

# Projet Analyse de données

Rudio et Léo-Paul

2023-05-09

## Présentation du projet et du jeu de données

Le jeu de données est constitué d'informations sur la vie d'étudiants dans une université du Portugal. Ces informations vont de leur résultats universitaires, leur vie familiale à leur consommation d'alcool. Le jeu a été construit à partir d'une enquête menée auprès d'étudiant 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 celle sur la consommation d'alcool proposée directement par le nom du jeu de données.

Voici les variables présentes dans ce jeu de données ;

- **school** - student's school (binary: 'GP' - Gabriel Pereira or 'MS' - Mousinho da Silveira)
- **sex** - student's sex (binary: 'F' - female or 'M' - male)
- **age** - student's age (numeric: from 15 to 22)
- **address** - student's home address type (binary: 'U' - urban or 'R' - rural)
- **famsize** - family size (binary: 'LE3' - less or equal to 3 or 'GT3' - greater than 3)
- **Pstatus** - parent's cohabitation status (binary: 'T' - living together or 'A' - apart)
- **Medu** - mother's education (numeric: 0 - none, 1 - primary education (4th grade), 2 - 5th to 9th grade, 3 - secondary education or 4 - higher education)
- **Fedu** - father's education (numeric: 0 - none, 1 - primary education (4th grade), 2 - 5th to 9th grade, 3 - secondary education or 4 - higher education)
- **Mjob** - mother's job (nominal: 'teacher', 'health' care related, civil 'services' (e.g. administrative or police), 'at\_home' or 'other')
- **Fjob** - father's job (nominal: 'teacher', 'health' care related, civil 'services' (e.g. administrative or police), 'at\_home' or 'other')
- **reason** - reason to choose this school (nominal: close to 'home', school 'reputation', 'course' preference or 'other')
- **guardian** - student's guardian (nominal: 'mother', 'father' or 'other')
- **traveltime** - home to school travel time (numeric: 1 - <15 min., 2 - 15 to 30 min., 3 - 30 min. to 1 hour, or 4 - >1 hour)
- **studytime** - weekly study time (numeric: 1 - <2 hours, 2 - 2 to 5 hours, 3 - 5 to 10 hours, or 4 - >10 hours)
- **failures** - number of past class failures (numeric: n if  $1 \leq n \leq 3$ , else 4)
- **schoolsup** - extra educational support (binary: yes or no)
- **famsup** - family educational support (binary: yes or no)
- **paid** - extra paid classes within the course subject (Math or Portuguese) (binary: yes or no)
- **activities** - extra-curricular activities (binary: yes or no)
- **nursery** - attended nursery school (binary: yes or no)
- **higher** - wants to take higher education (binary: yes or no)
- **internet** - Internet access at home (binary: yes or no)
- **romantic** - with a romantic relationship (binary: yes or no)
- **famrel** - quality of family relationships (numeric: from 1 - very bad to 5 - excellent)

- **freetime** - free time after school (numeric: from 1 - very low to 5 - very high)
- **goout** - going out with friends (numeric: from 1 - very low to 5 - very high)
- **Dalc** - workday alcohol consumption (numeric: from 1 - very low to 5 - very high)
- **Walc** - weekend alcohol consumption (numeric: from 1 - very low to 5 - very high)
- **health** - current health status (numeric: from 1 - very bad to 5 - very good)
- **absences** - number of school absences (numeric: from 0 to 93)

These grades are related with the course subject, Math or Portuguese: - **G1** - first period grade (numeric: from 0 to 20) - **G2** - second period grade (numeric: from 0 to 20) - **G3** - final grade (numeric: from 0 to 20, output target)

Au cours de ce projet, nous nous concentrons sur la variable G3 qui est la variable de sortie représentant la note finale des élèves. Il s'agirait donc d'un problème de régression sur la variables G3 ou même plus généralement un problème de classification.

Voici les étapes que nous allons suivre :

1. Identifier les variables significatives
2. Appliquer des méthodes de classification sur la réussite scolaire
3. Effectuer une regression linéaires pour prédire G3
4. Comparer des méthodes de machine learning pour prédire G3

## 1.Chargement des données

```
##  school sex age address famsize Pstatus Medu Fedu Mjob Fjob reason
## 1    GP   F  18      U    GT3      A    4    4  at_home teacher  course
## 2    GP   F  17      U    GT3      T    1    1  at_home  other  course
## 3    GP   F  15      U    LE3      T    1    1  at_home  other  other
## 4    GP   F  15      U    GT3      T    4    2  health services  home
## 5    GP   F  16      U    GT3      T    3    3  other  other  home
## 6    GP   M  16      U    LE3      T    4    3  services  other reputation
##  guardian traveltime studytime failures schoolsup famsup paid activities
## 1  mother          2          2          0        yes    no    no          no
## 2  father          1          2          0        no    yes    no          no
## 3  mother          1          2          3        yes    no    yes          no
## 4  mother          1          3          0        no    yes    yes          yes
## 5  father          1          2          0        no    yes    yes          no
## 6  mother          1          2          0        no    yes    yes          yes
##  nursery higher internet romantic famrel freetime goout Dalc Walc health
## 1    yes    yes      no      no      4      3    4    1    1    3
## 2    no    yes      yes      no      5      3    3    1    1    3
## 3    yes    yes      yes      no      4      3    2    2    3    3
## 4    yes    yes      yes     yes      3      2    2    1    1    5
## 5    yes    yes      no      no      4      3    2    1    2    5
## 6    yes    yes      yes      no      5      4    2    1    2    5
##  absences G1 G2 G3
## 1         6  5  6  6
## 2         4  5  5  6
## 3        10  7  8 10
## 4         2 15 14 15
## 5         4  6 10 10
## 6        10 15 15 15
```

## 2. Nettoyage et vérification des données

Le jeu est composé de 33 variables dont 17 qualitatives et 16 quantitatives. On calcule la moyenne pour chaque élève, et on rajoute une variable pour la réussite scolaire.

```
## 'data.frame': 1044 obs. of 33 variables:
## $ school : Factor w/ 2 levels "GP","MS": 1 1 1 1 1 1 1 1 1 1 ...
## $ sex : Factor w/ 2 levels "F","M": 1 1 1 1 1 2 2 1 2 2 ...
## $ age : int 18 17 15 15 16 16 16 17 15 15 ...
## $ address : Factor w/ 2 levels "R","U": 2 2 2 2 2 2 2 2 2 2 ...
## $ famsize : Factor w/ 2 levels "GT3","LE3": 1 1 2 1 1 2 2 1 2 1 ...
## $ Pstatus : Factor w/ 2 levels "A","T": 1 2 2 2 2 2 2 1 1 2 ...
## $ Medu : int 4 1 1 4 3 4 2 4 3 3 ...
## $ Fedu : int 4 1 1 2 3 3 2 4 2 4 ...
## $ Mjob : Factor w/ 5 levels "at_home","health",...: 1 1 1 2 3 4 3 3 4 3 ...
## $ Fjob : Factor w/ 5 levels "at_home","health",...: 5 3 3 4 3 3 3 5 3 3 ...
## $ reason : Factor w/ 4 levels "course","home",...: 1 1 3 2 2 4 2 2 2 2 ...
## $ guardian : Factor w/ 3 levels "father","mother",...: 2 1 2 2 1 2 2 2 2 2 ...
## $ traveltime: int 2 1 1 1 1 1 1 2 1 1 ...
## $ studytime : int 2 2 2 3 2 2 2 2 2 2 ...
## $ failures : int 0 0 3 0 0 0 0 0 0 0 ...
## $ schoolsup : Factor w/ 2 levels "no","yes": 2 1 2 1 1 1 1 2 1 1 ...
## $ famsup : Factor w/ 2 levels "no","yes": 1 2 1 2 2 2 1 2 2 2 ...
## $ paid : Factor w/ 2 levels "no","yes": 1 1 2 2 2 2 1 1 2 2 ...
## $ activities: Factor w/ 2 levels "no","yes": 1 1 1 2 1 2 1 1 1 2 ...
## $ nursery : Factor w/ 2 levels "no","yes": 2 1 2 2 2 2 2 2 2 2 ...
## $ higher : Factor w/ 2 levels "no","yes": 2 2 2 2 2 2 2 2 2 2 ...
## $ internet : Factor w/ 2 levels "no","yes": 1 2 2 2 1 2 2 1 2 2 ...
## $ romantic : Factor w/ 2 levels "no","yes": 1 1 1 2 1 1 1 1 1 1 ...
## $ famrel : int 4 5 4 3 4 5 4 4 4 5 ...
## $ freetime : int 3 3 3 2 3 4 4 1 2 5 ...
## $ goout : int 4 3 2 2 2 2 4 4 2 1 ...
## $ Dalc : int 1 1 2 1 1 1 1 1 1 1 ...
## $ Walc : int 1 1 3 1 2 2 1 1 1 1 ...
## $ health : int 3 3 3 5 5 5 3 1 1 5 ...
## $ absences : int 6 4 10 2 4 10 0 6 0 0 ...
## $ G1 : int 5 5 7 15 6 15 12 6 16 14 ...
## $ G2 : int 6 5 8 14 10 15 12 5 18 15 ...
## $ G3 : int 6 6 10 15 10 15 11 6 19 15 ...
## NULL

## [1] 1044
```

	school	sex	age	address	famsize	Pstatus	Medu	Fedu	Mjob	Fjob	reason
## 1	GP	F	18	U	GT3	A	4	4	at_home	teacher	course
## 2	GP	F	17	U	GT3	T	1	1	at_home	other	course
## 3	GP	F	15	U	LE3	T	1	1	at_home	other	other
## 4	GP	F	15	U	GT3	T	4	2	health	services	home
## 5	GP	F	16	U	GT3	T	3	3	other	other	home
## 6	GP	M	16	U	LE3	T	4	3	services	other	reputation

	guardian	traveltime	studytime	failures	schoolsup	famsup	paid	activities
## 1	mother	2	2	0	yes	no	no	no
## 2	father	1	2	0	no	yes	no	no
## 3	mother	1	2	3	yes	no	yes	no
## 4	mother	1	3	0	no	yes	yes	yes
## 5	father	1	2	0	no	yes	yes	no

```
## 6   mother      1      2      0      no   yes   yes      yes
##   nursery higher internet romantic famrel freetime goout Dalc Walc health
## 1    yes    yes      no      no      4      3      4      1      1      3
## 2    no     yes      yes      no      5      3      3      1      1      3
## 3    yes    yes      yes      no      4      3      2      2      3      3
## 4    yes    yes      yes      yes     3      2      2      1      1      5
## 5    yes    yes      no      no      4      3      2      1      2      5
## 6    yes    yes      yes      no      5      4      2      1      2      5
##   absences G1 G2 G3      Moy      RS
## 1         6 5 6 6  5.666667  exclusion
## 2         4 5 5 6  5.333333  exclusion
## 3        10 7 8 10  8.333333  exclusion
## 4         2 15 14 15 14.666667   admis
## 5         4 6 10 10  8.666667 redoublement
## 6        10 15 15 15 15.000000   admis
```

### 3. Exploration des données : études des variables

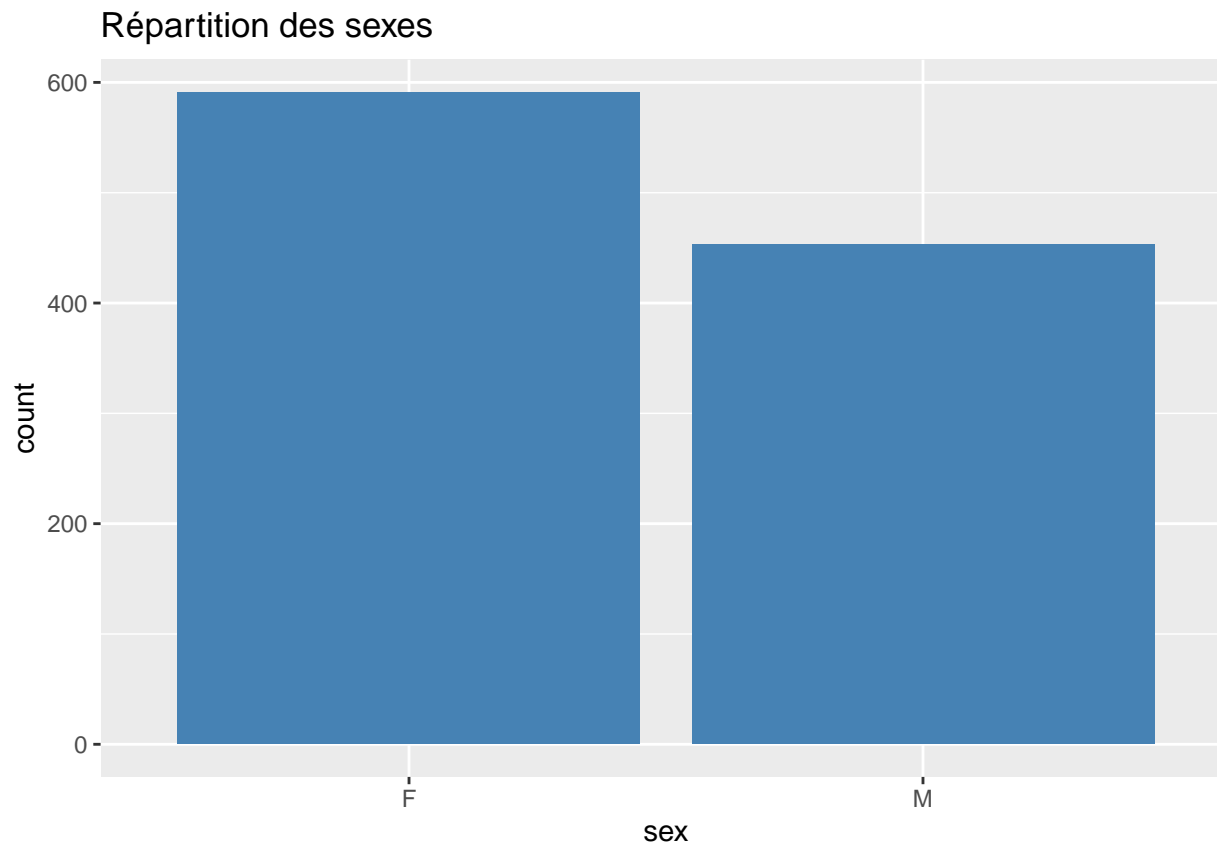
Cette partie consiste à appliquer des méthodes de statistiques descriptives afin de mieux comprendre le jeu de données. On se concentre sur l'analyse de la distribution des variables et leur corrélation avec les résultats scolaires.

#### Les variables qualitatives

##### 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,  $p\text{-value} > 5\%$  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.

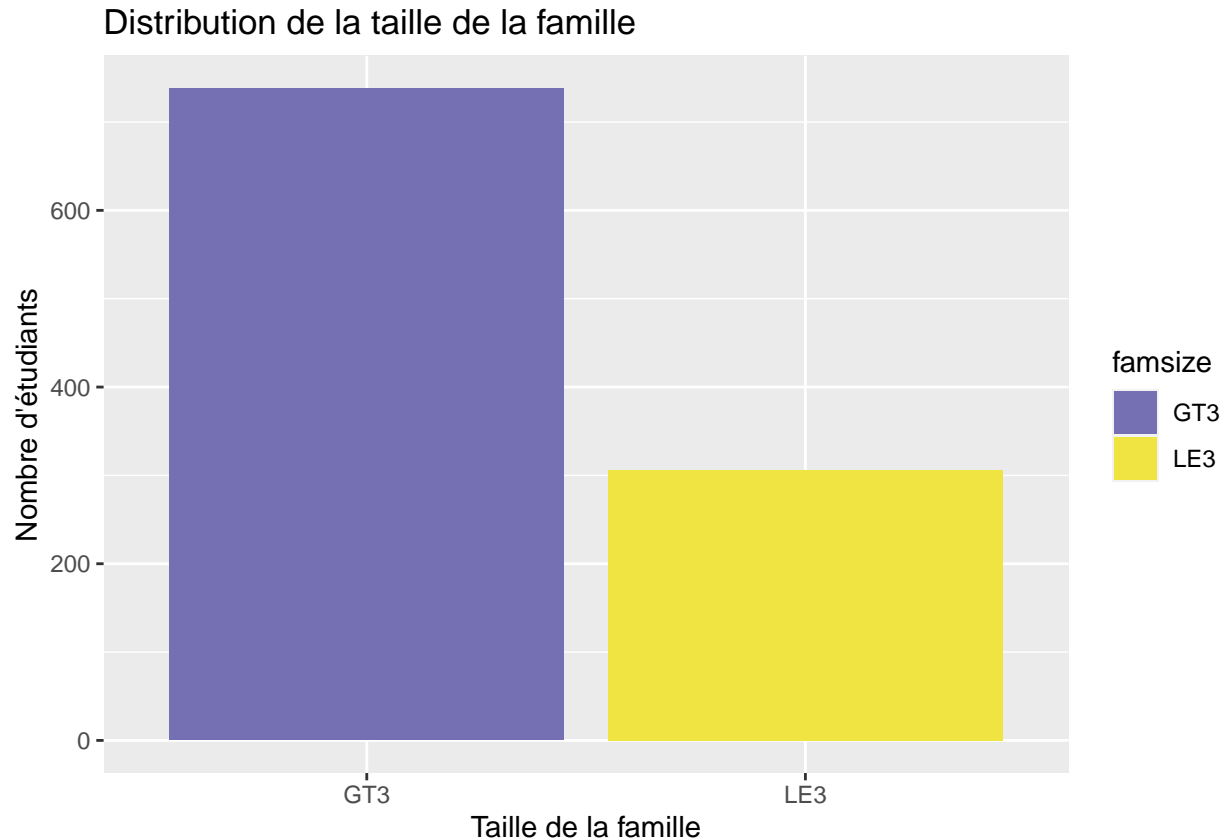
```
##
## Attachement du package : 'ggplot2'
## L'objet suivant est masqué depuis 'package:randomForest':
##
##   margin
```



```
##
## Call:
## lm(formula = Moy ~ sex, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.0152  -2.0152  -0.0152   2.1722   8.1722
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  11.3486     0.1324  85.706  <2e-16 ***
## sexM         -0.1874     0.2010  -0.932   0.351
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.219 on 1042 degrees of freedom
## Multiple R-squared:  0.0008335, Adjusted R-squared:  -0.0001254
## F-statistic: 0.8693 on 1 and 1042 DF, p-value: 0.3514
##
## Pearson's Chi-squared test
##
## data:  df$sex and df$RS
## X-squared = 1.1035, df = 2, p-value = 0.5759
```

## 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.



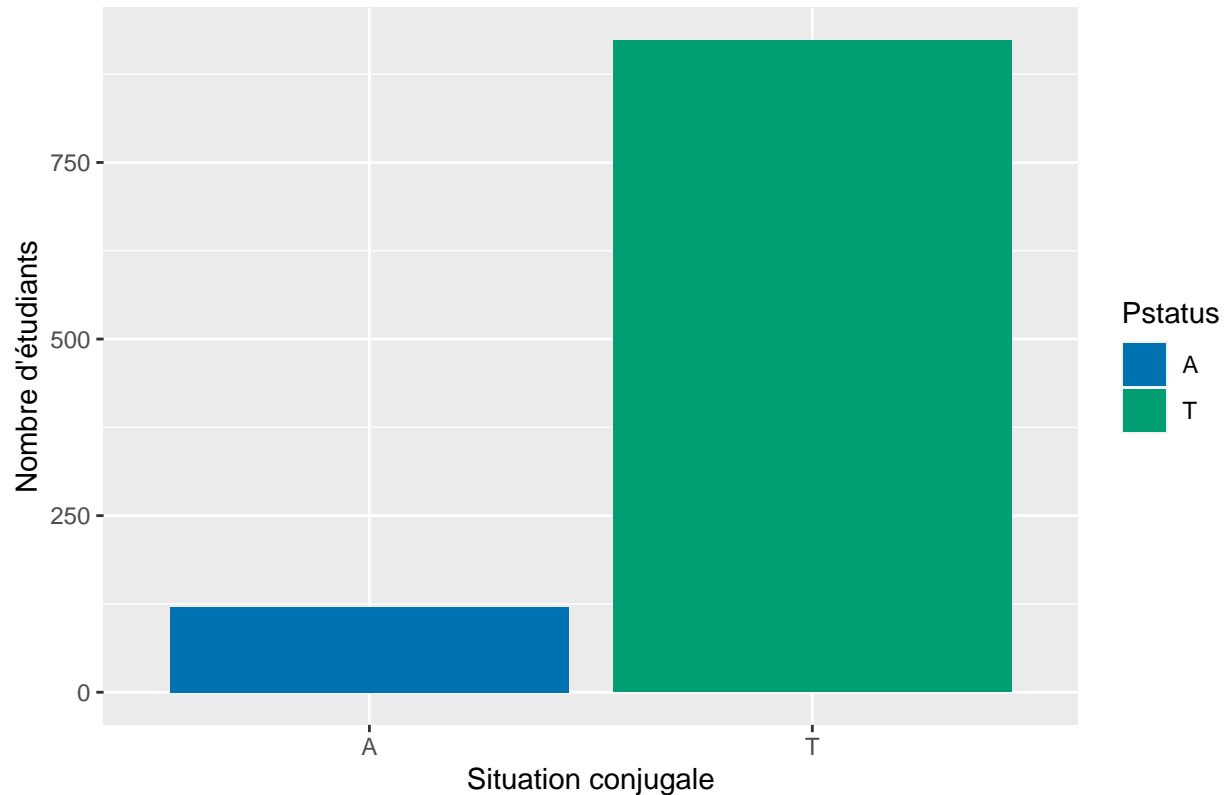
```
##
## Call:
## lm(formula = Moy ~ famsize, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.9096 -1.9096 -0.1391  2.1942  8.1942
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  11.1391     0.1183   94.15  <2e-16 ***
## famsizeLE3    0.4371     0.2185    2.00  0.0457 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.214 on 1042 degrees of freedom
## Multiple R-squared:  0.003825, Adjusted R-squared:  0.002869
## F-statistic: 4.001 on 1 and 1042 DF, p-value: 0.04573
##
## Pearson's Chi-squared test
```

```
##
## data:  df$famsize and df$RS
## X-squared = 4.5986, df = 2, p-value = 0.1003
```

### 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.

Distribution de la situation conjugale des parents



```
##
## Call:
## lm(formula = Moy ~ Pstatus, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.0744  -1.9155   0.0845   2.0845   8.0845
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  11.4077     0.2927   38.97  <2e-16 ***
## PstatusT     -0.1589     0.3113   -0.51    0.61
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.22 on 1042 degrees of freedom
## Multiple R-squared:  0.0002499, Adjusted R-squared: -0.0007095
```

```
## F-statistic: 0.2605 on 1 and 1042 DF, p-value: 0.6099
```

```
##
```

```
## Pearson's Chi-squared test
```

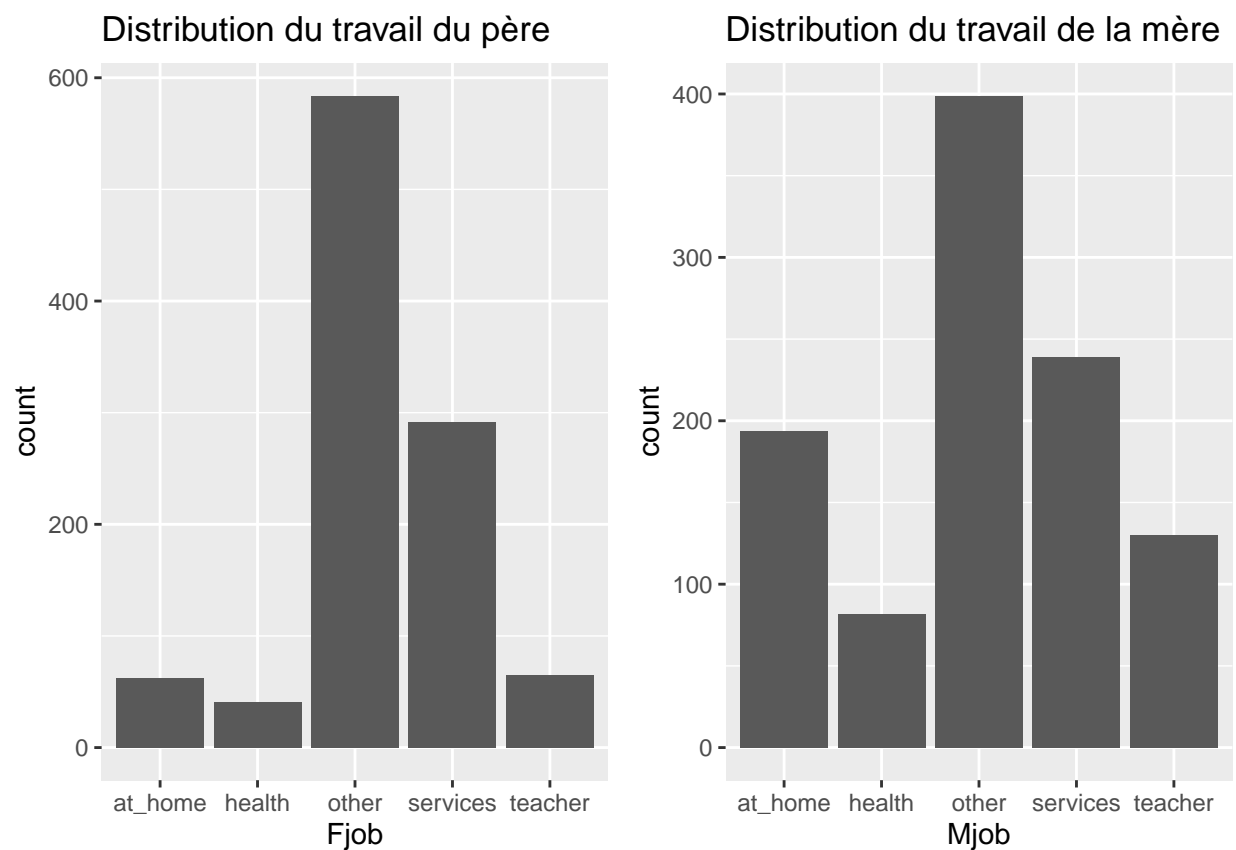
```
##
```

```
## data: df$Pstatus and df$RS
```

```
## X-squared = 0.9565, df = 2, p-value = 0.6199
```

## 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.



```
##
```

```
## Call:
```

```
## lm(formula = Moy ~ Medu + Fedu, data = df)
```

```
##
```

```
## Residuals:
```

```
##      Min       1Q   Median       3Q      Max
## -10.265  -1.732   0.068   2.126   7.852
```

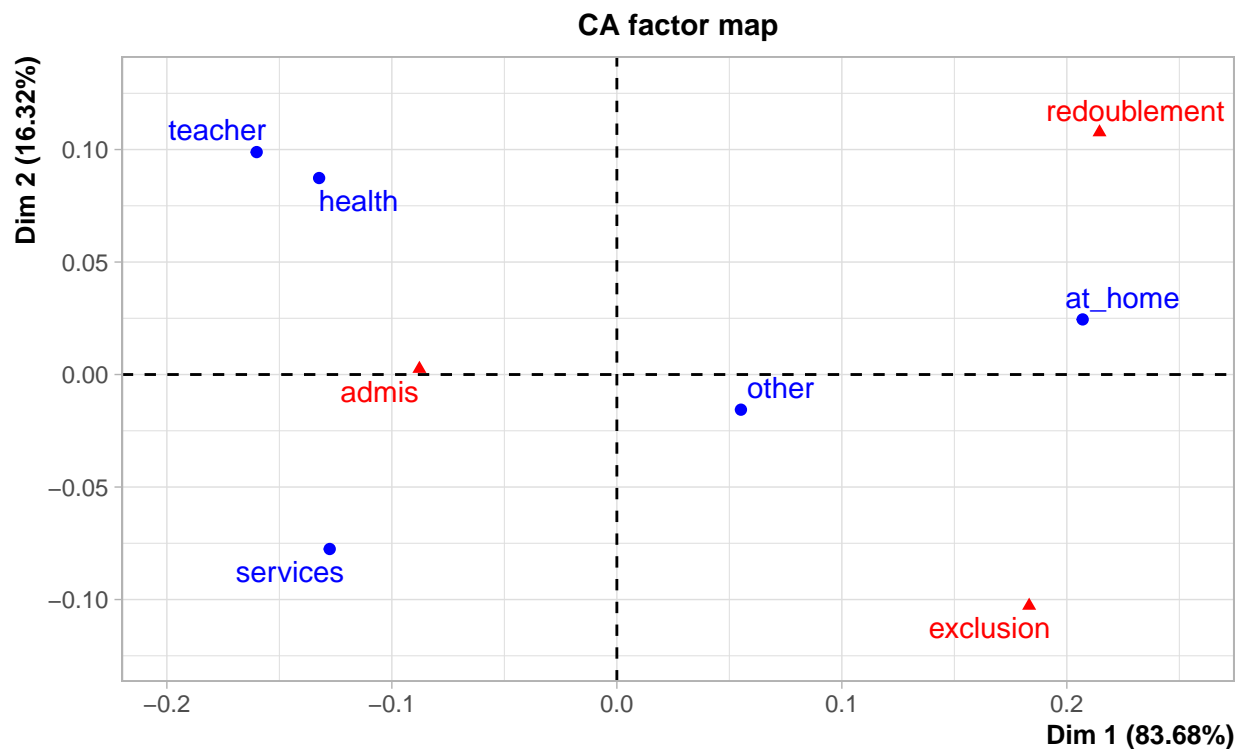
```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    9.4233     0.2595  36.316 < 2e-16 ***
```



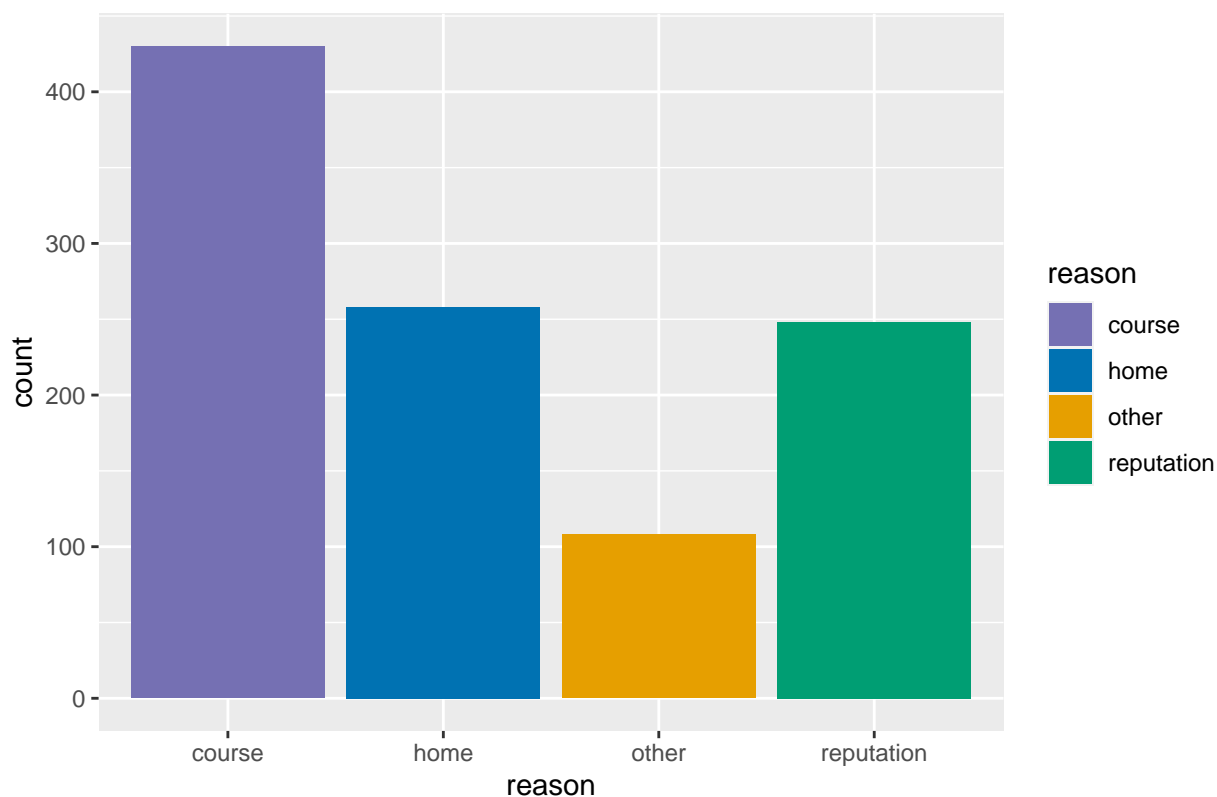
```
## Medu          0.5214      0.1125    4.635 4.02e-06 ***
## Fedu          0.2037      0.1150    1.771  0.0769 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.133 on 1041 degrees of freedom
## Multiple R-squared:  0.05434,    Adjusted R-squared:  0.05252
## F-statistic: 29.91 on 2 and 1041 DF,  p-value: 2.344e-13
##
## Pearson's Chi-squared test
##
## data:  df$Mjob and df$RS
## X-squared = 21.736, df = 8, p-value = 0.005429
##
## Pearson's Chi-squared test
##
## data:  df$Fjob and df$RS
## X-squared = 7.4964, df = 8, p-value = 0.4841
```



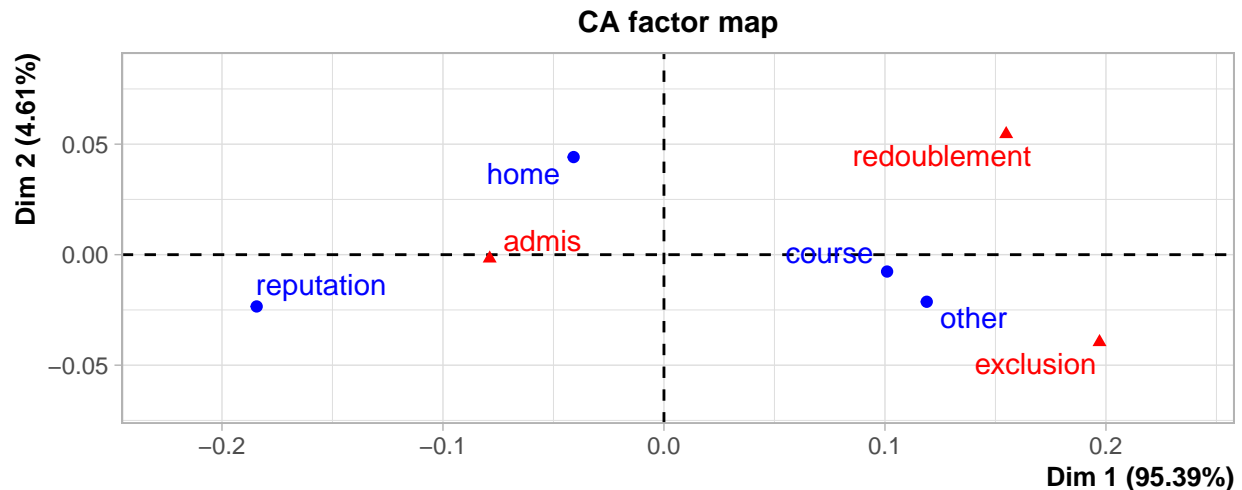
### Les raisons du choix d'école

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 ( $p\text{-value} < 5\%$ ). 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 du travail du père



```
##
## 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
##
##
## Pearson's Chi-squared test
##
## data:  df$reason and df$RS
## X-squared = 15.479, df = 6, p-value = 0.01684
```



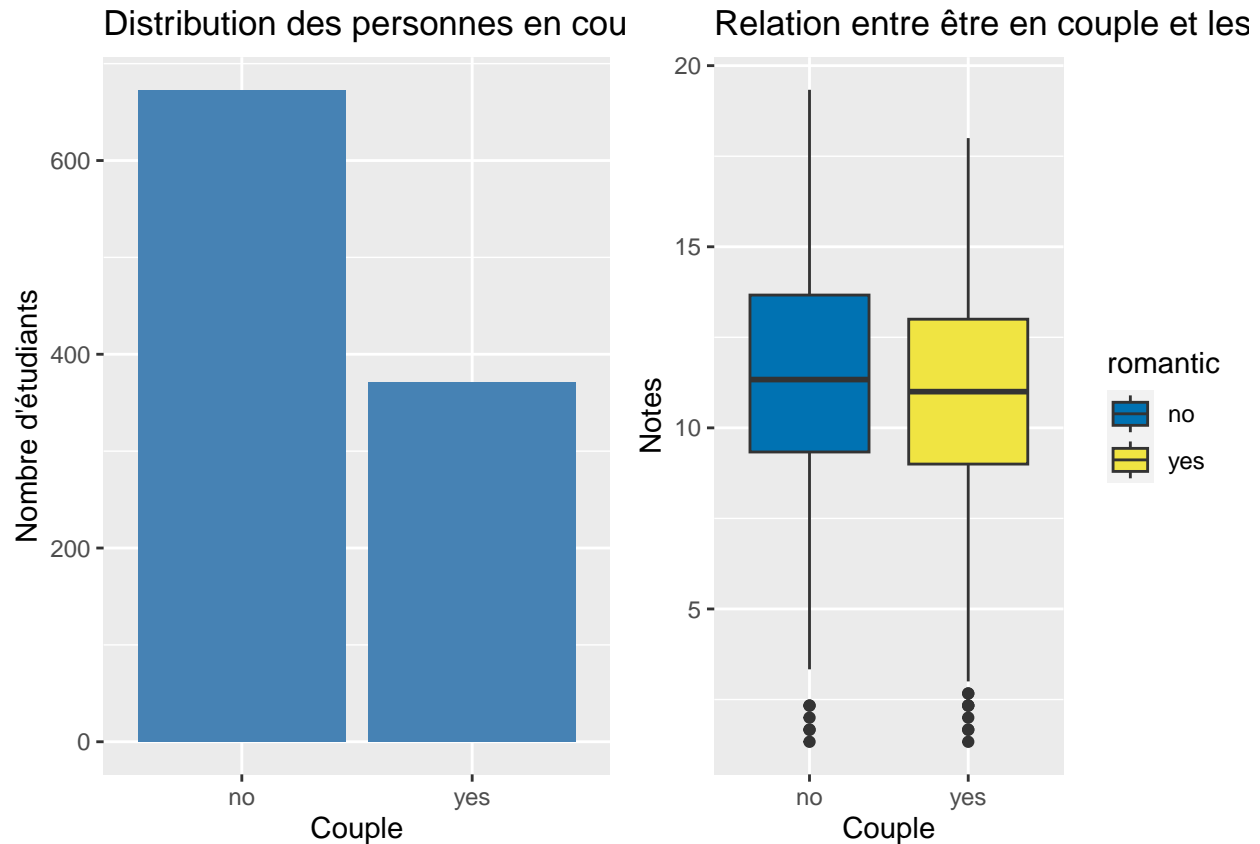
On voit bien avec l'AFC que les personnes étant admises sont celles qui choisissent l'école pour sa réputation et sa proximité par rapport à leur domicile. A l'inverse on voit que les étudiants qui ont échoués sont ceux qui ont choisis l'école pour les cours ou d'autres raisons. On voit ici une des limite 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 choisies l'école pour les cours plus que pour sa réputation.

### 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.

```
##
## Call:
## lm(formula = Moy ~ romantic, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.1486  -1.9455   0.1222   2.1847   7.8514
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   11.4819     0.1236  92.871  < 2e-16 ***
## romanticyes   -0.6041     0.2074  -2.913  0.00366 **
```

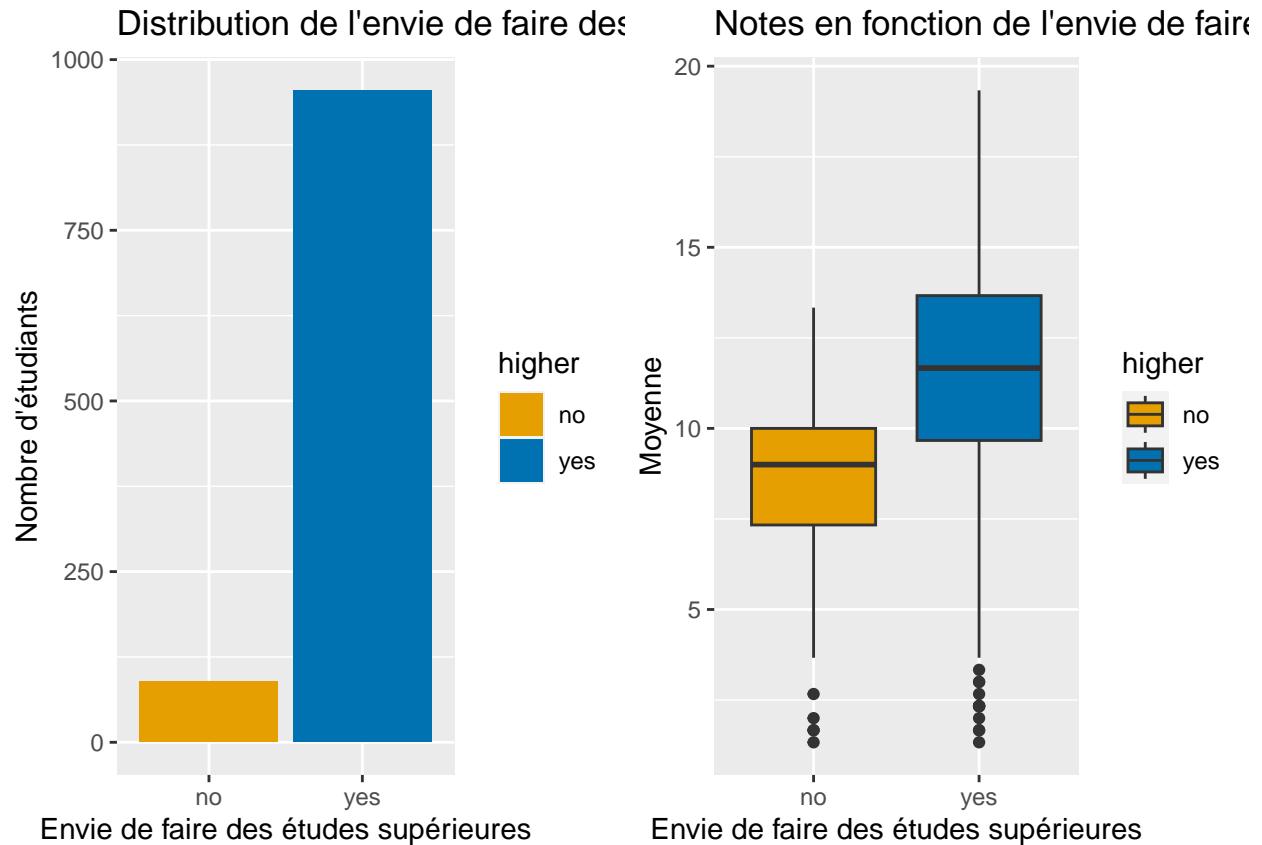
```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.207 on 1042 degrees of freedom
## Multiple R-squared:  0.008077,    Adjusted R-squared:  0.007125
## F-statistic: 8.485 on 1 and 1042 DF,  p-value: 0.003658
```



```
##
## Pearson's Chi-squared test
##
## data: df$romantic and df$RS
## X-squared = 5.5477, df = 2, p-value = 0.06242
```

### 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.

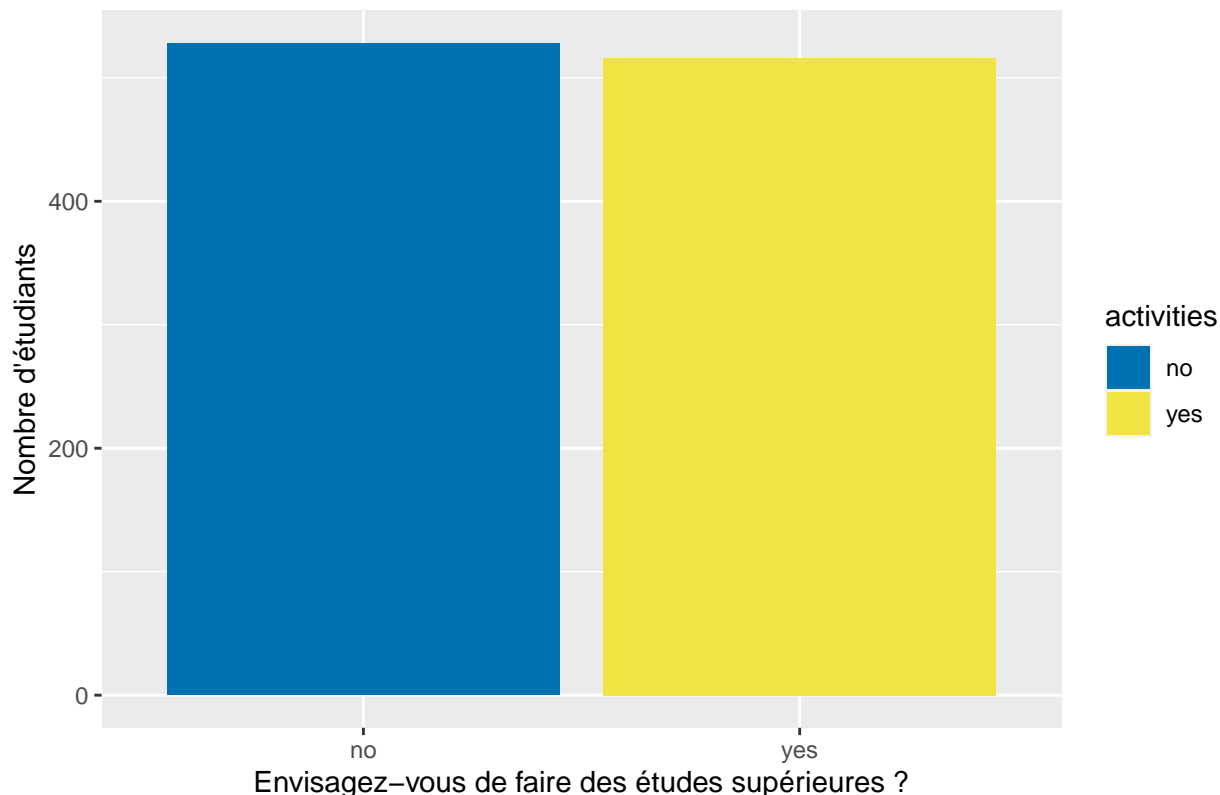


```
##
## Call:
## lm(formula = Moy ~ higher, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.1930  -1.8597   0.1403   2.1403   7.8070
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   8.4869     0.3293  25.775  <2e-16 ***
## higheryes     3.0395     0.3443   8.829  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.106 on 1042 degrees of freedom
## Multiple R-squared:  0.0696, Adjusted R-squared:  0.06871
## F-statistic: 77.95 on 1 and 1042 DF,  p-value: < 2.2e-16
##
## Pearson's Chi-squared test
##
## data:  df$higher and df$RS
## X-squared = 66.594, df = 2, p-value = 3.461e-15
```

## 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.

Distribution de la pratique d'activités extrascolaires

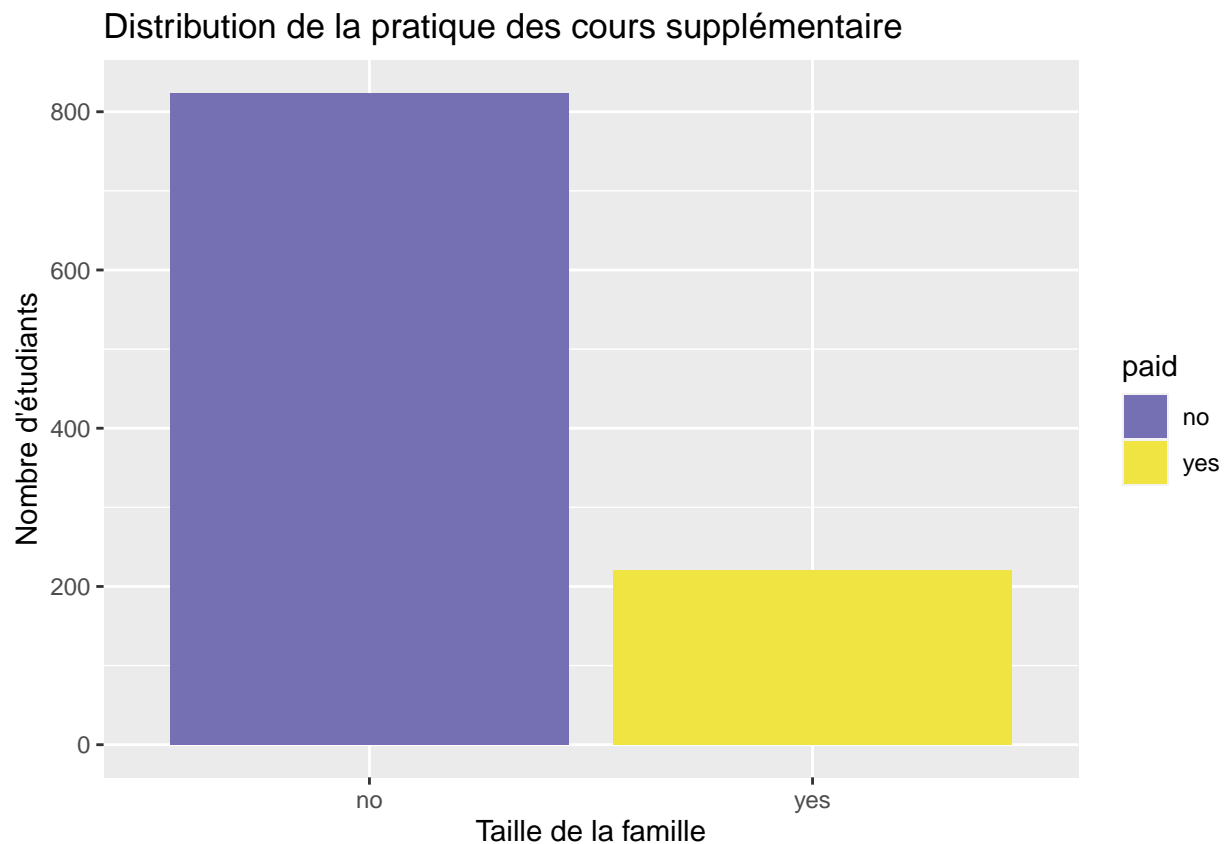


```
##
## Call:
## lm(formula = Moy ~ activities, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.1085  -2.0966  -0.0966   2.2248   7.8915
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    11.0966     0.1399  79.292  <2e-16 ***
## activitiesyes     0.3453     0.1991   1.734   0.0831 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.216 on 1042 degrees of freedom
## Multiple R-squared:  0.002879, Adjusted R-squared:  0.001922
## F-statistic: 3.008 on 1 and 1042 DF, p-value: 0.08313
```

```
##
## Pearson's Chi-squared test
##
## data: df$activities and df$RS
## X-squared = 2.5236, df = 2, p-value = 0.2831
```

### 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.



```
##
## Call:
## lm(formula = Moy ~ paid, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.9951 -1.9951  0.0049  2.0049  8.0049
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   11.3285     0.1121  101.05  <2e-16 ***
## paidyes       -0.2906     0.2442   -1.19   0.234
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## Residual standard error: 3.218 on 1042 degrees of freedom
## Multiple R-squared:  0.001357,    Adjusted R-squared:  0.0003986
## F-statistic: 1.416 on 1 and 1042 DF,  p-value: 0.2344

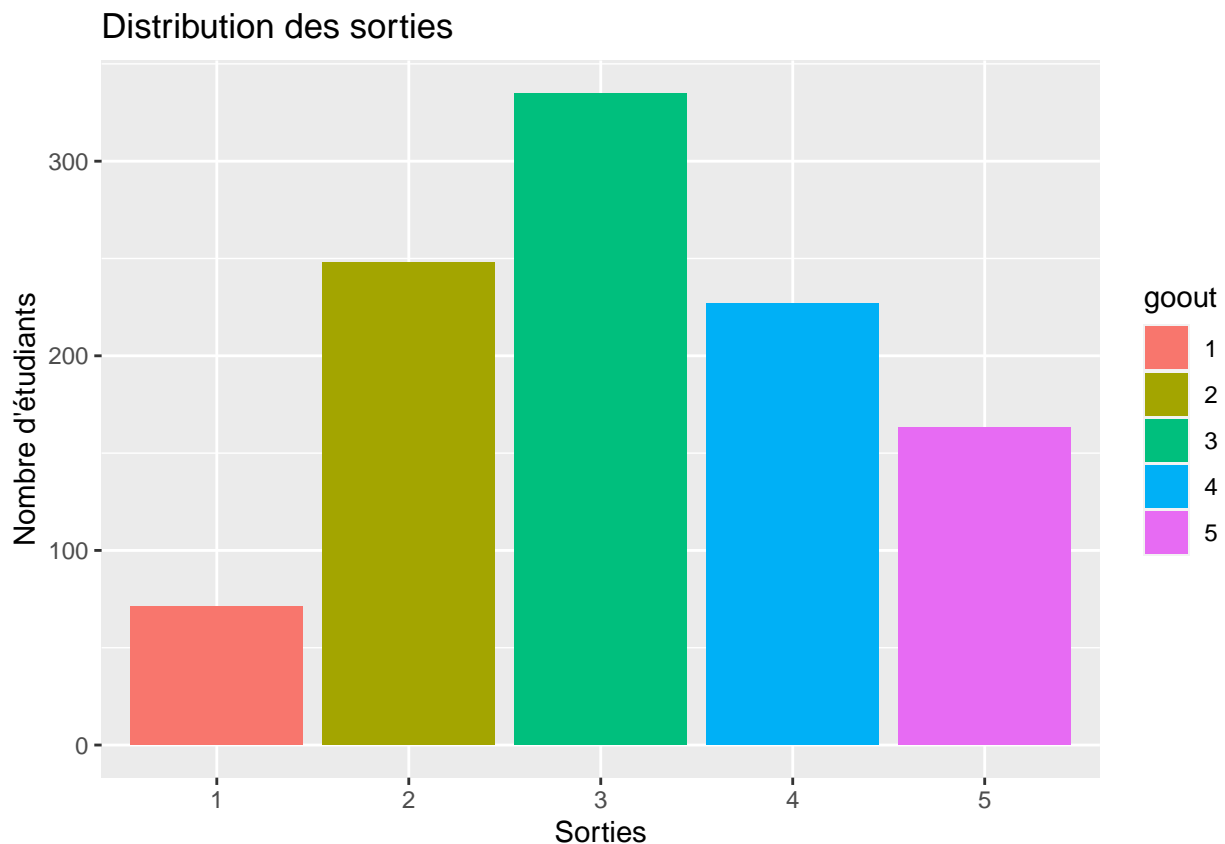
##
## Pearson's Chi-squared test
##
## data:  df$paid and df$RS
## X-squared = 4.8571, df = 2, p-value = 0.08816
```

## Les variables qualitatives à modalités numériques

### 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 aux 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.

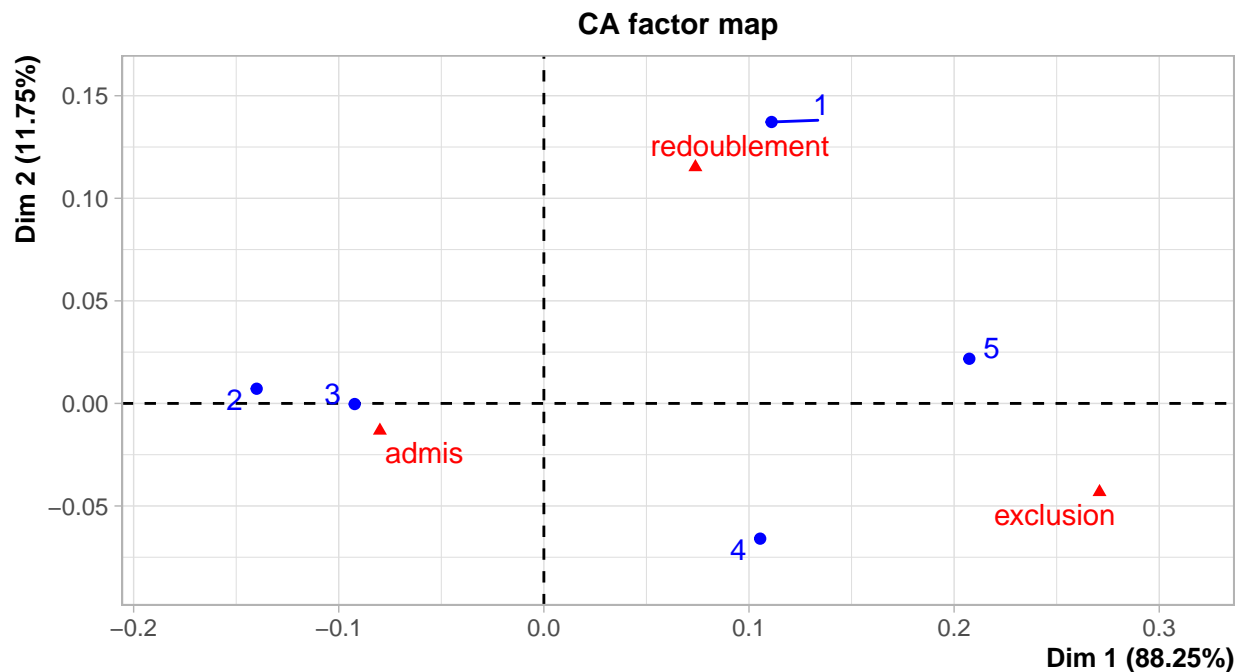
Etant donné, la corrélation entre RS et goout, on peut effectuer une AFC pour préciser. On peut remarquer que ceux qui sortent peu-moyennement auront tendance à être admis alors que ce qui ne sortent pas (retrait/exclusion sociale) vont plutôt redoubler et les autres vont avoir tendances à se faire exclure.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



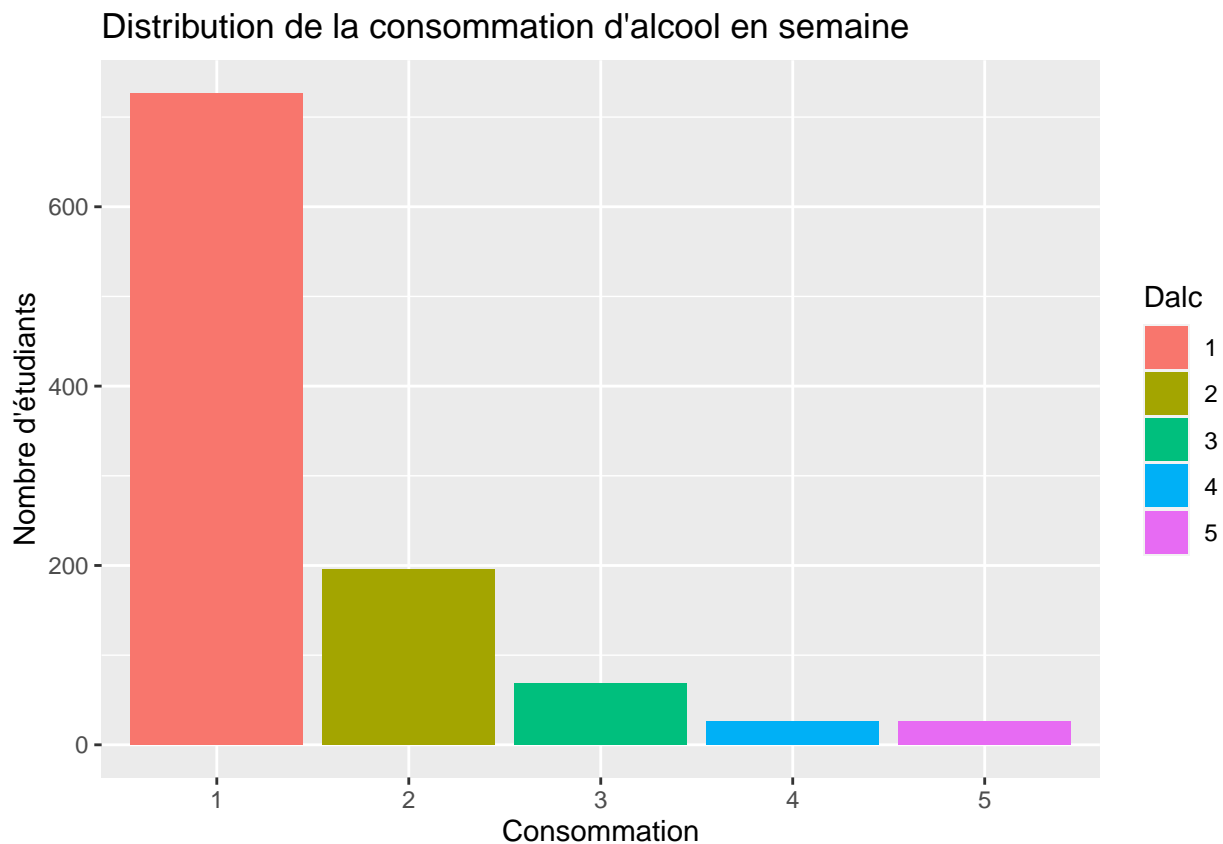
Avec un test du  $\chi^2$  on observe que les variables goout et RS sont corrélées (p-valeur petite devant 5%).

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.

### La consommation d'alcool

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

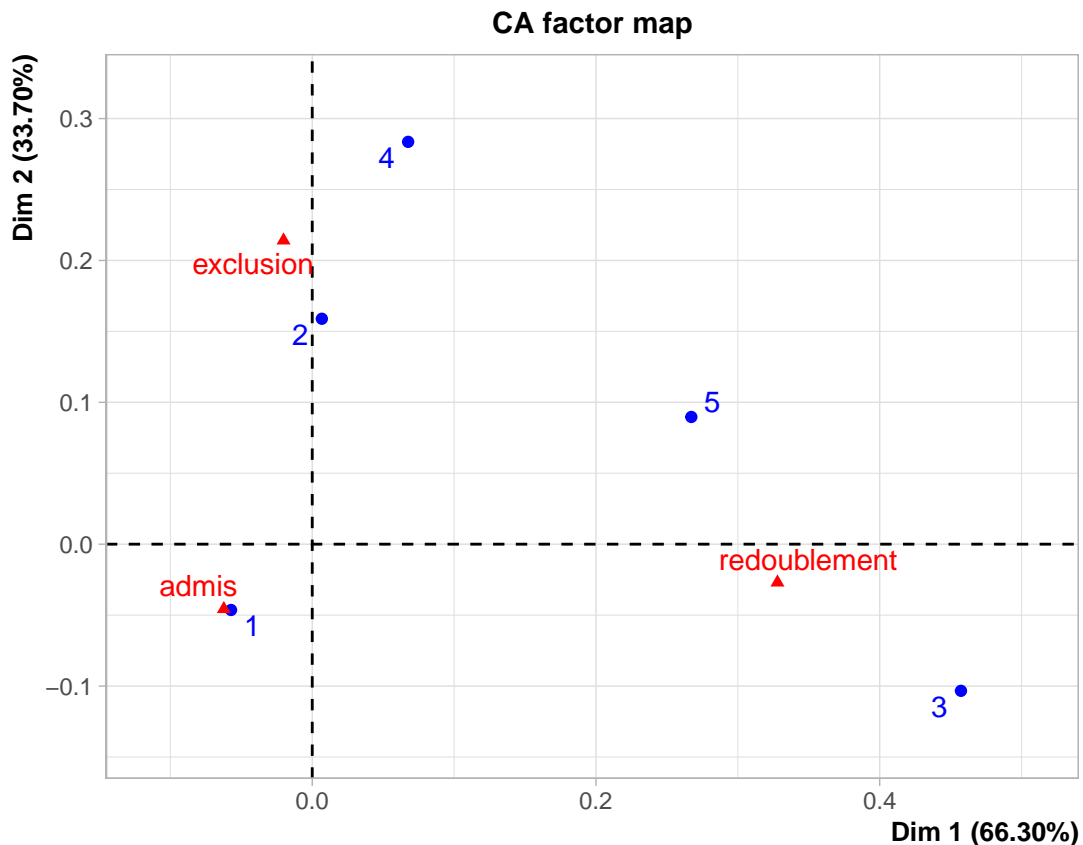


```
##
## Call:
## lm(formula = Moy ~ Dalc, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.244  -1.911   0.089   2.089   7.756
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  11.5777     0.1181  98.001  < 2e-16 ***
## Dalc2        -0.8889     0.2564  -3.467  0.000547 ***
## Dalc3        -0.8917     0.4013  -2.222  0.026475 *
```

```
## Dalc4      -2.0008      0.6358  -3.147 0.001696 **
## Dalc5      -1.3982      0.6358  -2.199 0.028079 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.185 on 1039 degrees of freedom
## Multiple R-squared:  0.02443,    Adjusted R-squared:  0.02068
## F-statistic: 6.505 on 4 and 1039 DF,  p-value: 3.594e-05

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

##
## Pearson's Chi-squared test
##
## data:  df$Dalc and df$RS
## X-squared = 28.342, df = 8, p-value = 0.0004134
```

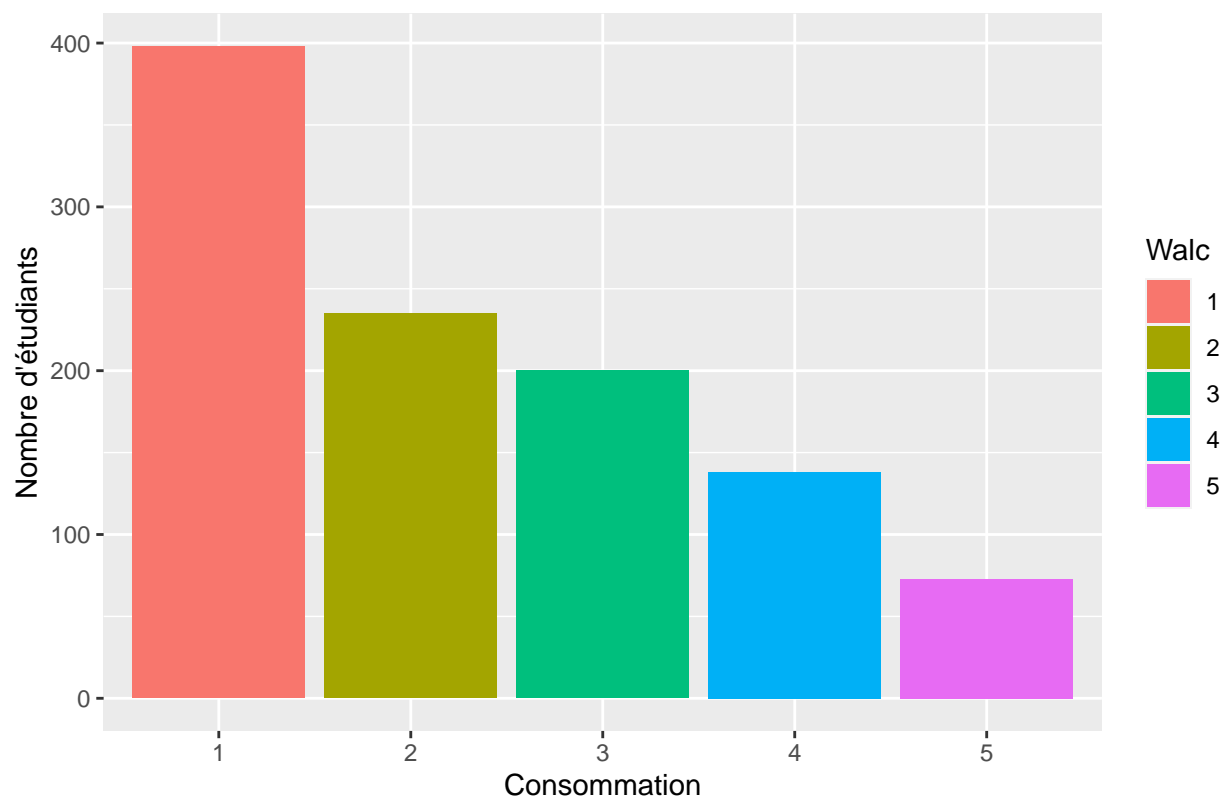


Avec le test du  $\chi^2$  on voit que les variables Dalc et RS sont corrélées. On va donc réaliser une AFC dessus. De même avec la p-valeur du test de Fisher sur les variables Moy et Dalc, on voit que ces variables sont aussi corrélées (et c'est logique au vu du test du  $\chi^2$ ).

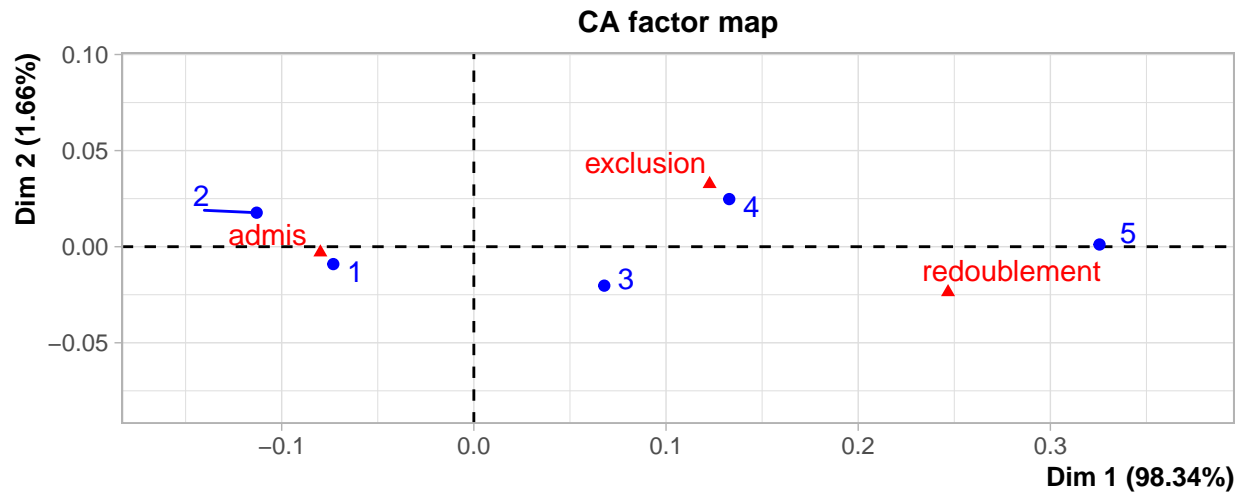
On voit très clairement avec l'AFC que les étudiants qui consomment le plus d'alcool sont ceux qui réussissent le moins. En effet, une forte consommation d'alcool témoigne d'un grand nombre de sortie ou bien d'un grave problème de santé (alcoolisme). Ceux qui réussissent le plus sont ceux qui consomment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

Distribution de la consommation d'alcool le week-end



```
##
## Call:
## lm(formula = Moy ~ Walc, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.2831  -1.9050   0.0503   2.0614   7.7169
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  11.6164     0.1600  72.593 < 2e-16 ***
## Walc2        -0.1398     0.2626  -0.532 0.594577
## Walc3        -0.3781     0.2767  -1.366 0.172115
## Walc4        -1.1768     0.3154  -3.731 0.000201 ***
## Walc5        -1.2831     0.4065  -3.157 0.001642 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.192 on 1039 degrees of freedom
## Multiple R-squared:  0.0201, Adjusted R-squared:  0.01633
## F-statistic: 5.328 on 4 and 1039 DF,  p-value: 0.0003004
##
## Pearson's Chi-squared test
##
## data:  df$Walc and df$RS
## X-squared = 16.5, df = 8, p-value = 0.03576
```



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 bâtons 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).

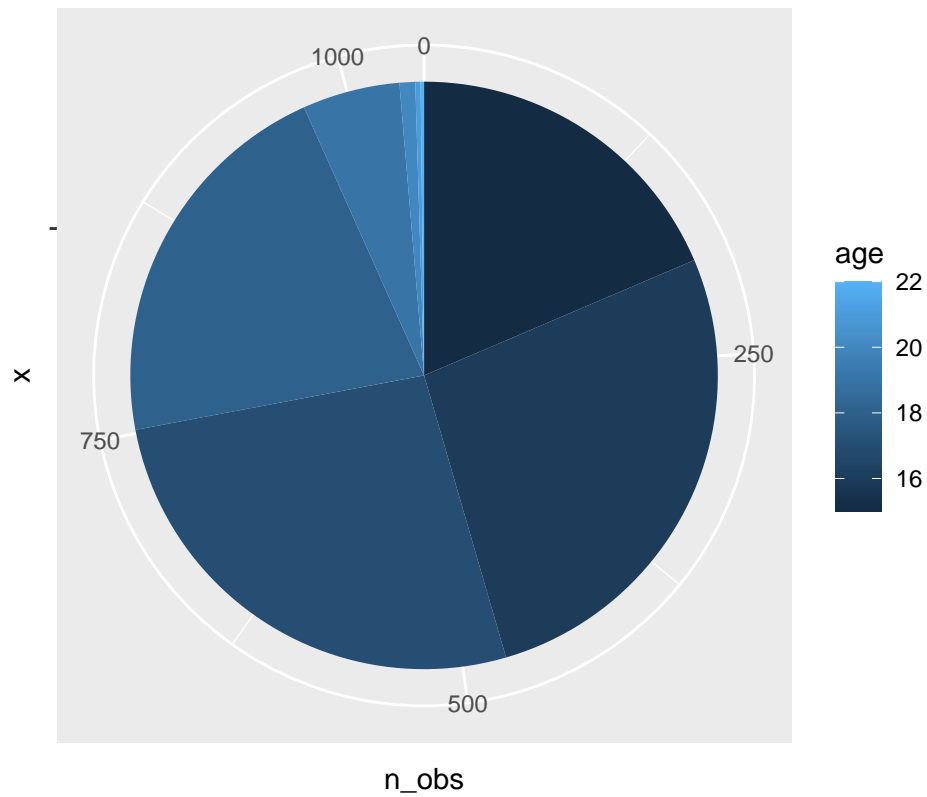
## Les variables quantitatives

##	age	Medu	Fedu	traveltime	studytime	failures	freetime	goout	Dalc	Walc	health
## 1	18	4	4	2	2	0	3	4	1	1	3
## 2	17	1	1	1	2	0	3	3	1	1	3
## 3	15	1	1	1	2	3	3	2	2	3	3
## 4	15	4	2	1	3	0	2	2	1	1	5
## 5	16	3	3	1	2	0	3	2	1	2	5
## 6	16	4	3	1	2	0	4	2	1	2	5
##	absences	G1	G2	G3							
## 1		6	5	6	6						
## 2		4	5	5	6						
## 3		10	7	8	10						
## 4		2	15	14	15						

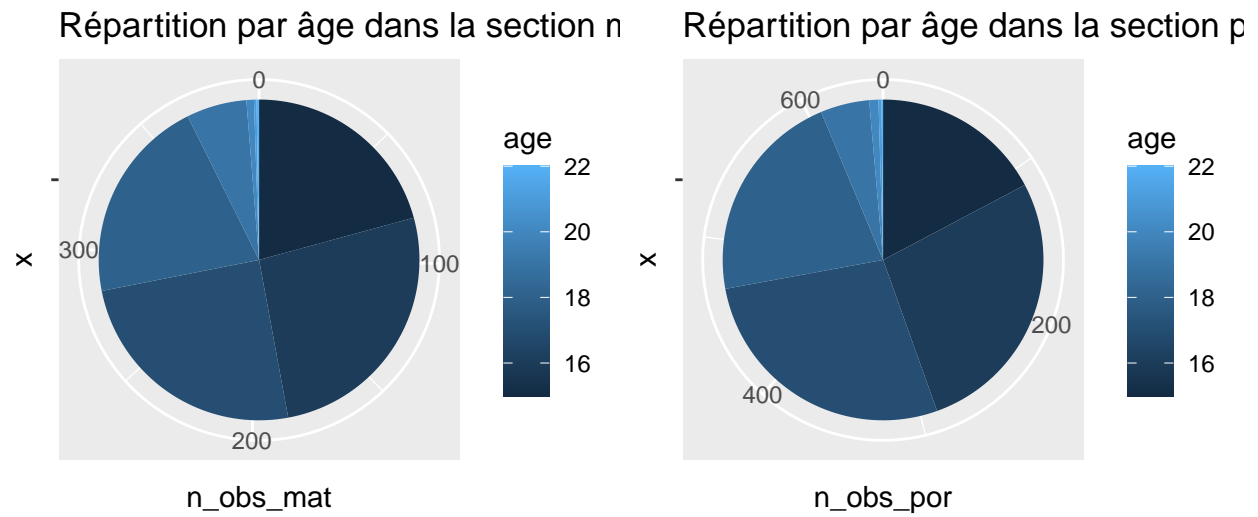
```
## 5      4  6 10 10
## 6     10 15 15 15
```

L'âge des élèves

Répartition par âge toutes filière confondue

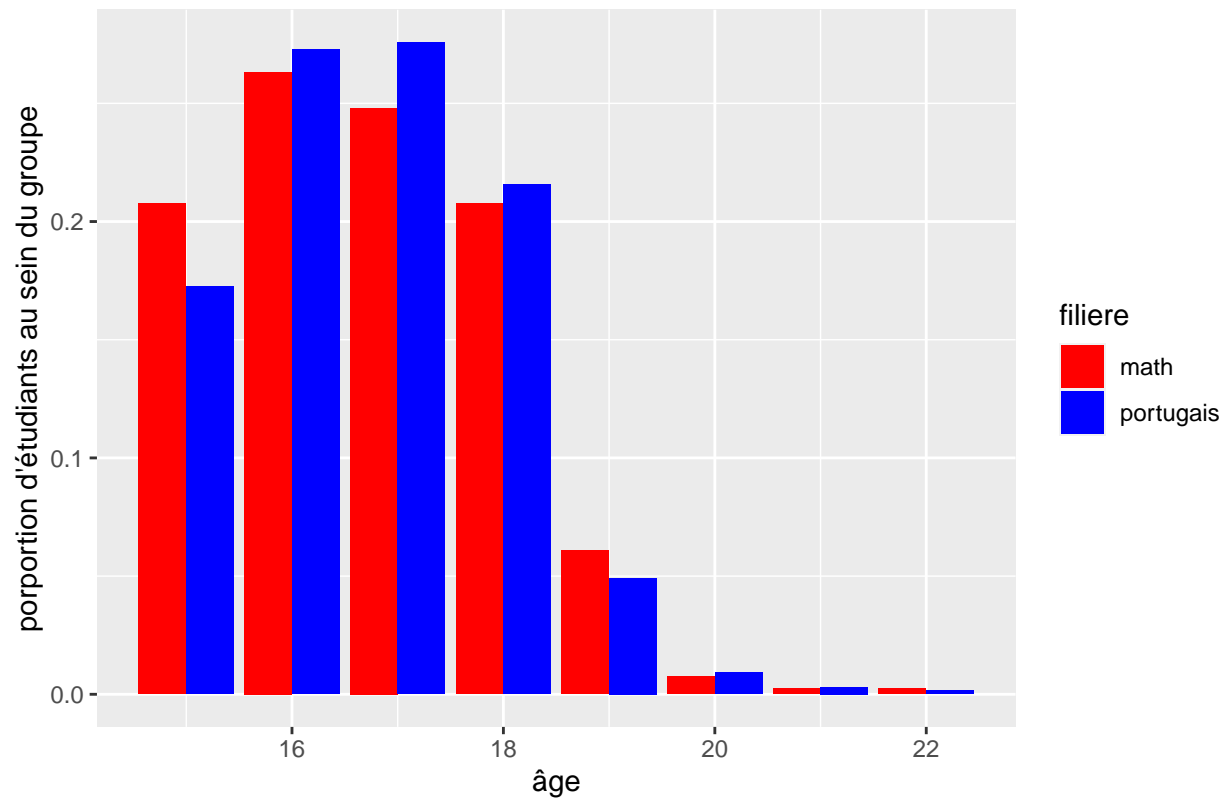


La couleur la plus claire correspond à l'âge le plus grand (22 ans), dès que l'on passe à une couleur plus foncée, on diminue l'âge de 1. On voit clairement ici que la majorité des étudiants ont entre 15 et 19 ans.



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

Comparaison des âges dans chaque filière

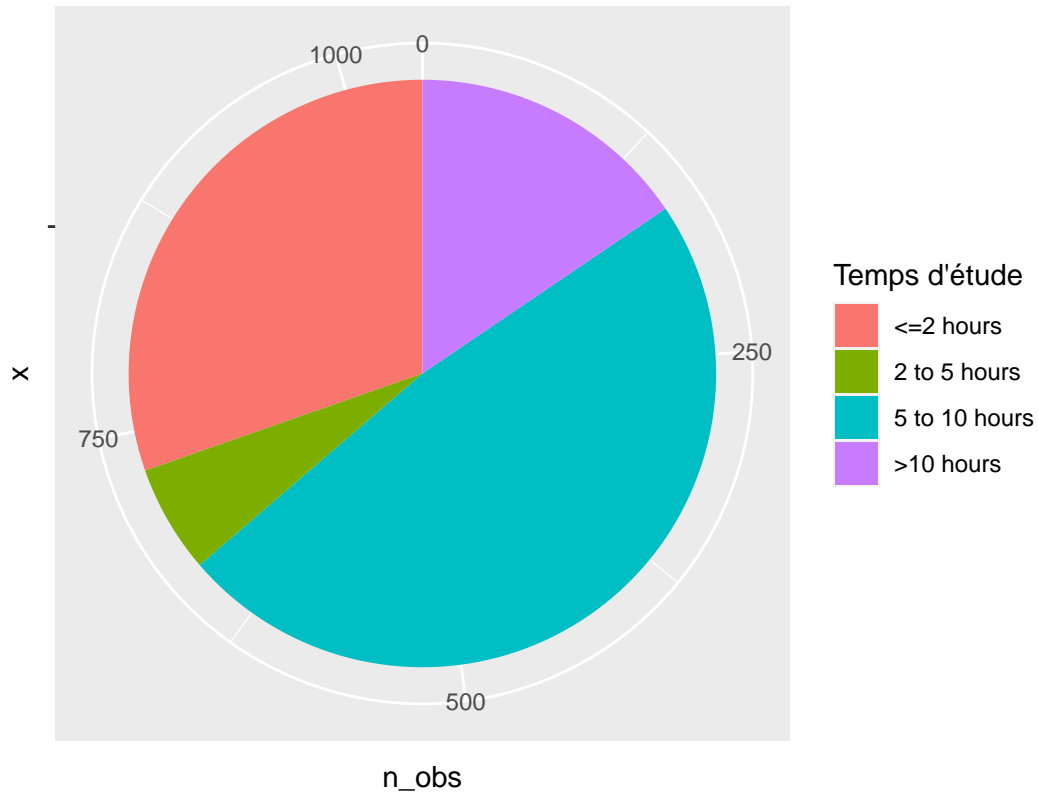


On voit que la répartition d'âge est la même dans chaque filière



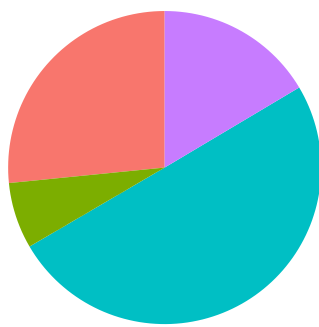
## Quantité de travail

### Répartition des temps d'étude toutes filières confondues

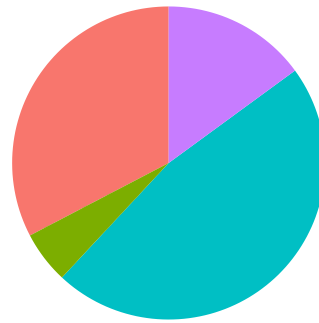


On voit clairement que les étudiants travaillent majoritairement moins de 2h00 ou entre 5h00 et 10h00 par semaines.

Temps d'étude par semaine dans la section maths (à gauche) et portugaise (à droite)



factor(studytime)

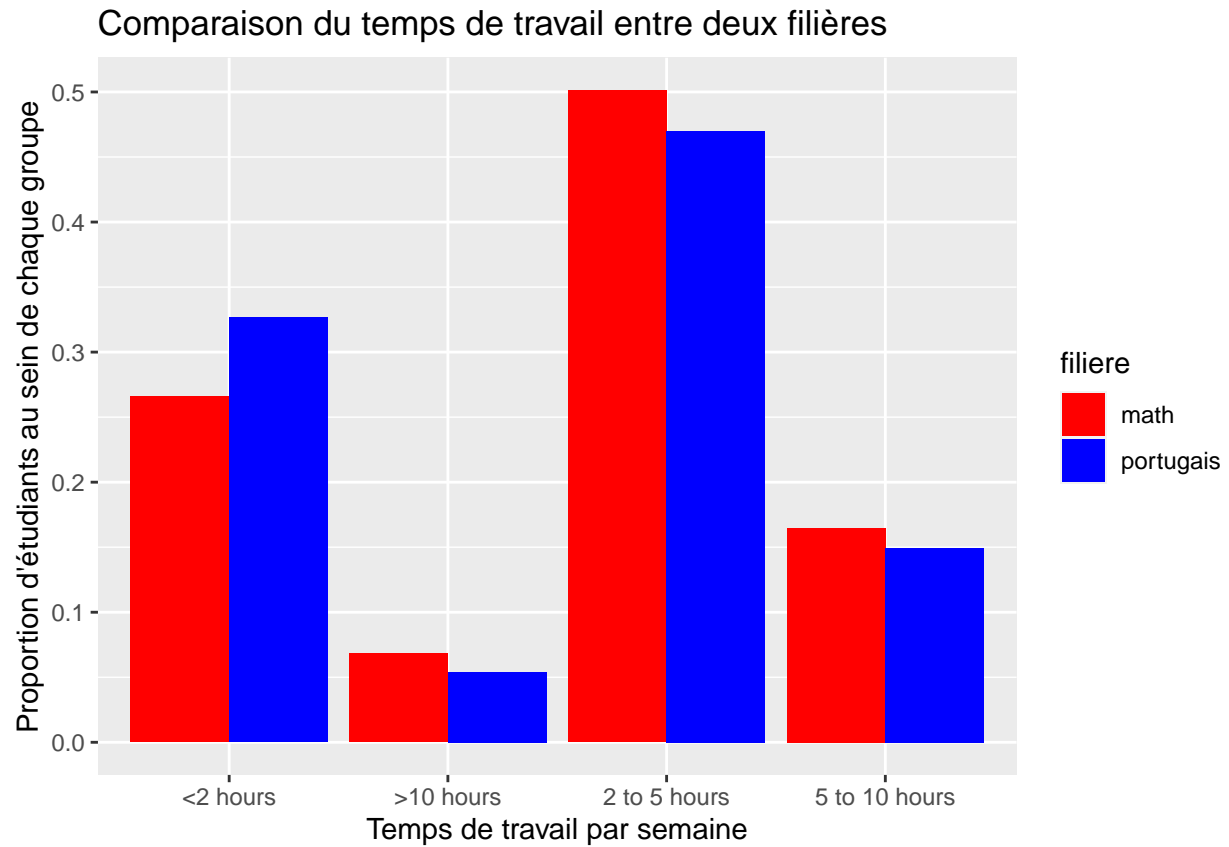


factor(studytime)

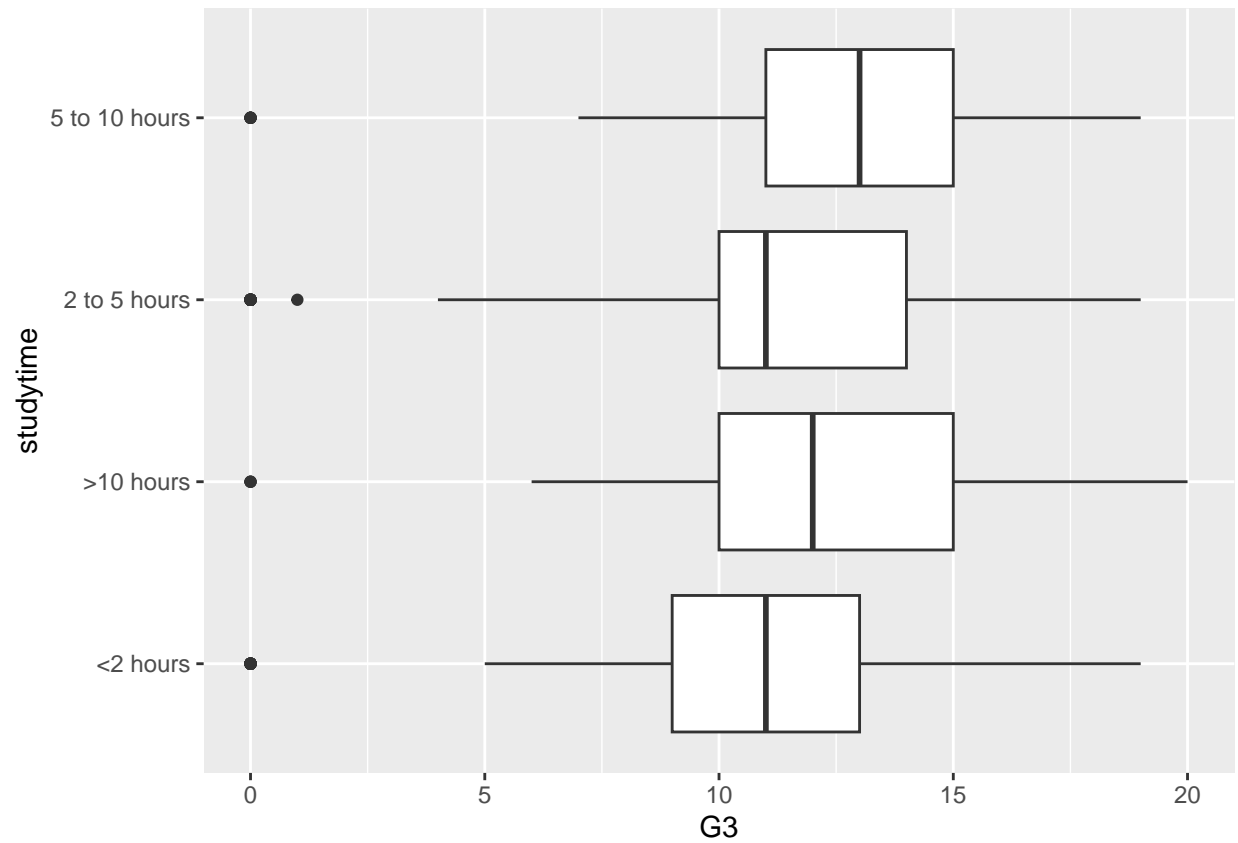


```
## # A tibble: 4 x 2
##   studytime    n_obs_mat
##   <chr>        <int>
## 1 <2 hours      105
## 2 2 to 5 hours  198
## 3 5 to 10 hours   65
## 4 >10 hours     27
```

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

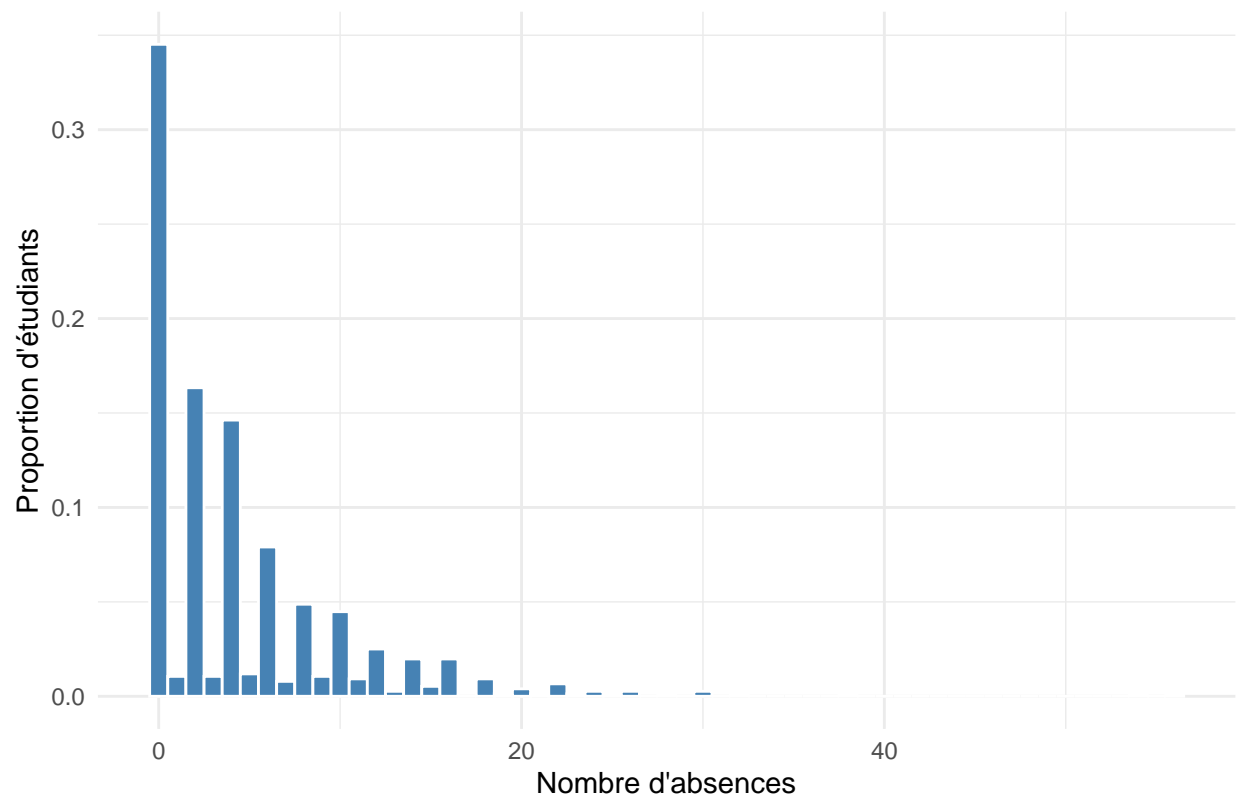


On voit que globalement, les élèves qui travaillent plus ont de meilleures notes (comportement bizarre à vérifier)

### Absences des étudiants

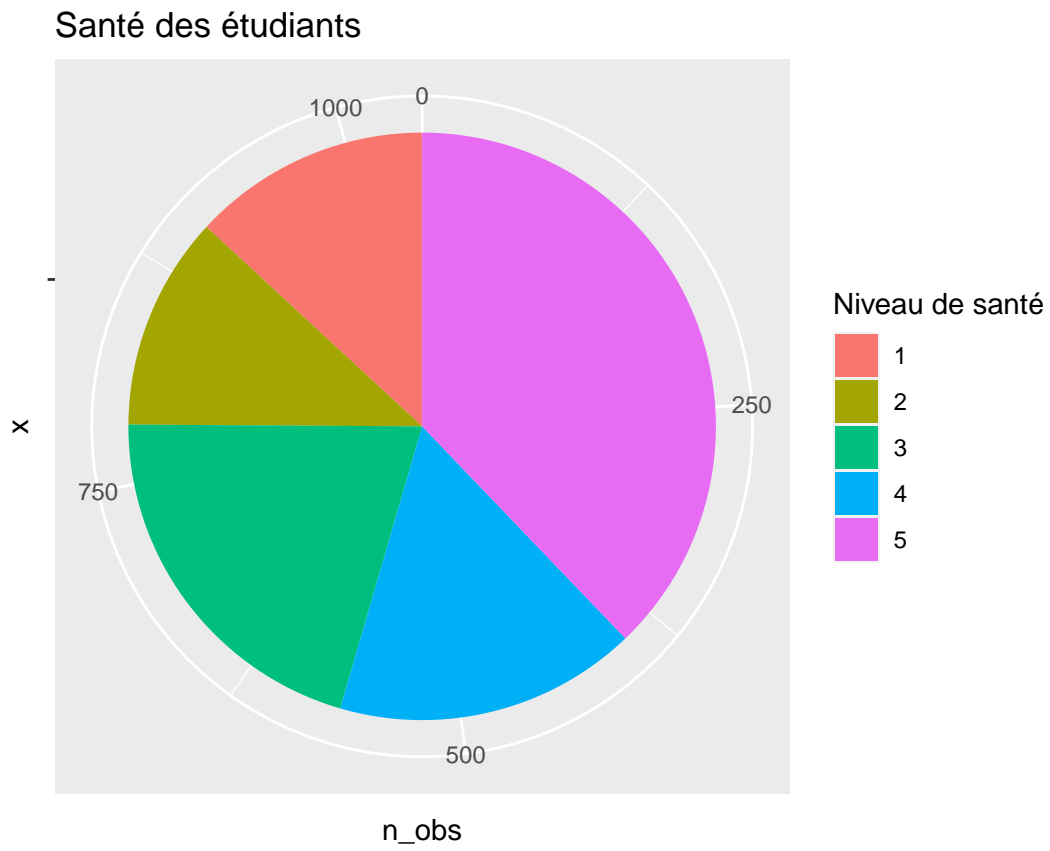
```
## Warning: The dot-dot notation (`..count..`) was deprecated in ggplot2 3.4.0.
## i Please use `after_stat(count)` instead.
```

Distribution des absences des étudiants vivants en ville

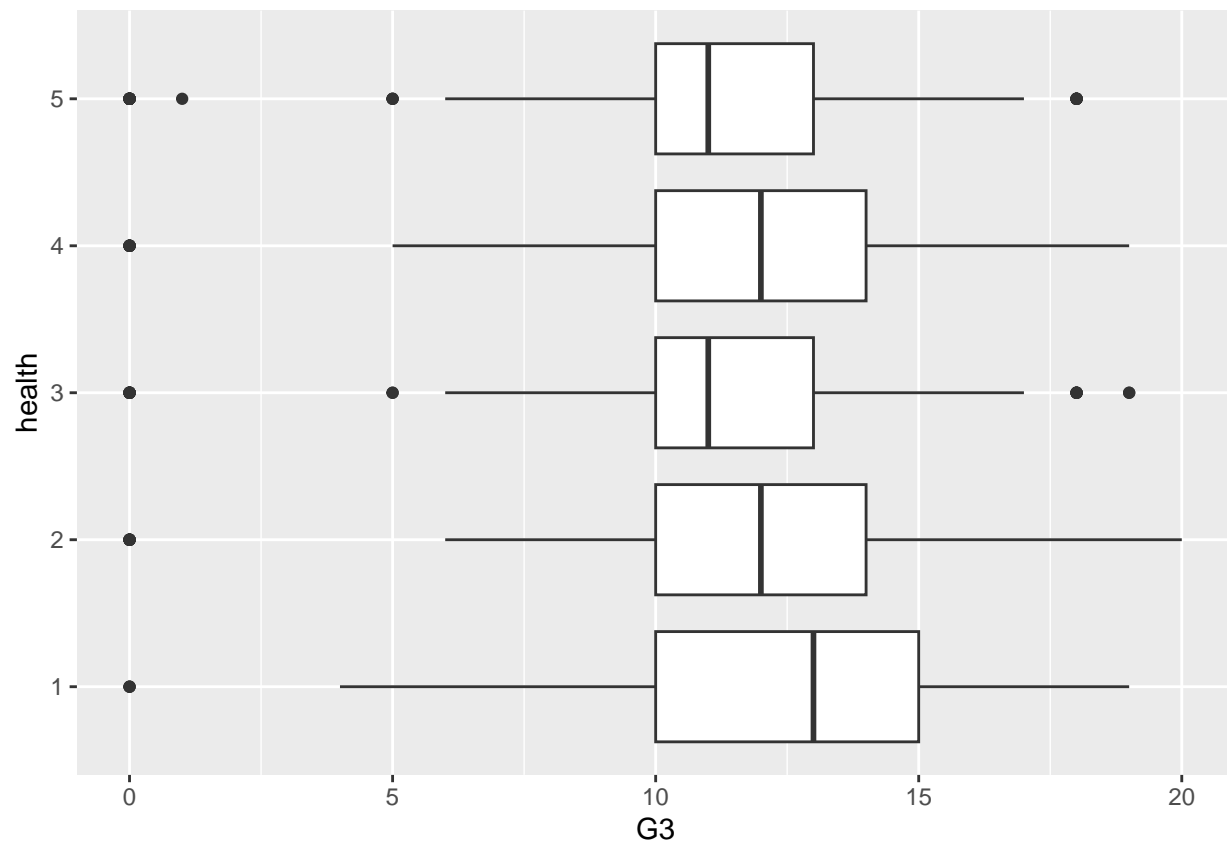


## Variables qualitatives à modalités numériques

### Santé des étudiants

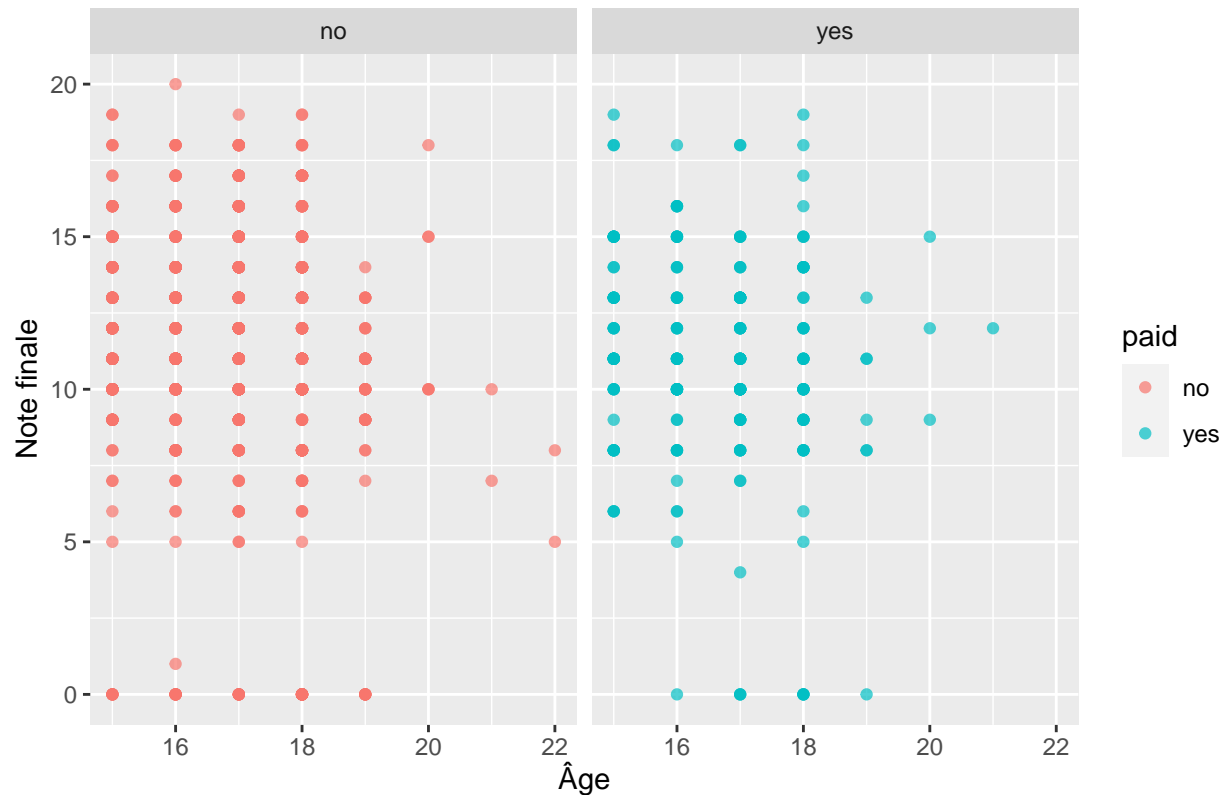


On voit que la plupart des étudiant sont en bonne santé



On voit que les étudiants en meilleure santé ont une meilleure réussite

Distribution de l'âge et de la note finale en fonction cours particuliers et de l'



Curieusement, les résultats semblent meilleur pour ceux qui n'on pas pris de cours

## 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). Nous nous intéresserons essentiellement à la comparaison des résultats de chacune des méthodes. Les méthodes utilisées seront évaluées avec leur accuracy et leur courbe ROC.

### a) Séparation du jeu de données

Ici, nous découpons notre dataset en jeu d'entraînement et jeu de test. Le ratio utilisé est  $\frac{1}{5}$  pour le jeu de test. Tout d'abord on modifie notre jeu de données pour le préparer pour la classification en retirant les notes.

```
## [1] 209
```

### b) LDA

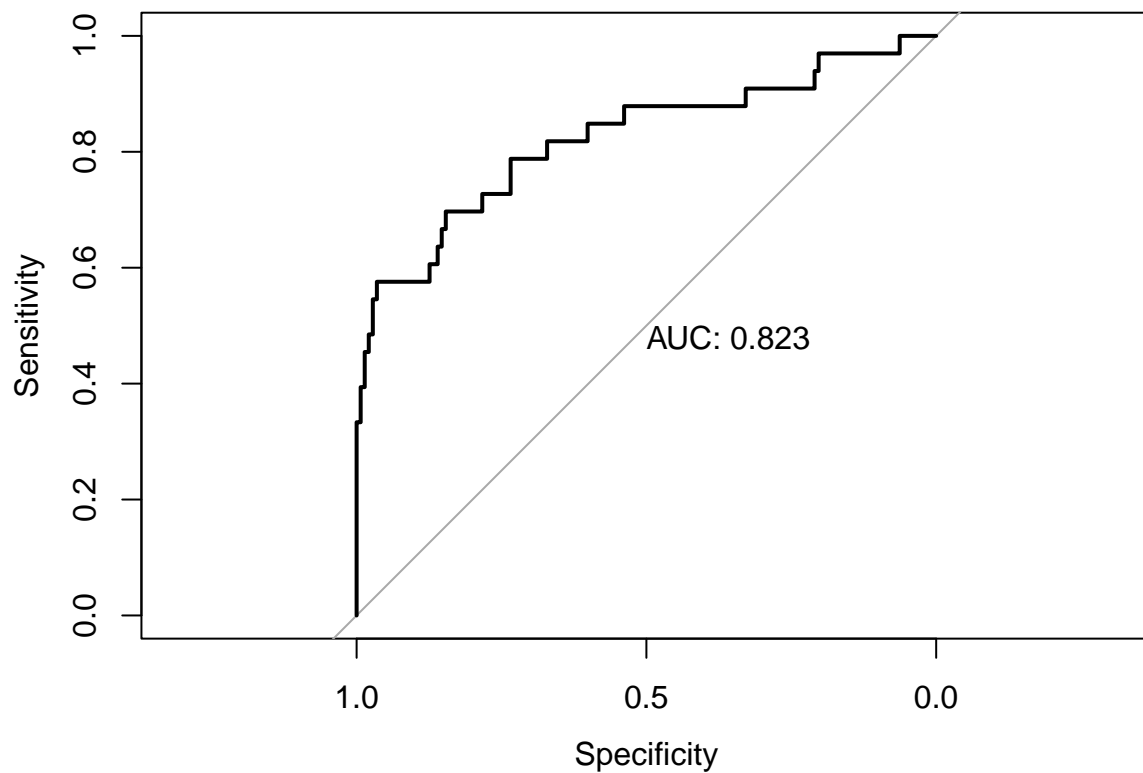
```
##
##          admis exclusion redoublement
##  admis          137         1         5
##  exclusion         17        14         2
##  redoublement      23         4         6

## Warning in roc.default(df.test$RS, pred_lda): 'response' has more than two
## levels. Consider setting 'levels' explicitly or using 'multiclass.roc' instead

## Setting levels: control = admis, case = exclusion
```



```
## Setting direction: controls < cases
```



```
## Area under the curve: 0.8226
```

```
## [1] "accuracy lda = "
```

```
## [1] 0.7511962
```

### c) QDA

```
##
```

```
##      admis exclusion redoublement
```

```
##  admis      128      12          3
```

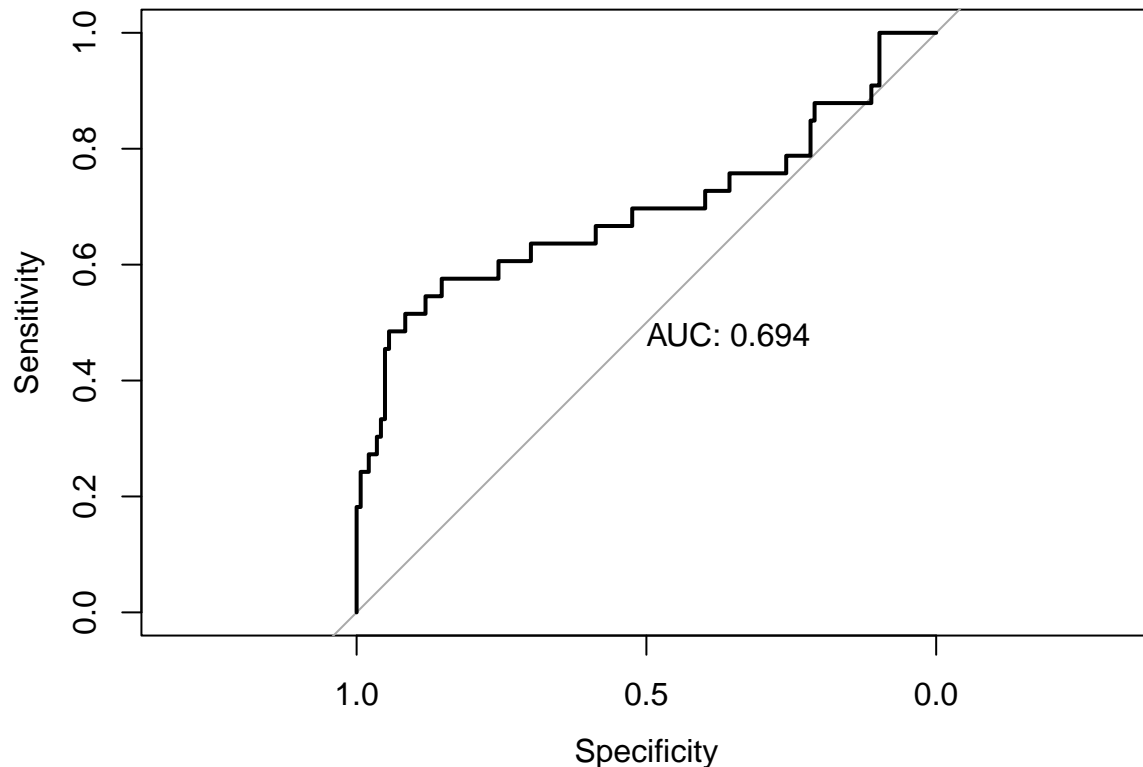
```
##  exclusion    12      17          4
```

```
##  redoublement  19       7          7
```

```
## Warning in roc.default(df.test$RS, pred_qda): 'response' has more than two
## levels. Consider setting 'levels' explicitly or using 'multiclass.roc' instead
```

```
## Setting levels: control = admis, case = exclusion
```

```
## Setting direction: controls < cases
```



```
## Area under the curve: 0.6944
```

```
## [1] "accuracy qda = "
```

```
## [1] 0.7272727
```

#### d) Stepwise

```
## `stepwise classification', using 10-fold cross-validated correctness rate of method lda'.
```

```
## 835 observations of 30 variables in 3 classes; direction: backward
```

```
## stop criterion: improvement less than 5%.
```

```
## Warning in cv.rate(vars = start.vars, data = data, grouping = grouping, :
```

```
## error(s) in modeling/prediction step
```

```
## correctness rate: 0; starting variables (30): school, sex, age, address, famsize, Pstatus, Medu, Fe
```

```
## Warning in cv.rate(trymodel, data = data, grouping = grouping, method =
```

```
## method, : error(s) in modeling/prediction step
```

```
## Warning in cv.rate(trymodel, data = data, grouping = grouping, method =
```

```
## method, : error(s) in modeling/prediction step
```

```
## Warning in cv.rate(trymodel, data = data, grouping = grouping, method =
```

```
## method, : error(s) in modeling/prediction step
```

```
## Warning in cv.rate(trymodel, data = data, grouping = grouping, method =
```

```
## method, : error(s) in modeling/prediction step
```



```

## Warning in cv.rate(trymodel, data = data, grouping = grouping, method =
## method, : error(s) in modeling/prediction step

## Warning in cv.rate(trymodel, data = data, grouping = grouping, method =
## method, : error(s) in modeling/prediction step

## Warning in cv.rate(trymodel, data = data, grouping = grouping, method =
## method, : error(s) in modeling/prediction step

## Warning in cv.rate(trymodel, data = data, grouping = grouping, method =
## method, : error(s) in modeling/prediction step

## Warning in cv.rate(trymodel, data = data, grouping = grouping, method =
## method, : error(s) in modeling/prediction step

## Warning in cv.rate(trymodel, data = data, grouping = grouping, method =
## method, : error(s) in modeling/prediction step

## Warning in cv.rate(trymodel, data = data, grouping = grouping, method =
## method, : error(s) in modeling/prediction step

##
##   hr.elapsed min.elapsed sec.elapsed
##         0.0         0.0         2.2

## method      : lda
## final model : RS ~ school + sex + age + address + famsize + Pstatus + Medu +
##   Fedu + Mjob + Fjob + reason + guardian + traveltime + studytime +
##   failures + schoolsup + famsup + paid + activities + nursery +
##   higher + internet + romantic + famrel + freetime + goout +
##   Dalc + Walc + health + absences
## <environment: 0x0000000030e34158>
##
## correctness rate = 0

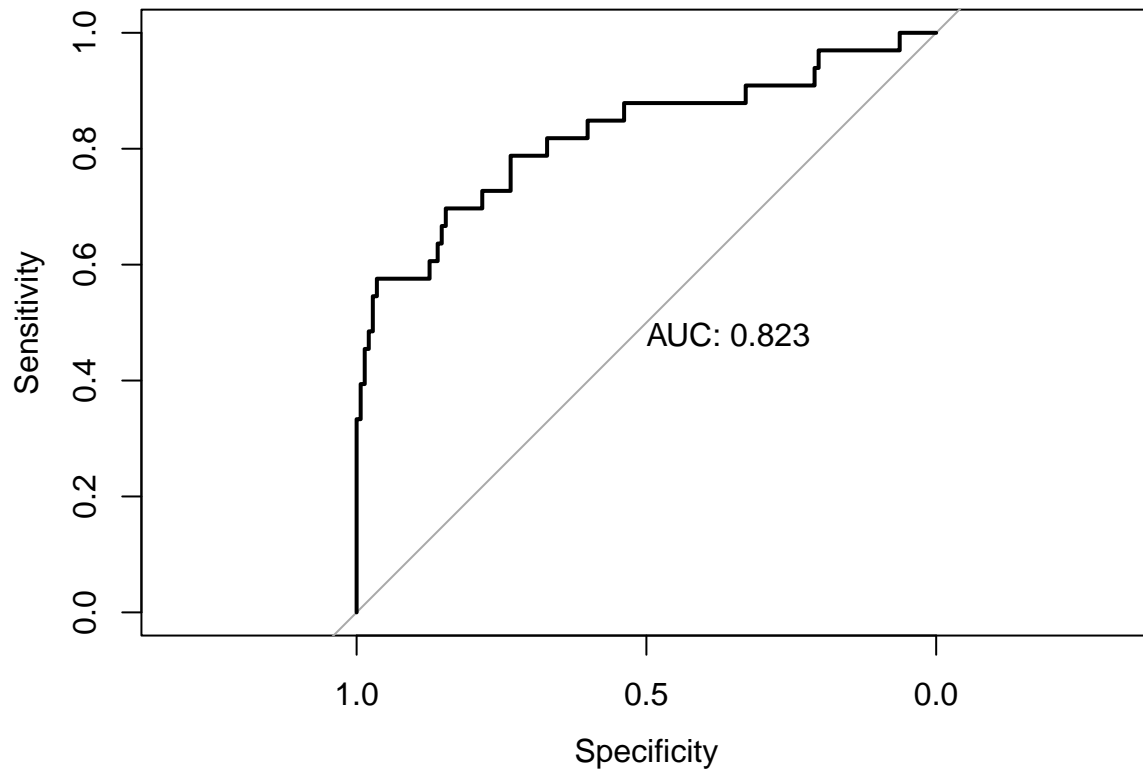
##
##           admis exclusion redoublement
##   admis           137           1           5
##   exclusion          17          14           2
##   redoublement       23           4           6

## Warning in roc.default(df.test$RS, pred_lda): 'response' has more than two
## levels. Consider setting 'levels' explicitly or using 'multiclass.roc' instead

## Setting levels: control = admis, case = exclusion

## Setting direction: controls < cases

```

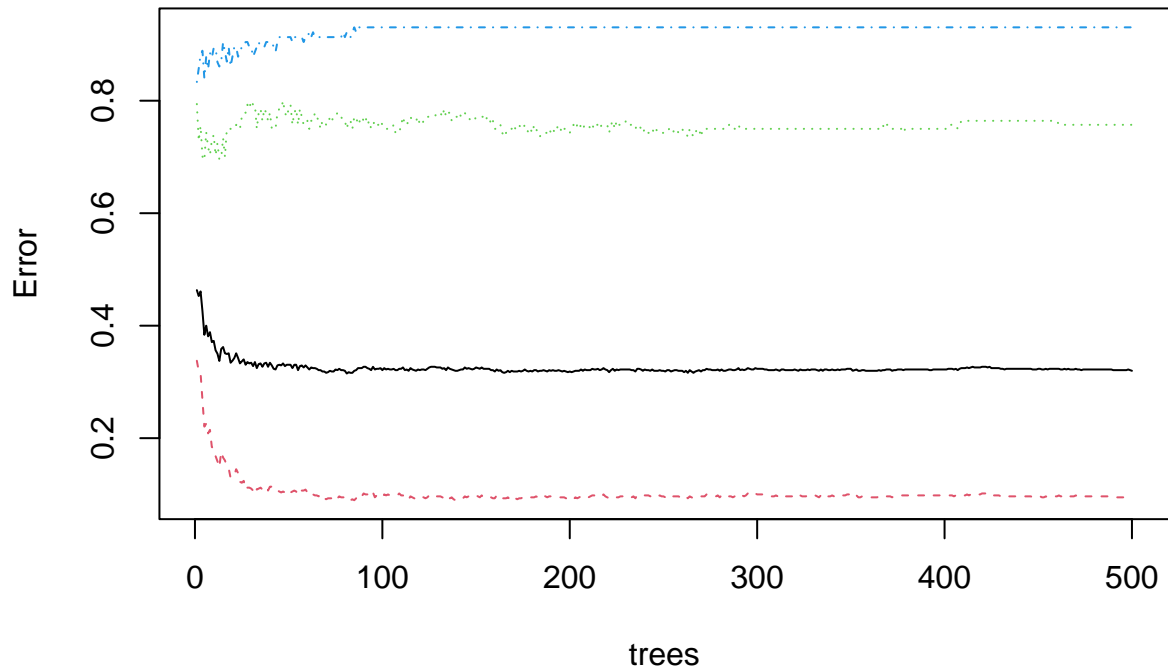


```
## Area under the curve: 0.8226
## [1] "accuracy lda stepwise = "
## [1] 0.7511962
```

#### e) Random Forest

```
##
## Call:
## randomForest(formula = RS ~ ., data = df.train)
##           Type of random forest: classification
##           Number of trees: 500
## No. of variables tried at each split: 5
##
##           OOB estimate of  error rate: 31.98%
## Confusion matrix:
##           admis exclusion redoublement class.error
## admis          526         34          20 0.09310345
## exclusion        95         34          11 0.75714286
## redoublement     87         20           8 0.93043478
```

## res\_RF



```
##
##          admis exclusion redoublement
##  admis          135         7         1
##  exclusion        16        14         3
##  redoublement     26         6         1

## Warning in roc.default(df.test$RS, pred_RF): 'response' has more than two
## levels. Consider setting 'levels' explicitly or using 'multiclass.roc' instead

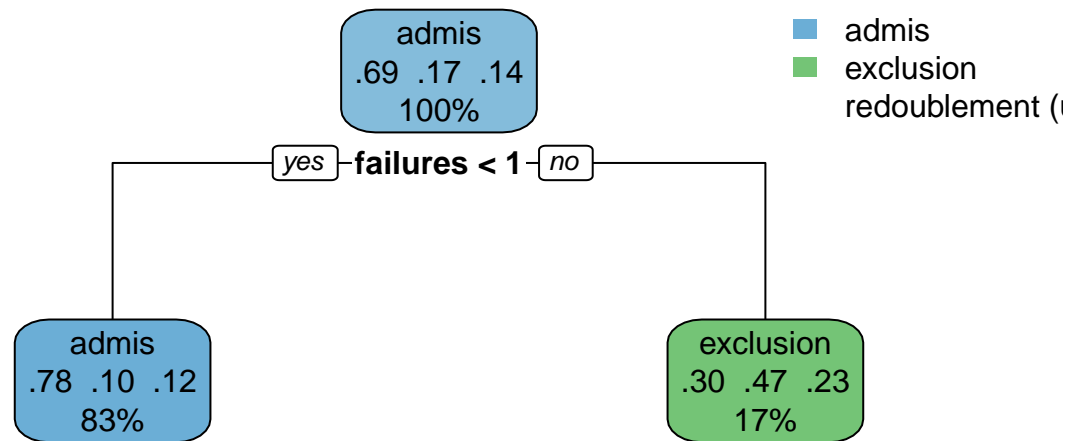
## Setting levels: control = admis, case = exclusion

## Setting direction: controls < cases

## Area under the curve: 0.8535

## [1] "accuracy RF = "
## [1] 0.7177033
```

## f) CART



```
##
##               admis exclusion redoublement
##   admis           133         10           0
##   exclusion        13         20           0
##   redoublement     24          9           0

## Warning in roc.default(df.test$RS, pred_cart): 'response' has more than two
## levels. Consider setting 'levels' explicitly or using 'multiclass.roc' instead
## Setting levels: control = admis, case = exclusion
## Setting direction: controls < cases
## Area under the curve: 0.7681
## [1] "accuracy cart = "
## [1] 0.7320574
```

## h) Adaboost

## i) Regression Logistique

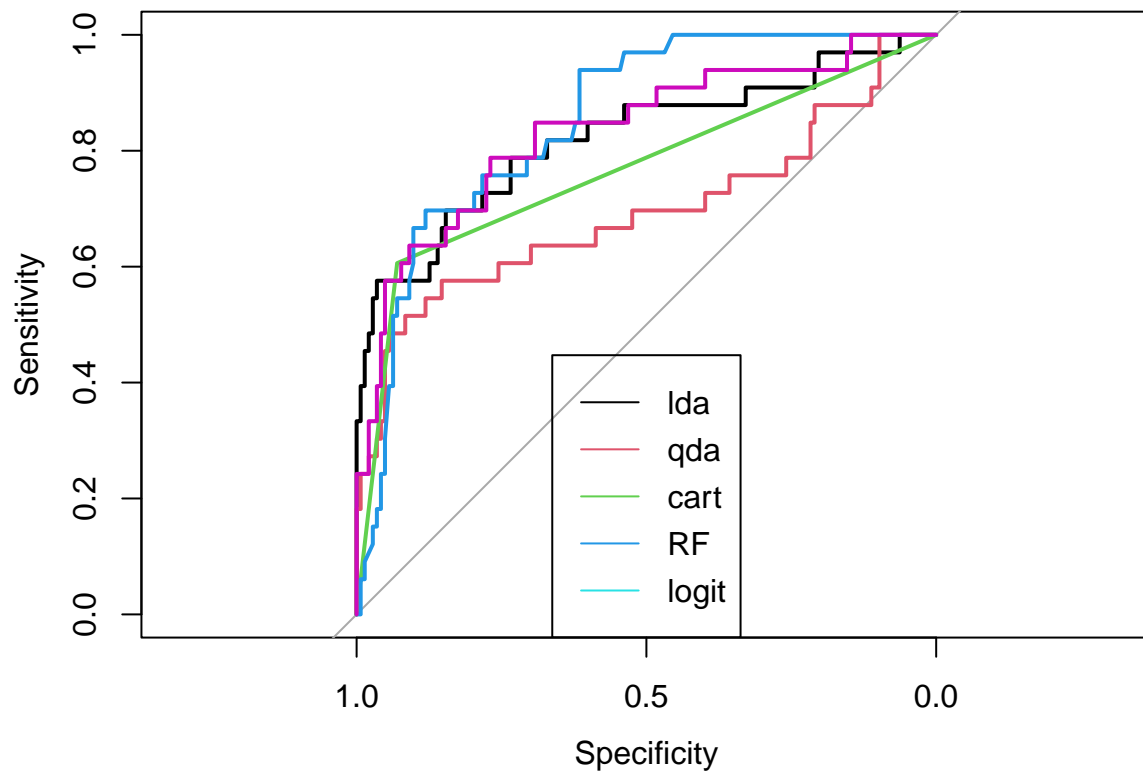
```
##               class
##               0   1
##   admis           137  6
##   exclusion        19 14
##   redoublement     23 10
```

```
## Warning in roc.default(df.test$RS, pred_logit): 'response' has more than two
## levels. Consider setting 'levels' explicitly or using 'multiclass.roc' instead
## Setting levels: control = admis, case = exclusion
## Setting direction: controls < cases
## [1] "accuracy regression logistique = "
## [1] 0.722488
## Area under the curve: 0.8339
```

## Comparison

```
##          lda      qda      cart      RF      logit
## accuracy 0.7511962 0.7272727 0.7320574 0.7177033 0.7224880
## AUC       0.8226319 0.6944268 0.7680653 0.8534647 0.8338631

## accuracy      AUC
##           1      4
```



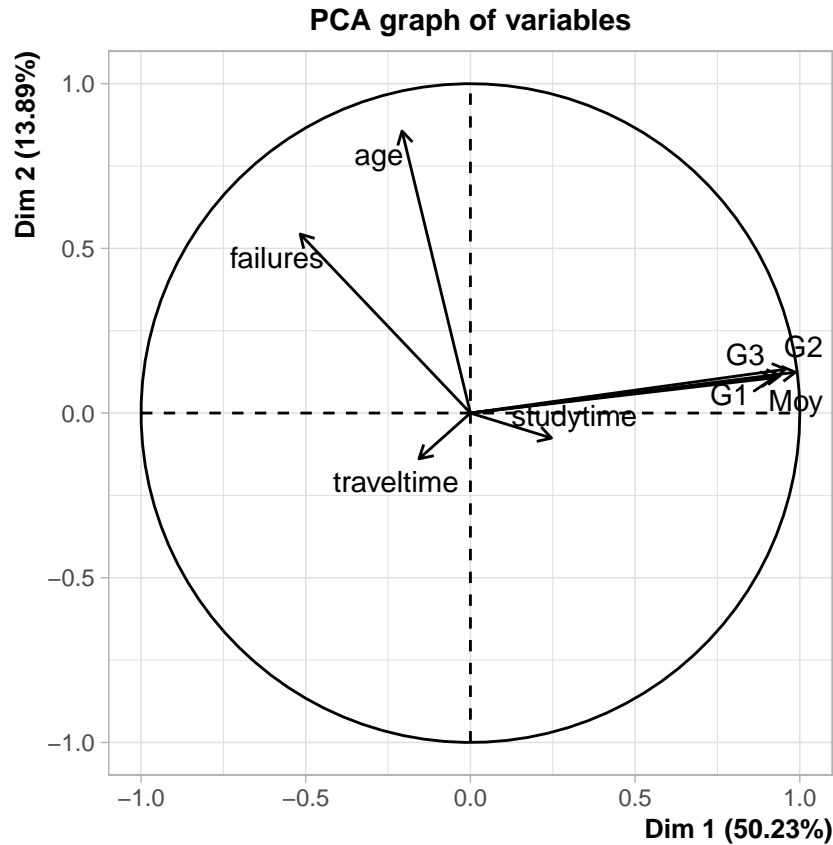
## ##ACP

```
##   age traveltime studytime failures G1 G2 G3      Moy
## 1  18          2          2         0  5  6  6  5.666667
## 2  17          1          2         0  5  5  6  5.333333
## 3  15          1          2         3  7  8 10  8.333333
## 4  15          1          3         0 15 14 15 14.666667
## 5  16          1          2         0  6 10 10  8.666667
## 6  16          1          2         0 15 15 15 15.000000
```



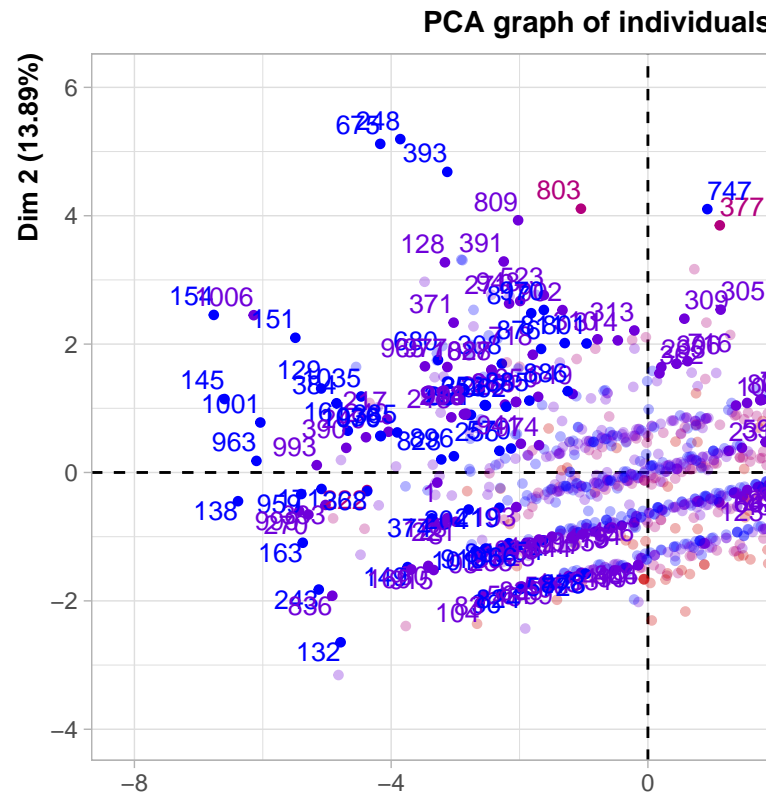
##	age	traveltime	studytime	failures	G1	G2	G3	Moy
## 1	18	22.5	210	0	5	6	6	5.666667
## 2	17	7.5	210	0	5	5	6	5.333333
## 3	15	7.5	210	3	7	8	10	8.333333
## 4	15	7.5	450	0	15	14	15	14.666667
## 5	16	7.5	210	0	6	10	10	8.666667
## 6	16	7.5	210	0	15	15	15	15.000000





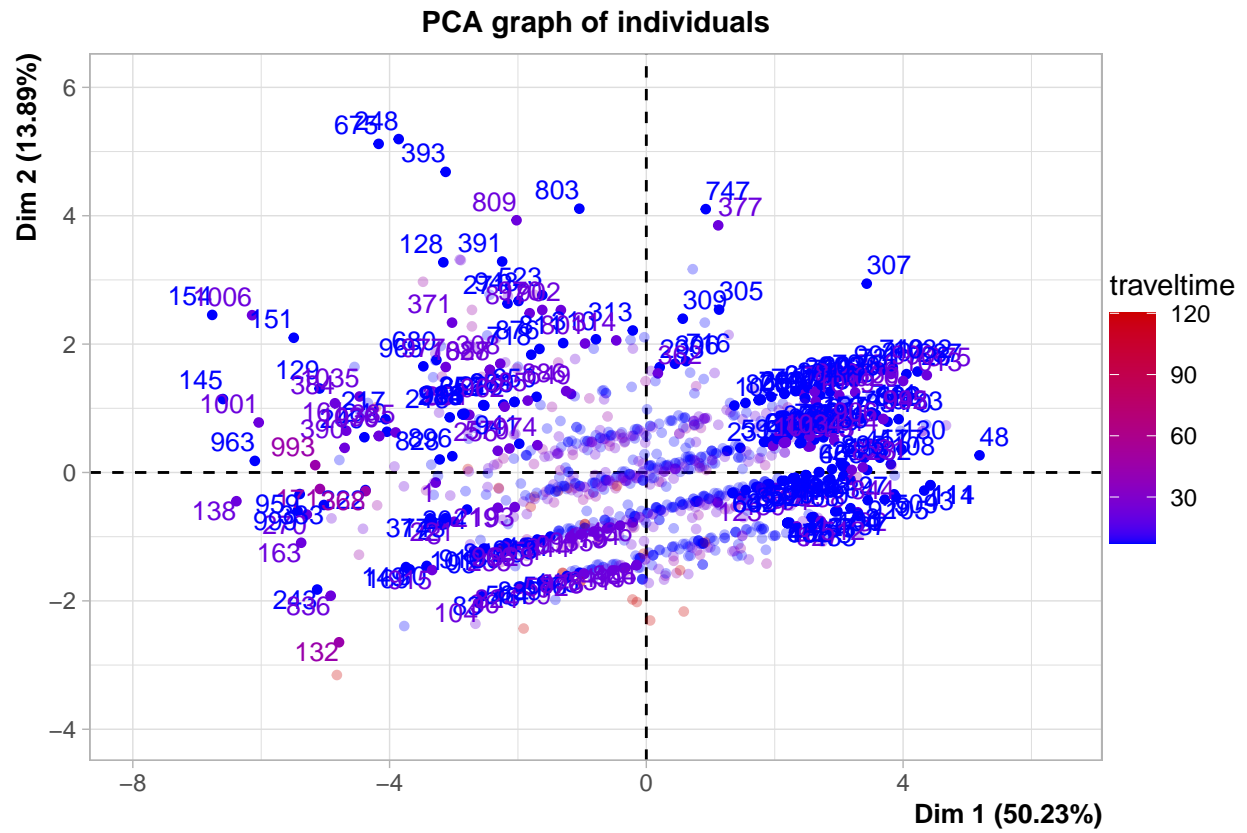
On voit que les variables study time et travel time sont mal projetées, on ne peut donc pas les interpréter. De manière logique on retrouve que les élèves ayant une bonne moyenne ont eu une bonne note à chaque semestre. Vers la gauche se trouvent les paramètres ayant une influence négative que la moyenne comme les échecs et plus curieusement l'âge (peut-être s'agit-il de personnes ayant redoublé).

##	eigenvalue	percentage of variance	cumulative percentage of variance
## comp 1	4.018158e+00	5.022697e+01	50.22697
## comp 2	1.111100e+00	1.388875e+01	64.11572
## comp 3	9.810219e-01	1.226277e+01	76.37850
## comp 4	9.605297e-01	1.200662e+01	88.38512
## comp 5	6.518265e-01	8.147831e+00	96.53295
## comp 6	1.970800e-01	2.463500e+00	98.99645
## comp 7	8.028406e-02	1.003551e+00	100.00000
## comp 8	1.593875e-30	1.992344e-29	100.00000

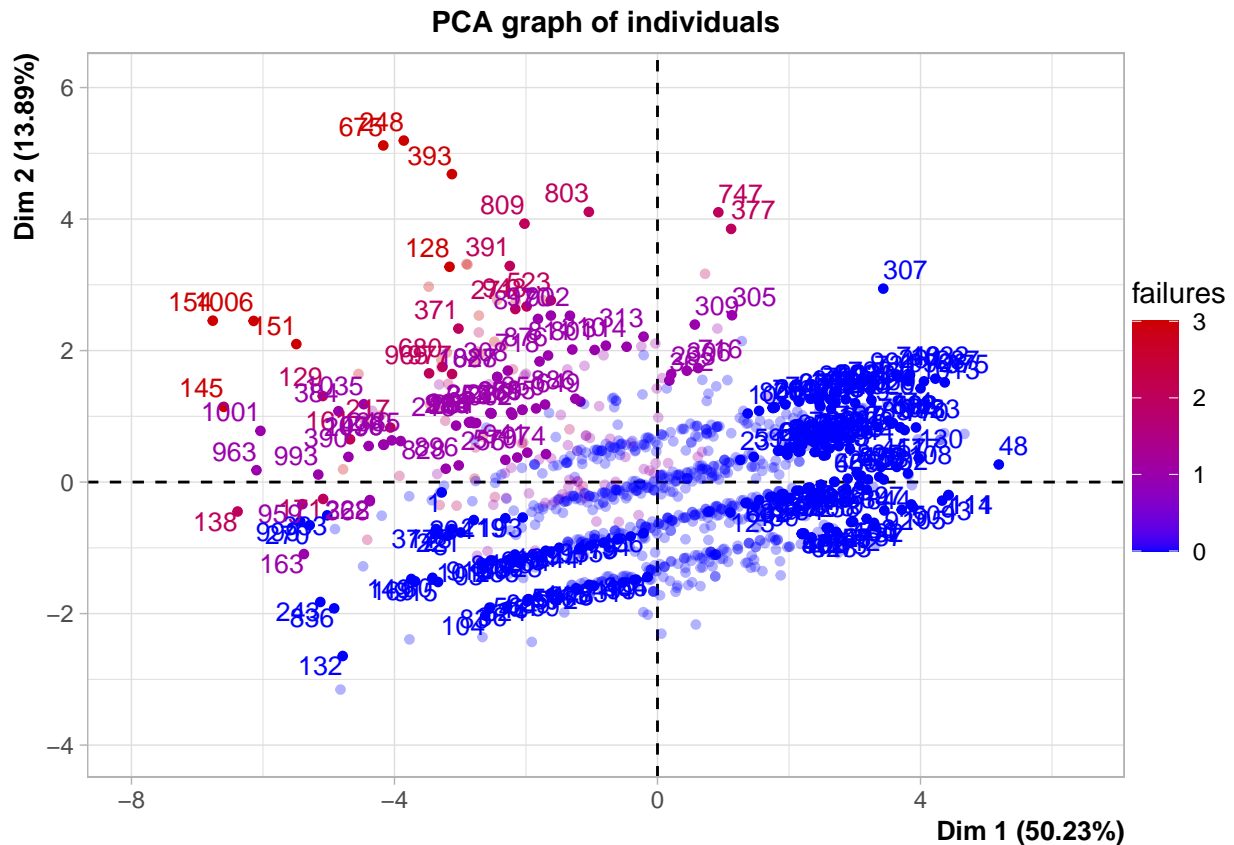


On ne garde que deux dimensions ici, d'où l'analyse ci dessus

On retrouve bien que les personnes ayant le plus travaillé se situent du côté des bonnes notes.



De la même manière on voit que le temps de trajet a une influence négative sur la réussite.



On voit aussi que les élèves qui ont les meilleurs résultats sont ceux qui ont le moins d'échecs. Par ailleurs l'acp ici ne semble pas très pertinente car la plupart des variables du jeu de données sont quantitatives, nous avons donc été obligés de les rendre (lorsque cela a un sens) qualitatives. Néanmoins on voit par exemple que pour ces variables transformées, leur projection est très mauvaise et ne peuvent donc pas être interprétés à l'aide de l'ACP (comme studytime et traveltime). Egalement peut être qu'il y a une meilleure de les rendre qualitatives. C'est pour quoi l'on va réaliser par la suite un anova 2 sur les variables quatitatives studytime et traveltime afin de pouvoir expliqués la variable Moy avec. ##Anova 2 sur les variables studytime et traveltime

```
## Les objets suivants sont masqués depuis data_quanti:
```

```
##
```

```
##      studytime, traveltime
```

```
##      traveltime studytime      Moy
```

```
## 1           2           2  5.666667
```

```
## 2           1           2  5.333333
```

```
## 3           1           2  8.333333
```

```
## 4           1           3 14.666667
```

```
## 5           1           2  8.666667
```

```
## 6           1           2 15.000000
```

```
##
```

```
##           1    2    3    4
```

```
## 1 165 314 108 36
```

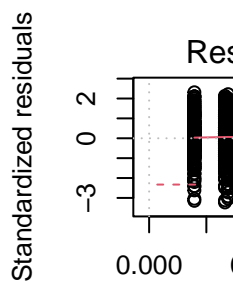
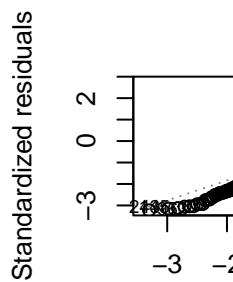
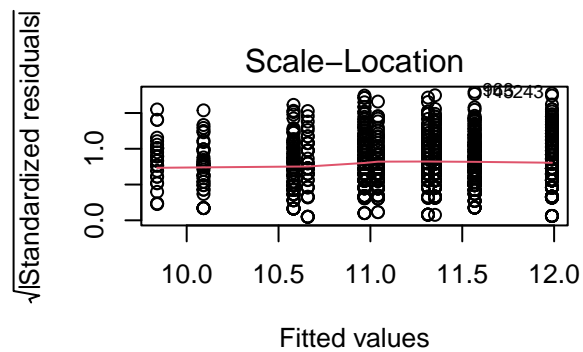
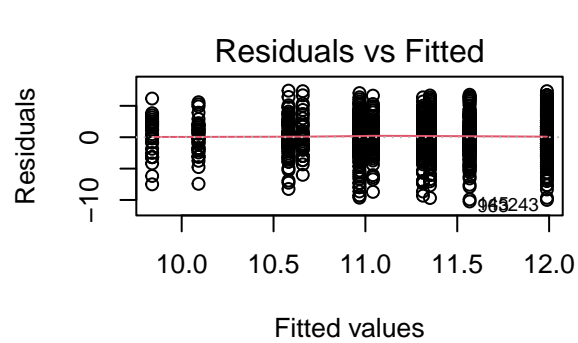
```
## 2 110 143  46 21
```

```
## 3  34  37   4  2
```

```
## 4   8   9   4  3
```

Le plan est trop déséquilibré pour faire un anova ##Anova 2 sur les variables romantic et Walc

##	school	sex	age	address	famsize	Pstatus	Medu	Fedu	Mjob	Fjob	reason													
## 1	GP	F	18	U	GT3	A	4	4	at_home	teacher	course													
## 2	GP	F	17	U	GT3	T	1	1	at_home	other	course													
## 3	GP	F	15	U	LE3	T	1	1	at_home	other	other													
## 4	GP	F	15	U	GT3	T	4	2	health	services	home													
## 5	GP	F	16	U	GT3	T	3	3	other	other	home													
## 6	GP	M	16	U	LE3	T	4	3	services	other	reputation													
##	guardian	traveltime	studytime	failures	schoolsup	famsup	paid	activities																
## 1	mother	2	2	0	yes	no	no	no																
## 2	father	1	2	0	no	yes	no	no																
## 3	mother	1	2	3	yes	no	yes	no																
## 4	mother	1	3	0	no	yes	yes	yes																
## 5	father	1	2	0	no	yes	yes	no																
## 6	mother	1	2	0	no	yes	yes	yes																
##	nursery	higher	internet	romantic	famrel	freetime	goout	Dalc	Walc	health														
## 1	yes	yes	no	no	4	3	4	1	1	3														
## 2	no	yes	yes	no	5	3	3	1	1	3														
## 3	yes	yes	yes	no	4	3	2	2	3	3														
## 4	yes	yes	yes	yes	3	2	2	1	1	5														
## 5	yes	yes	no	no	4	3	2	1	2	5														
## 6	yes	yes	yes	no	5	4	2	1	2	5														
##	absences	G1	G2	G3	Moy	RS																		
## 1	6	5	6	6	5.666667	exclusion																		
## 2	4	5	5	6	5.333333	exclusion																		
## 3	10	7	8	10	8.333333	exclusion																		
## 4	2	15	14	15	14.666667	admis																		
## 5	4	6	10	10	8.666667	redoublement																		
## 6	10	15	15	15	15.000000	admis																		
##	Walc	romantic	Moy																					
## 1	1	no	5.666667																					
## 2	1	no	5.333333																					
## 3	3	no	8.333333																					
## 4	1	yes	14.666667																					
## 5	2	no	8.666667																					
## 6	2	no	15.000000																					
##																								
##	no	yes																						
## 1	253	145																						
## 2	151	84																						
## 3	127	73																						
## 4	98	40																						
## 5	44	29																						



Le modèle est complet et n'est pas trop déséquilibré.

```
##
## Shapiro-Wilk normality test
##
## data: res$residuals
## W = 0.98902, p-value = 4.858e-07
```

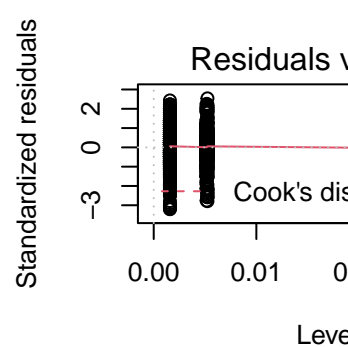
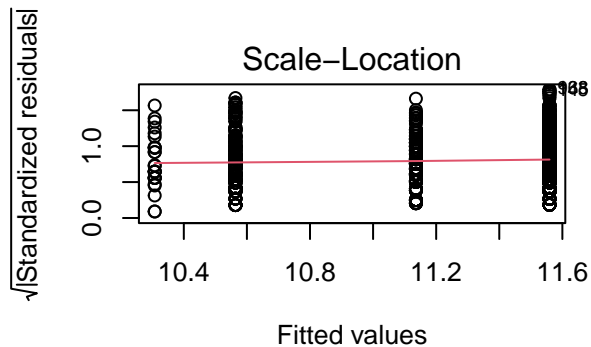
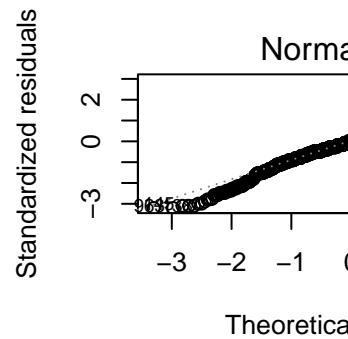
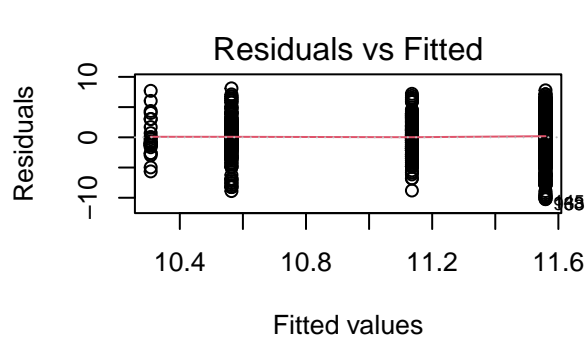
Les données ne sont pas du tout gaussiennes.

```
## school sex age address famsize Pstatus Medu Fedu Mjob Fjob reason
## 1 GP F 18 U GT3 A 4 4 at_home teacher course
## 2 GP F 17 U GT3 T 1 1 at_home other course
## 3 GP F 15 U LE3 T 1 1 at_home other other
## 4 GP F 15 U GT3 T 4 2 health services home
## 5 GP F 16 U GT3 T 3 3 other other home
## 6 GP M 16 U LE3 T 4 3 services other reputation
## guardian traveltime studytime failures schoolsup famsup paid activities
## 1 mother 2 2 0 yes no no no
## 2 father 1 2 0 no yes no no
## 3 mother 1 2 3 yes no yes no
## 4 mother 1 3 0 no yes yes yes
## 5 father 1 2 0 no yes yes no
## 6 mother 1 2 0 no yes yes yes
## nursery higher internet romantic famrel freetime goout Dalc Walc health
## 1 yes yes no no 4 3 4 1 1 3
## 2 no yes yes no 5 3 3 1 1 3
## 3 yes yes yes no 4 3 2 2 3 3
## 4 yes yes yes yes 3 2 2 1 1 5
```

```
## 5      yes      yes      no      no      4      3      2      1      2      5
## 6      yes      yes      yes      no      5      4      2      1      2      5
##      absences G1 G2 G3      Moy      RS
## 1          6 5 6 6 5.666667      exclusion
## 2          4 5 5 6 5.333333      exclusion
## 3         10 7 8 10 8.333333      exclusion
## 4          2 15 14 15 14.666667      admis
## 5          4 6 10 10 8.666667      redoublement
## 6         10 15 15 15 15.000000      admis

##      Moy paid internet
## 1 5.666667      no      no
## 2 5.333333      no      yes
## 3 8.333333      yes      yes
## 4 14.666667      yes      yes
## 5 8.666667      yes      no
## 6 15.000000      yes      yes

##
##      no yes
##      no 191 26
##      yes 633 194
```



Le plan est complet et quasiment équilibré

```
##
##      Shapiro-Wilk normality test
##
## data:  res$residuals
```



```
## W = 0.99028, p-value = 2.162e-06
```

On obtiens encore que les données ne sont pas gaussiennes

```
##Modèle linéaire Gaussien: Régression mutiple
```

```
## Le chargement a nécessité le package : carData
```

```
##
```

```
## Attachement du package : 'car'
```

```
## L'objet suivant est masqué depuis 'package:dplyr':
```

```
##
```

```
##      recode
```

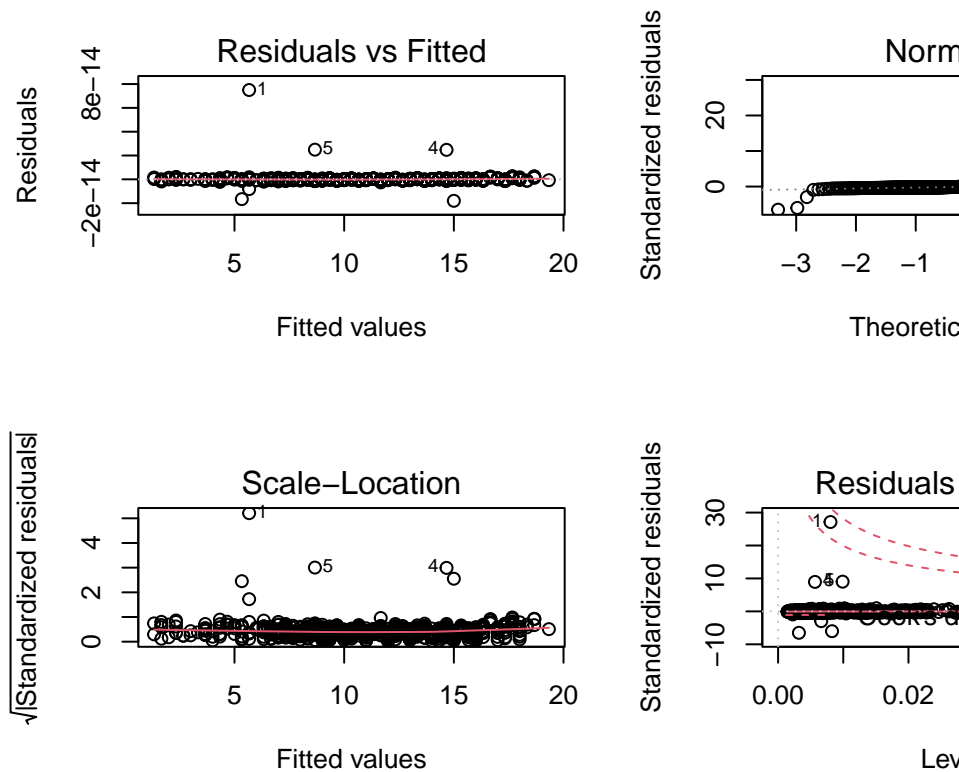
```
## Warning in summary.lm(object, ...): essentially perfect fit: summary may be
```

```
## unreliable
```

```
##      age traveltime studytime failures      G1      G2      G3
```

```
##  1.087837  1.027492  1.048104  1.278165  3.943316  7.959427  6.069691
```

Aucune valeur n'est plus grande que 10, la matrice est donc de plein rang. On va maintenant vérifier si les



résidus sont iid, gaussiens centrée et réduits

On voit qu'il n'y a pas de forme de trompette sur le graphe des résidus donc l'hypothèse d'homoscédasticité est vérifiée. Néanmoins il semble y avoir plusieurs points avec des résidus trop grands.

```
##      1      2      4      5      6      8
```

```
## 50.412602  6.120782  9.311399  9.394551  6.639145  2.980612
```

En effet, on voit qu'il y en a huit. Il faudrait enlever le point le plus éloigné. Néanmoins, on voit en regardant le qqplot nos variables n'ont aucune chance d'être gaussiennes. En effet, avec la p-valeur du test de Shapiro qui est très petite devant 5%, on rejette  $H_0$ , les données ne sont donc pas gaussiennes. Le modèle n'est donc pas adapté.

```
##  
## Shapiro-Wilk normality test  
##  
## data: reg$residuals  
## W = 0.17122, p-value < 2.2e-16
```