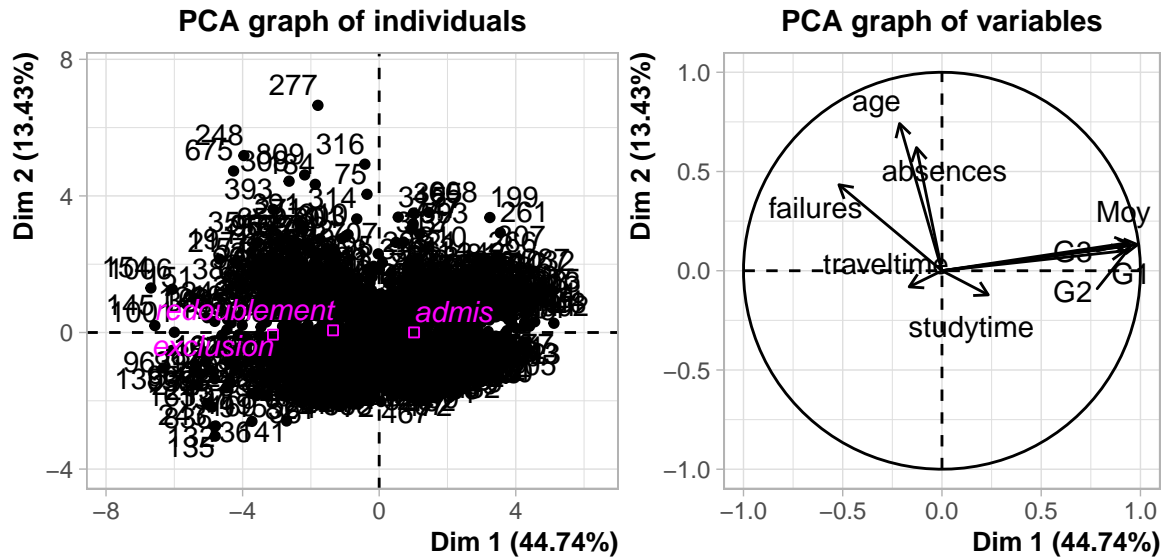
Le graphe nous montre la répartition bivariée des variables quantitatives ainsi que la corrélation entre les variables. On s'intéressera notamment à la corrélation par rapport à la moyenne.

On en déduit les informations suivantes:

- La majorité des étudiants ont entre 15 et 19 ans. L'âge est corrélé négativement avec la moyenne ce qui est plutôt surprenant.
- La plupart des étudiants ont moins de 30 minutes de temps de trajet. On a une corrélation négative pour le temps de trajet, ce qui paraît logique.
- Les étudiants travaillent généralement moins de 4h. Corrélation positive prévisible.
- globalement peu d'absences et d'échecs. Pas de corrélation pour les absences et forte corrélation pour les échecs.
- quasiment la même distribution de notes aux 3 semestres et donc de même pour la moyenne. Les élèves ont des moyennes qui tournent majoritairement entre 10-12. On remarque également que toutes les notes sont très fortement corrélées entre elles.

### 3.2.2 ACP

On effectue une Analyse en composantes principales (ACP) sur nos données quantitatives afin d'avoir plus d'informations sur nos données.



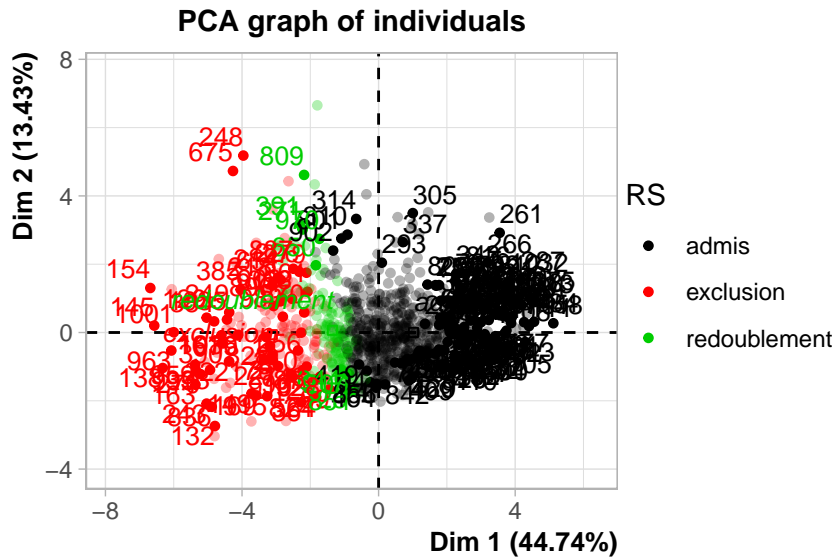
On voit que les variables study time et travel time sont très mal représentées, on ne peut donc pas les interpréter. Échec, âge et absences ne sont pas particulièrement bien représentées non plus, mais elles sont interprétables. En interprétant le cercle de corrélation, on remarque que le premier axe sépare les bons des mauvais élèves : bons élèves à droite. Le second axe sépare les élèves plus âgés et absents (vers le haut) des élèves moins âgés et plus assidus.

On peut afficher les pourcentages d'explications pour chacun des axes :

##	eigenvalue	percentage of variance	cumulative percentage of variance
## comp 1	4.026781e+00	4.474201e+01	44.74201
## comp 2	1.208335e+00	1.342594e+01	58.16796
## comp 3	1.011244e+00	1.123605e+01	69.40401
## comp 4	9.600833e-01	1.066759e+01	80.07160
## comp 5	8.730811e-01	9.700901e+00	89.77250
## comp 6	6.445814e-01	7.162015e+00	96.93451
## comp 7	1.961581e-01	2.179535e+00	99.11405
## comp 8	7.973555e-02	8.859506e-01	100.00000
## comp 9	1.041806e-30	1.157562e-29	100.00000

On remarque que seul le premier axe explique une grande partie de l'inertie, les axes 2,3 et 4 expliquent chacun 10% d'inertie. Il faut donc 4 axes pour expliquer 80% d'inertie dans notre cas.

On affiche alors le graphe des individus pour avoir une meilleure idée.



On distingue assez clairement les 3 groupes d'individus sur le graphe. De manière logique on retrouve que les élèves ayant une bonne moyenne vers la droite. De la même manière on voit que malgré tout le temps de trajet semble quand même avoir une influence négative sur la réussite. On pourra noter que les personnes qui vont être exclues sont celles qui ont le plus grand nombre d'échecs, ce qui est cohérent. On remarque que les élèves les plus âgés se dirigent vers un redoublement ou une exclusion, ce qui est plutôt curieux mais qui pourrait s'expliquer par la présence d'élèves redoublants en difficulté.

Malheureusement, à notre niveau, nous ne disposons pas de réel moyen d'évaluer l'influence des variables quantitatives sur la réussite scolaire. Nous partirons donc du principe que les variables qui impactent le plus les notes auront un réel impact sur la réussite scolaire. On pensera notamment aux échecs.

## 4 Machine Learning : Classification de la réussite scolaire

Dans cette partie, nous nous concentrons sur la mise en place de méthodes de classification afin de prédire la variable RS (réussite scolaire).

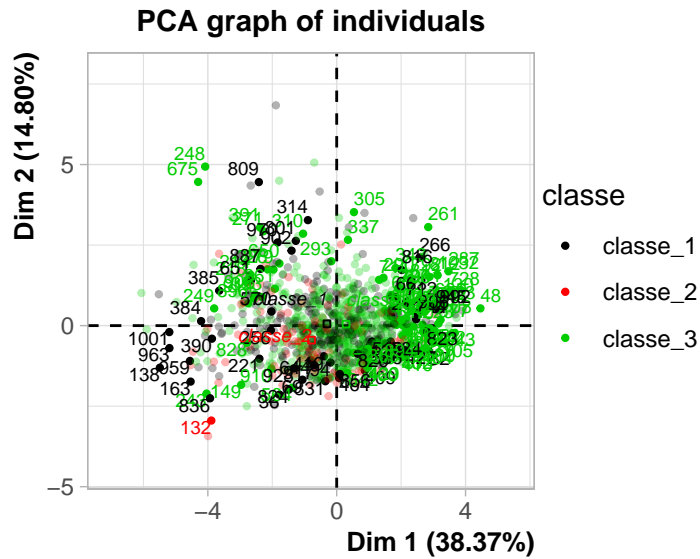
### 4.1 Classification non supervisée

L'intérêt d'effectuer de la classification non supervisé serait de voir comment seraient répartis les étudiants à partir des données quantitatives. On peut imaginer notamment que la classification pourrait s'effectuer sur les notes des élèves, mais selon quelles frontières ?

#### 4.1.1 Kmeans

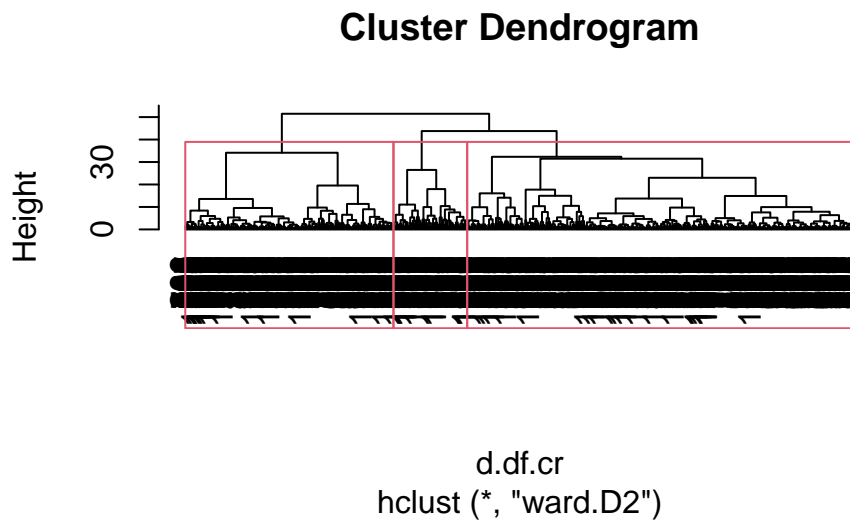
On applique l'algorithme de Kmeans avec 100 itérations et sans initialisation sur les données pour 3 groupes. Les résultats obtenus sont bien loin de la classification déjà présente. On en déduit donc que l'algorithme ne classe pas en fonction de la réussite au final. Le graphe d'ACP associé aux clusters montre qu'il n'y a pas de signification concrète à ces groupes

```
##
##      admis exclusion redoublement
##  1      217         66           45
##  2       58         19           24
##  3      448         88           79
```

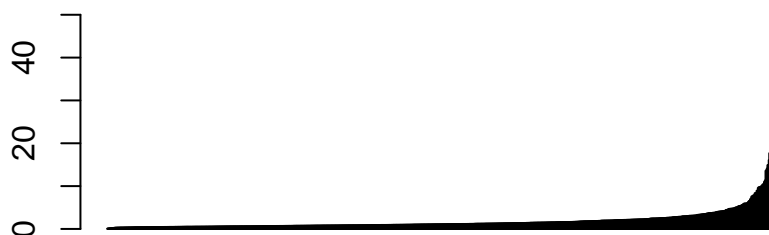


#### 4.1.2 CAH

On lance un algorithme de cah sur nos données quantitatives. On obtient le dendrogramme suivant duquel on extrait nos 3 classes qui nous intéressent.

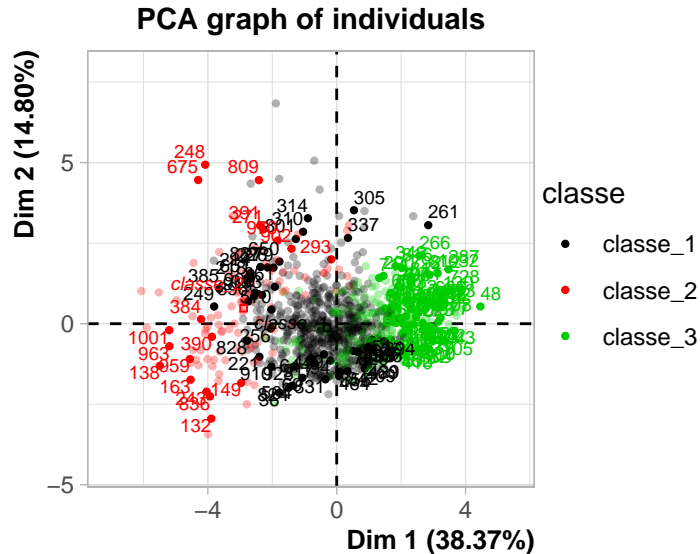


On peut alors observer le graphe des hauteurs pour chaque branche. En se basant sur la perte d'inertie, il est clair que l'on aurait gardé plus de 3 classes. Cela laisse penser que l'algorithme classerait pourrait donc raffiner les classes plus que nous l'avons déjà fait en distinguant 3 classes de réussite. Pour avoir une meilleure interprétation sur ces classes, on peut alors effectuer une ACP.



En observant l'ACP, on se rend compte que contrairement à la méthode de kmeans, les classes de la cah se distinguent bien sur le graphe des individus. On pourra également remarquer que la classification donnée est

similaire à celle que nous avons mis en place pour RS (ACP section 3.2.2) ce qui est assez intéressant. On pourrait donc penser que l'algorithme classe selon les notes, mais ce qui est le plus intéressant est que les frontières pour chaque groupe sont vraiment proches de celles que nous avons mis en place.



## 4.2 Classification supervisée

Il s'agit ici de la partie la plus intéressante : prédire la réussite scolaire d'un élève à partir de données. Notre objectif est donc de trouver un modèle qui aurait des résultats fiables pour cette tâche. Nous nous sommes donc essentiellement intéressé à la comparaison des résultats de chacune des méthodes. Les méthodes utilisées seront évaluées avec leur accuracy et leur courbe ROC.

Pour évaluer les modèles de manière plus précise, on calcule la moyenne d'accuracy sur  $N$  configurations de jeu de données et de jeu de test. Cela nous permet d'avoir des résultats plus généraux sur les performances des modèles.

Nous avons mis en place une procédure pour évaluer nos modèles. Afin de mettre en place notre procédure d'évaluation, nous avons donc implémenté des fonctions qui prennent en entrée le jeu d'entraînement et le jeu de test. Ces fonctions entraînent les modèles correspondants, effectuent les prédictions sur le jeu de test et renvoient une liste contenant l'accuracy, la table de confusion et la courbe ROC. Les modèles ont ensuite été évalués  $N=10$  fois avec à chaque fois une séparation différente dont le ratio est  $\frac{1}{5}$ . Cela permet d'avoir des résultats plus généraux étant donné que l'on teste les modèles dans plusieurs conditions différentes.

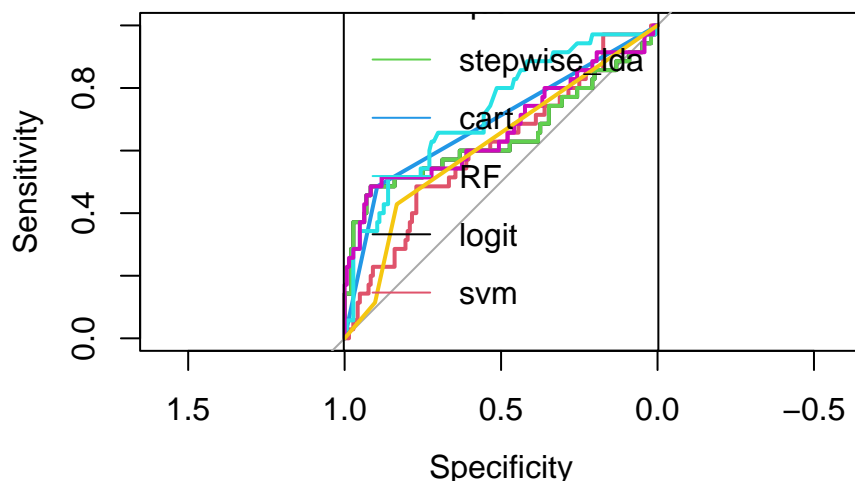
Voici les modèles que nous avons testés :

1. LDA
2. QDA
3. Stepwise lda
4. Random forest
5. Cart avec l'arbre optimal
6. Regression logistique
7. Support vector machine (bonus) avec un noyau radial avec  $c=10$ .

Notons également que nous avons retiré les variables de note du jeu de données puisque celles-ci réduiraient l'intérêt de la classification (dans notre cas la réussite scolaire sur l'année est calculée à partir des notes).

### 4.3 Comparaison

```
##          lda          qda stepwise_lda          cart          RF          logit
## accuracy 0.7200957 0.6626794 0.7177033 0.7153110 0.6985646 0.7224880
## AUC      0.6466270 0.6079365 0.6466270 0.6907738 0.7263889 0.6694444
##
##          svm
## accuracy 0.6339713
## AUC      0.6196429
```



Après calculs, on obtient des résultats plutôt décevants pour cette classification. Pour l'accuracy, les résultats tournent autour des 80% ce qui est assez faible. Le meilleur résultat s'obtient avec une régression logistique et le pire avec une SVM. Globalement, on en déduit que ces données ne se prêtent pas bien à la classification de la réussite scolaire dans leur état actuel.

## 5 Conclusion

Au final, ce jeu de données est vraiment intéressant étant donné qu'il renferme une grande quantité d'informations sur les étudiants. Il nous a permis d'appliquer une grande partie des méthodes statistiques apprises cette année de manière plus ou moins pertinente. On pourra noter également que la présence d'une grande quantité de variables qualitatives et de moins de variables quantitatives a rendu l'application de certaines méthodes plus compliquées. Cependant, globalement il s'agit d'un dataset intéressant pour apprendre à appliquer ces méthodes dans un cas moins trivial.

De plus, il nous a permis d'identifier les facteurs qui ont un impact sur les notes et la réussite scolaire des étudiants. On a ainsi pu observer des relations plutôt inattendues notamment par rapport à nos propres idées.

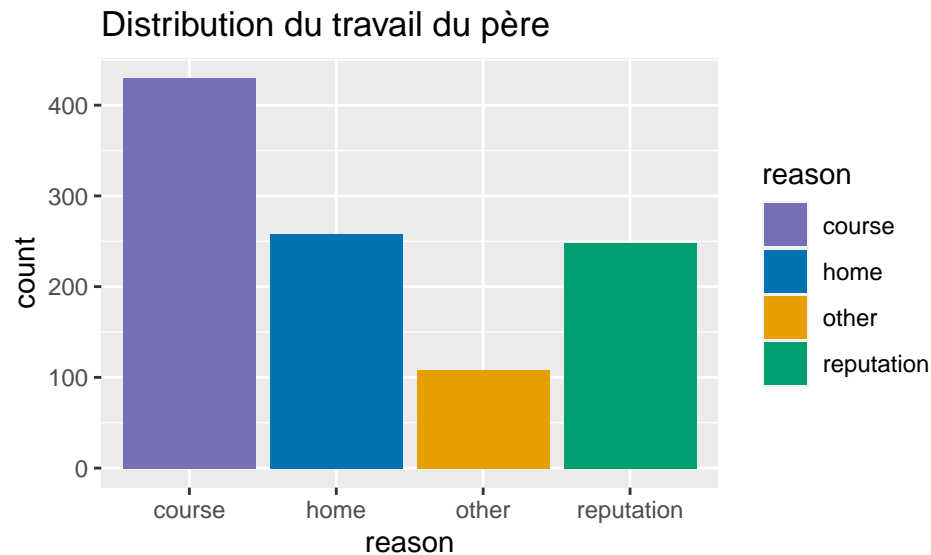
Cependant, les données du jeu ne sont pas forcément adaptées à la classification. On obtient des résultats assez bas pour toutes les méthodes utilisées.

## A Suite de l'étude des variables

### A.1 Les raisons du choix d'école

D'après le digramme circulaire, seule "other" possède un petit effectif alors que "course" domine. Ainsi, les élèves vont majoritairement en cours car ils les apprécient. D'après l'ANOVA1, il est clair que la raison d'aller en cours impacte les notes des étudiants ( $p$ -value  $< 5\%$ ). Cela paraît cohérent étant donné que cela détermine leur motivation à avoir de bonnes notes. De la même manière, la raison est bien corrélée avec la réussite scolaire, ce qui paraît bien cohérent.

```
# Distribution
ggplot(data = df, aes(x = reason, fill = reason)) +
  geom_bar() +
  labs(title="Distribution du travail du père") +
  scale_fill_manual(values = c("#7570b3", "#0072B2", "#E69F00", "#009E73", "#F0E442"))
```



```
# Lien avec les notes
summary(lm(Moy~ reason,data=df))
```

```
##
## Call:
## lm(formula = Moy ~ reason, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.3858  -1.8791  -0.0052   2.1209   7.7876
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   10.87907    0.15372  70.771 < 2e-16 ***
## reasonhome     0.45943    0.25103   1.830  0.0675 .
## reasonother    -0.03956    0.34309  -0.115  0.9082
## reasonreputation 1.17335    0.25417   4.616 4.39e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.188 on 1040 degrees of freedom
## Multiple R-squared:  0.02209,    Adjusted R-squared:  0.01927
## F-statistic: 7.832 on 3 and 1040 DF,  p-value: 3.587e-05
```

```
# Lien avec la réussite
chisq.test(df$reason,df$RS)
```

```
##
## Pearson's Chi-squared test
##
## data:  df$reason and df$RS
```

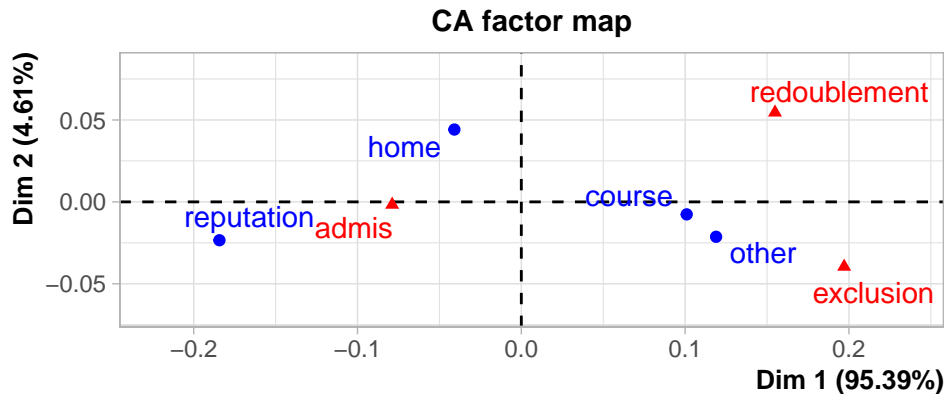
```
## X-squared = 15.479, df = 6, p-value = 0.01684
```

```
# AFC sur le travail de la mère
```

```
df.reason = data.frame(df$reason,df$RS)
```

```
table.reason = table(df.reason)
```

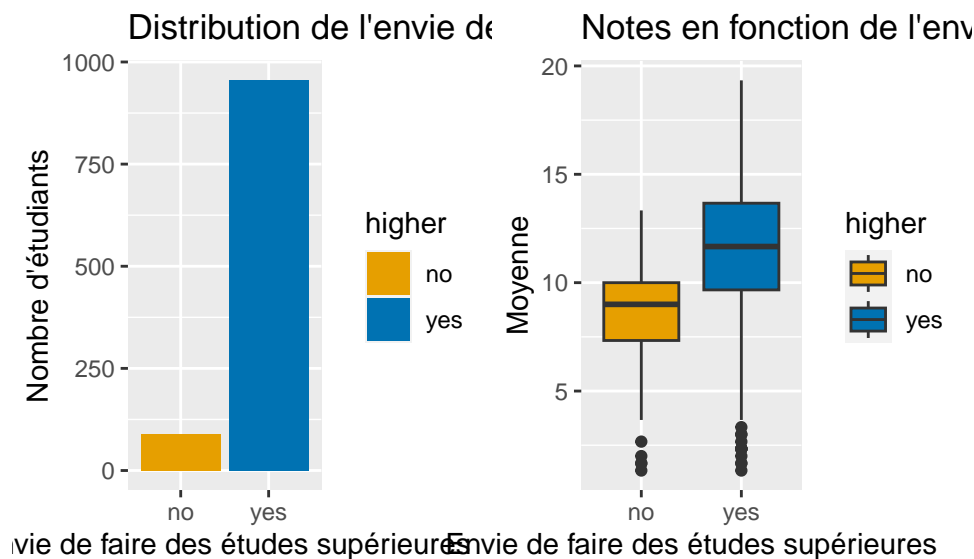
```
res = CA(table.reason)
```



On voit bien avec l'AFC que les personnes étant admises sont celles qui choisissent l'école pour sa réputation et sa proximité par rapport à leur domicile. A l'inverse on voit que les étudiants qui ont échoués sont ceux qui ont choisis l'école pour les cours ou d'autres raisons. On voit ici une des limites de cette méthode, en effet, on peut penser que les élèves qui réussissent le mieux sont ceux qui sont le plus motivés et donc qui ont choisis l'école pour les cours plus que pour sa réputation.

## A.2 Volonté de faire des études supérieures

On observe qu'au moins 80% des élèves veulent continuer leur études après le lycée, ce qui est plutôt rassurant. De plus, d'après le test de Fisher, les deux variables sont corrélées. On peut également annoncer que ceux qui veulent faire des études supérieures tendent à avoir de meilleures notes grâce au test unilatéral. A priori, la volonté de faire des études supérieures est corrélée à la réussite scolaire. Donc, ceux qui veulent poursuivre leurs études auront de meilleures notes et tendance à ne pas être en échec.

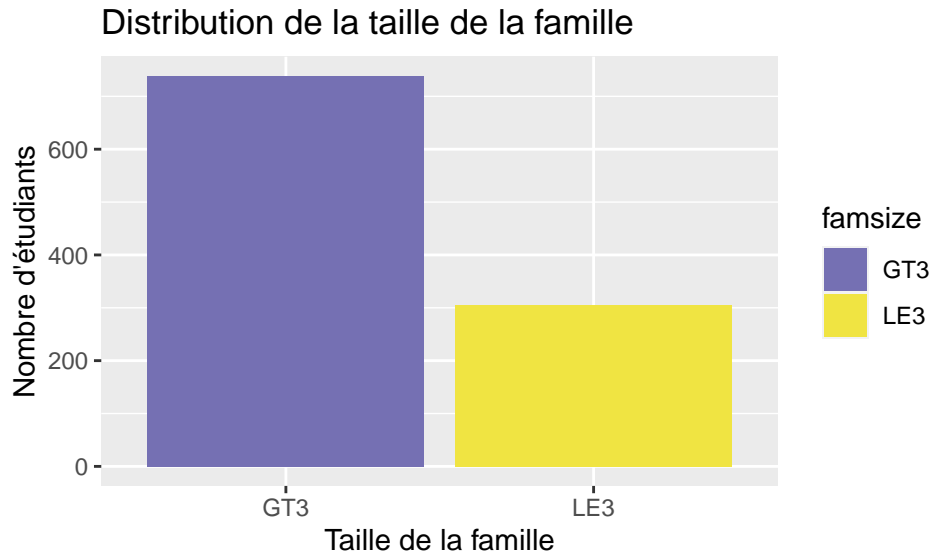


## A.3 La taille de la famille

On a deux fois plus de grandes familles que de petites familles. D'après le test de Fisher, il y a bien un impact de la taille de la famille sur les notes. Le test d'indépendance avec la réussite indique cependant que la

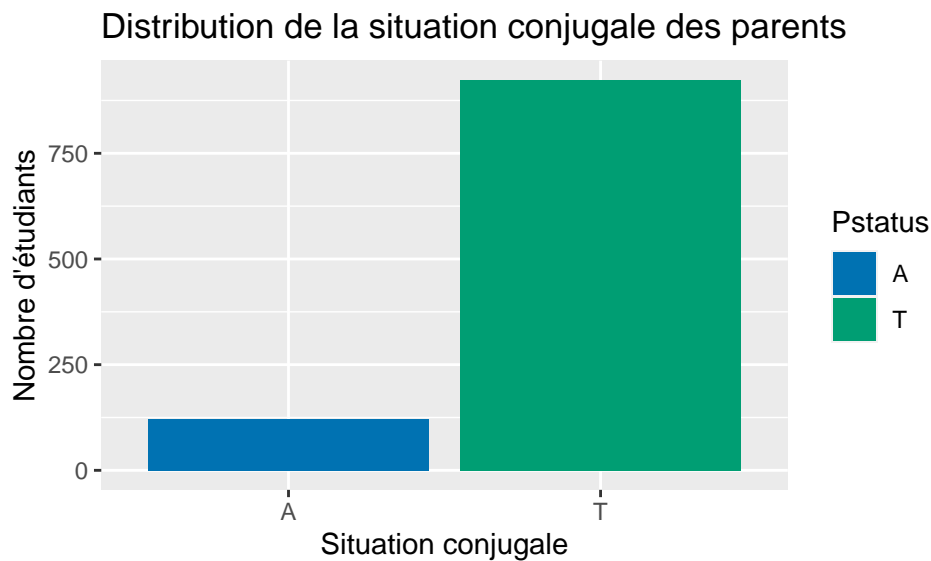


taille de la famille n'est pas liée à la réussite scolaire.



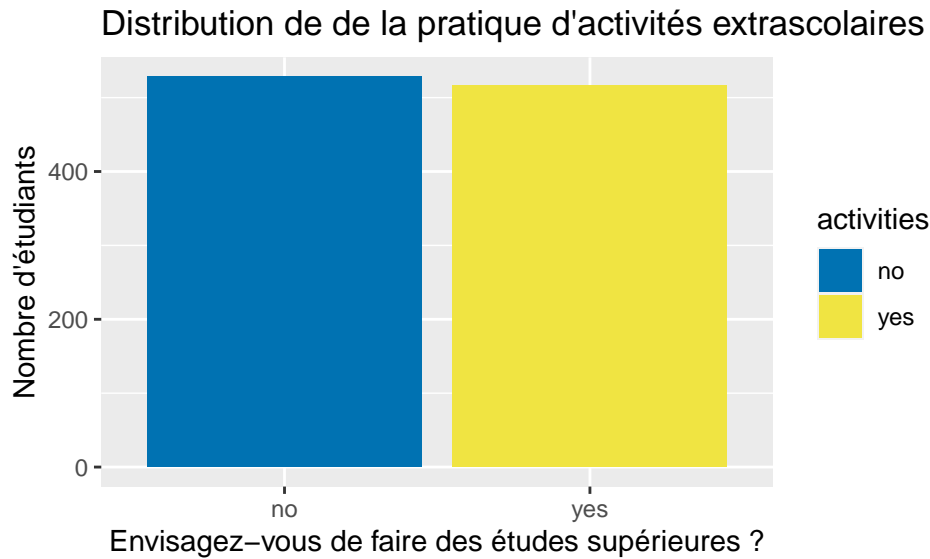
#### A.4 Situation familiale : séparation des parents

Le jeu est très déséquilibré au sujet de la situation famille : il y a 4 fois plus d'étudiants qui ont leurs parents qui vivent ensemble. De plus, le test de Fisher indique que la situation familiale n'a pas d'impact sur les notes. Le test de Chi2 soutient que le status des parents et la réussite scolaire sont indépendants.



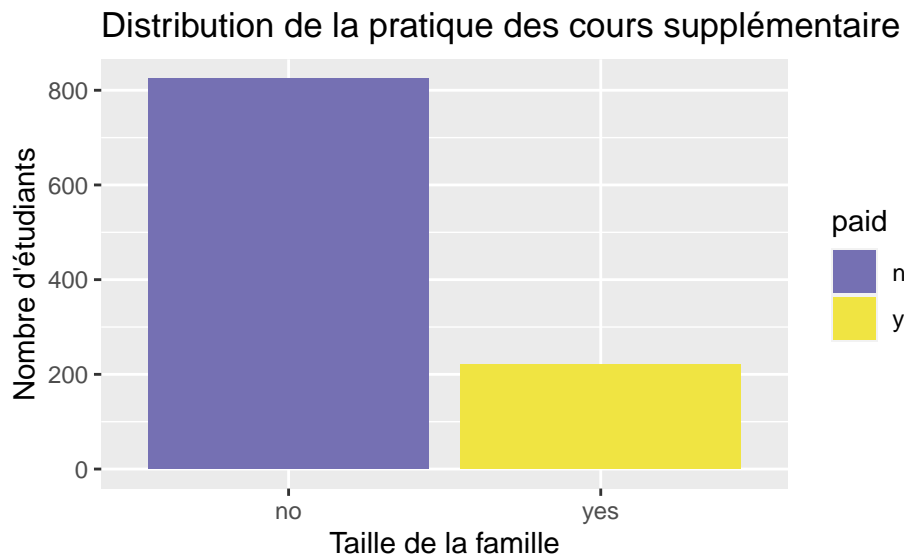
#### A.5 Activités extrascolaires

On a autant d'élèves qui pratiquent des activités extrascolaires que d'élèves qui n'en pratiquent pas, ce qui est plutôt intéressant. De plus, le test de Fisher indique plutôt qu'il n'y a pas de liens entre les activités extrascolaires et les notes, ce qui est plutôt surprenant étant donné que l'on aurait tendance à penser que les étudiants ayant des activités, ont moins de temps pour étudier. Dans la même lignée, les activités sont plutôt indépendantes de la réussite d'après le test de Chi2.



## A.6 Cours supplémentaires

Il y a bien plus d'élèves qui ne suivent pas de cours supplémentaires que d'élèves qui en suivent. Cette distribution est cohérente avec l'idée qu'on peut se faire. Le test de Fisher indique plutôt que les suivis de cours supplémentaires n'a pas d'impact sur la moyenne. De même, le suivi de cours supplémentaire n'est pas lié à la réussite.



On voit que la répartition d'âge est la même dans chaque filière

## A.7 Quantité de travail

On voit clairement que les étudiants travaillent majoritairement moins de 2h00 ou entre 5h00 et 10h00 par semaines.

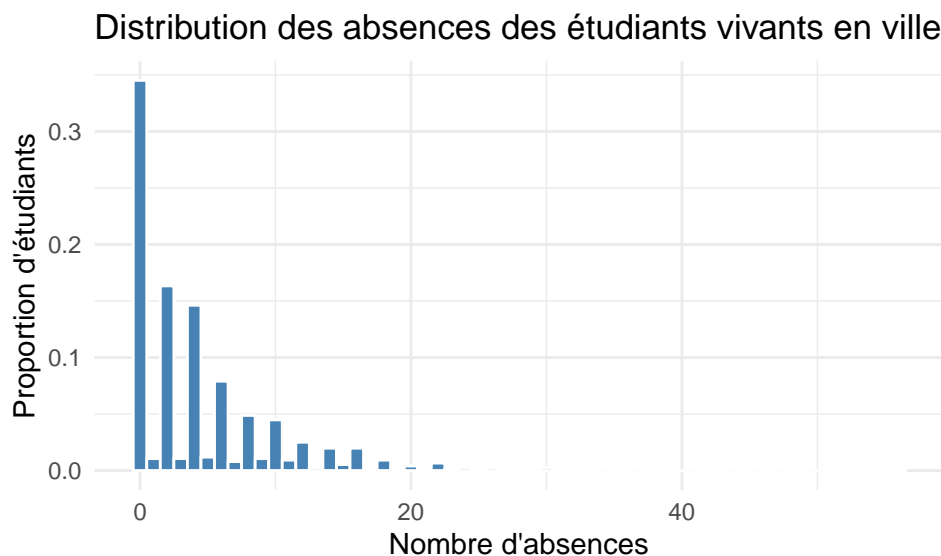
On voit qu'il y a plus de personnes qui travaillent moins de deux heures par semaine dans la section portugaise tandis qu'il y a moins de personnes qui travaillent plus de 10h00 dans cette même section. Le nombre d'étudiants travaillant entre 5 et 10 heures semble être à peu près le même. En effet:

On s'aperçoit donc que les élèves dans la filière mathématiques travaillent plus

On voit que globalement, les élèves qui travaillent plus ont de meilleures notes (comportement bizarre à vérifier)

#### A.7.1 Absences des étudiants

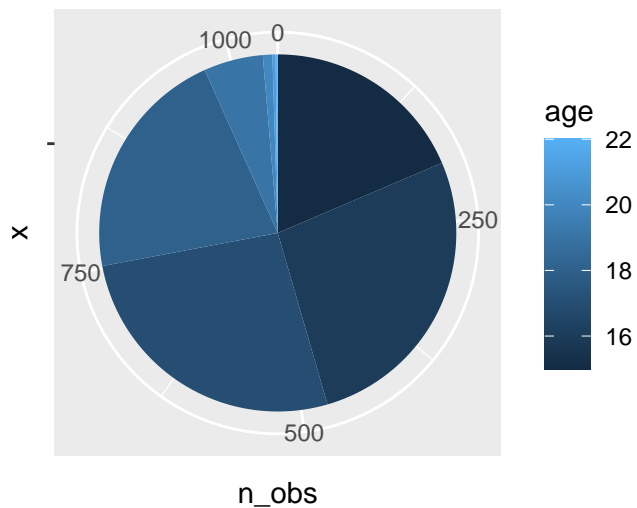
```
## Warning: The dot-dot notation (`..count..`) was deprecated in ggplot2 3.4.0.  
## i Please use `after_stat(count)` instead.
```



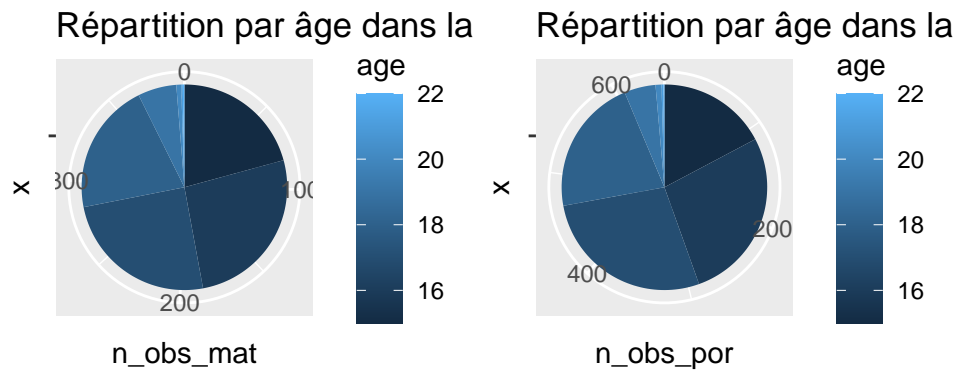
#### A.7.2 L'âge des élèves

On peut par exemple

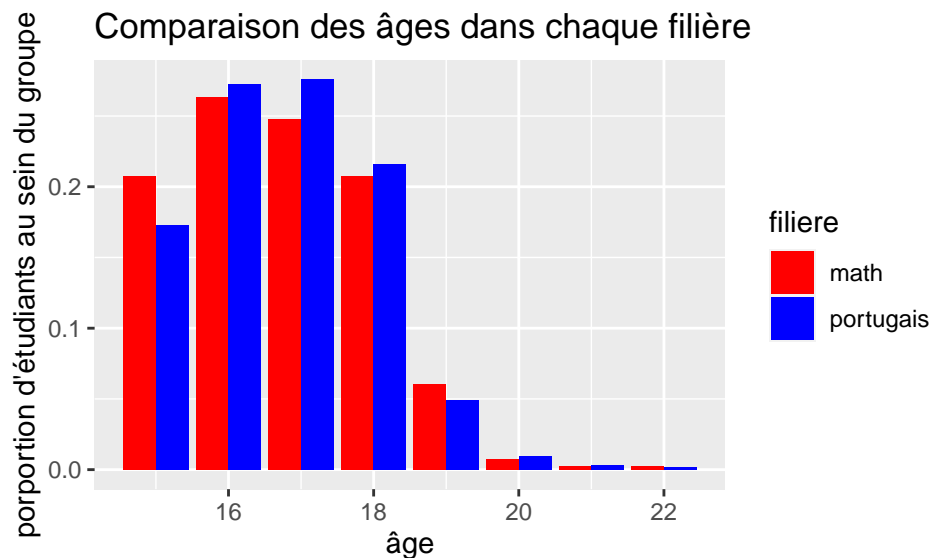
##### Répartition par âge toutes filière confondue



La couleur la plus claire correspond à l'âge le plus grand (22 ans), dès que l'on passe à une couleur plus foncée, on diminue l'âge de 1. On voit clairement ici que la majorité des étudiants ont entre 15 et 19 ans.



On voit que la répartition semble être la grossièrement la même, en effet:



## B Codes pour la partie Machine Learning

Voici un exemple de code pour l'architecture des fonctions d'évaluation:

```
LDA = function(data.train,data.test)
{
  res_lda=lda(data.train$RS ~., data=data.train)
  pred_lda <- predict(res_lda,newdata=data.test)$posterior[,2]

  # Table de confusion
  tab = table(data.test$RS,predict(res_lda,newdata=data.test)$class)

  # Courbe ROC
  ROC_lda <- roc(data.test$RS, pred_lda)

  # Accuracy
  accuracy_lda = mean(data.test$RS==predict(res_lda,newdata=data.test)$class)

  res = list(accuracy_lda,tab,ROC_lda)
  return(res)
}
```

Voici le code pour l'évaluation des modèles :

```
# Suppression des colonnes
X = subset(df, select = -c(G1,G2,G3,Moy) )

for(i in 1:N)
{
  # Génération des jeux d'entraînement et de test
  set.seed(i)
  n <- nrow(X)
  p <- ncol(X)-1
  test.ratio <- .2 # ratio of test/train samples
  n.test <- round(n*test.ratio)
  tr <- sample(1:n,n.test)
  df.test <- X[tr,]
  df.train <- X[-tr,]

  # LDA
  res = LDA(df.train,df.test)
  a_lda[i] = res[[1]]
  ROC_lda = res[[3]]

  # QDA
  res = QDA(df.train,df.test)
  a_qda[i] = res[[1]]
  ROC_qda = res[[3]]

  # Stepwise
  res = stepwise(df.train,df.test)
  a_step = res[[1]]
  ROC_step = res[[3]]

  # cart
  res = Cart(df.train,df.test)
  a_cart[i] = res[[1]]
  ROC_cart = res[[3]]

  # Random forest
  res = RF(df.train,df.test)
  a_RF[i] = res[[1]]
  ROC_RF = res[[3]]

  # Regression logistique
  res = logit(df.train,df.test)
  a_logit[i] = res[[1]]
  ROC_logit = res[[3]]

  # Regression logistique
  res = SVM(df.train,df.test)
  a_svm[i] = res[[1]]
  ROC_svm = res[[3]]
}
```