

# Rapport CYS

20 février, 2020

## Contents

Etude des données	1
Les données . . . . .	1
Test d'un modèle ANOVA de 3 facteurs(CYS S3, CMI, TP S3)	8
Test d'un modèle ANCOVA de 3 facteurs qualitatives (CYS S3, CMI, TP S3) et 1 facteur quantitative(Snapshot1) pour modéliser le Snapshot2	13
Train+ Test par le modèle ANCOVA . . . . .	16
Conclusion:	18

## Etude des données

### Les données

```
CYS3 = read.csv("Semestre3_Complet.csv")
summary(CYS3)
```

```
##      prénom      Semestre  Filière      snapshot.1      snapshot.2
## ALEXANDRE:  4      S3 2017-18: 54      EEA:242      : 30      Min.      : 1.75
## ALEXIS      :  3      S3 2018-19:188      5.5      : 11      1st Qu.: 7.00
## GUILLAUME:  3      7.75      : 11      Median : 8.75
## HUGO      :  3      3.5      : 10      Mean   : 8.84
## LUCAS      :  3      7      : 10      3rd Qu.:10.75
## NICOLAS    :  3      7.5      : 10      Max.   :15.75
## (Other)    :223      (Other):160      NA's   :27
## Snapshot.2...4m CYS.S3      CYS.S4      TP.S3      TP.S4      CMI      Groupe.S3
##      :169      ?      : 1      :185      : 23      :192      : 19      Siuban :59
## -      : 21      non:165      non: 28      FR:155      FR: 21      ?      : 2      Akane  :44
## 11      :  4      oui: 76      oui: 29      GB: 64      GB: 29      non:184      Steven :40
## 11.75    :  4      oui: 37      Alba    :38
## 14.5     :  4      Nadia   :29
## 6.75     :  4      Yolanda:20
## (Other): 36      (Other):12
##      Groupe.S4      Prof.TP
##      :210      :220
## Nadia      : 17      Didier: 10
## Virginia: 15      Pierre: 12
##
##
##
##
```

```
N=242
```

```
table(CYS3$CYS.S3)/N*100
```

```
##
```

```
##      ?      non      oui
```

```
## 0.4132231 68.1818182 31.4049587
```

```
#pie(table(CYS3$CYS.S3))
```

Commentaire:

On trouve que presque un tiers des étudiants utilisent l'outil CheckYourSmile en semestre 3.

```
table(CYS3$CMI)/N*100
```

```
##
```

```
##      ?      non      oui
```

```
## 7.8512397 0.8264463 76.0330579 15.2892562
```

```
#pie(table(CYS3$CMI))
```

Commentaire:

On trouve que presque 15% des étudiants sont CMI en semestre 3. Or, presque 9% des étudiants ne donnent pas une telle réponse pertinente.

```
table(CYS3$TP.S3)/N*100
```

```
##
```

```
##      FR      GB
```

```
## 9.504132 64.049587 26.446281
```

```
#pie(table(CYS3$TP.S3))
```

Commentaire:

On trouve que presque 64% des étudiants pratiquent les TP en français en semestre 3 et 26% pratiquent les TP en anglais. Or, presque 10% des étudiants ne donnent pas une telle réponse pertinente.

Commentaire:

On trouve que presque un tiers

```
library(gmodels)
```

```
CrossTable(CYS3$CYS.S3,CYS3$CMI)
```

```
CrossTable(CYS3$CYS.S3,CYS3$TP.S3)
```

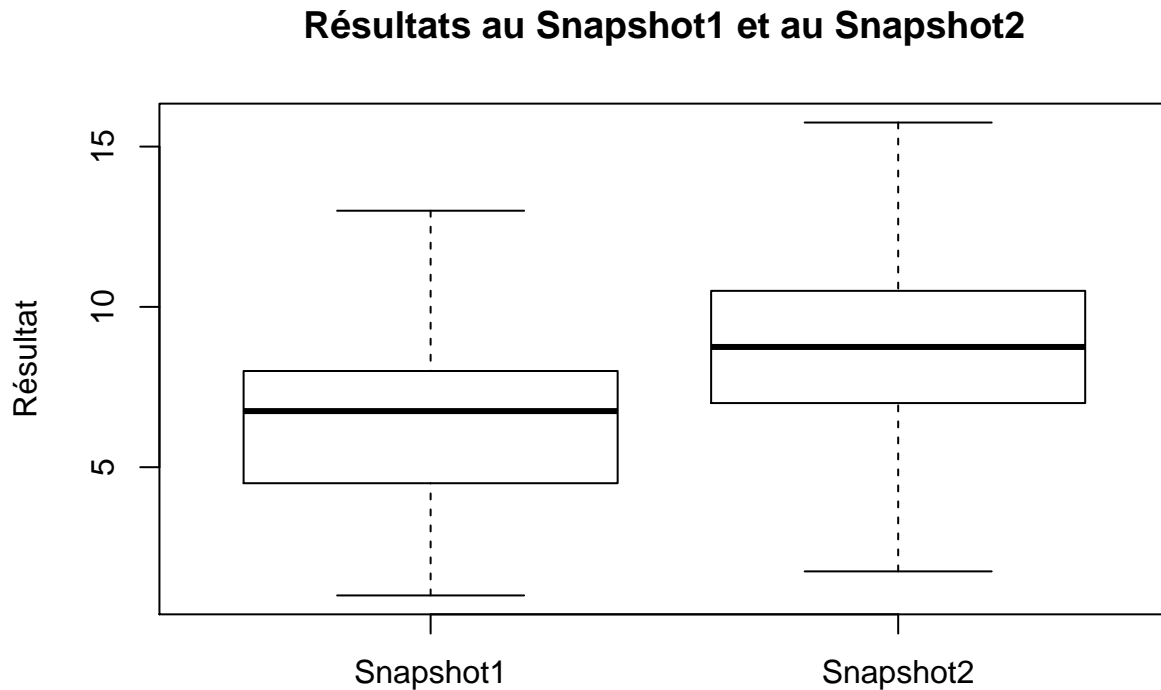
```
CYS = read.csv("W1_version1.csv")
```

```
summary(CYS)
```

```
##      prénom      Semestre      Filière      snapshot.1      snapshot.2
## ALEXANDRE: 3      S3 2017-18: 38      EEA:181      Min. : 1.000      Min. : 1.75
## ALEXIS : 3      S3 2018-19:143
## HUGO : 3
## LUCAS : 3
## NICOLAS : 3
## VINCENT : 3
## (Other) :163
## Snapshot.2...4m CYS.S3      CYS.S4      TP.S3      TP.S4      CMI      Groupe.S3
## :127      non:120      :137      FR:132      :143      non:150      Siuban :49
## - : 15      oui: 61      non: 21      GB: 49      FR: 15      oui: 31      Akane :34
```

```
## 11      : 3          oui: 23          GB: 23          Alba   :29
## 11.75    : 3
## 14.5     : 3
## 6.75     : 3
## (Other): 27          Nadia  :24
##          Groupe.S4   Prof.TP
##          :154        :168
## Nadia    : 15        Didier: 6
## Virginia: 12        Pierre: 7
##
##
##
##
```

```
# Stat. descriptives à compléter
boxplot(CYS$snapshot.1,CYS$snapshot.2,names=c("Snapshot1","Snapshot2"), ylab="Résultat",
        main="Résultats au Snapshot1 et au Snapshot2")
```



Au regard de boxplot, on constate que le Snapshot2 prend souvent les valeurs plus grandes que le Snapshot1 d'où la progression obtenue en résultat.

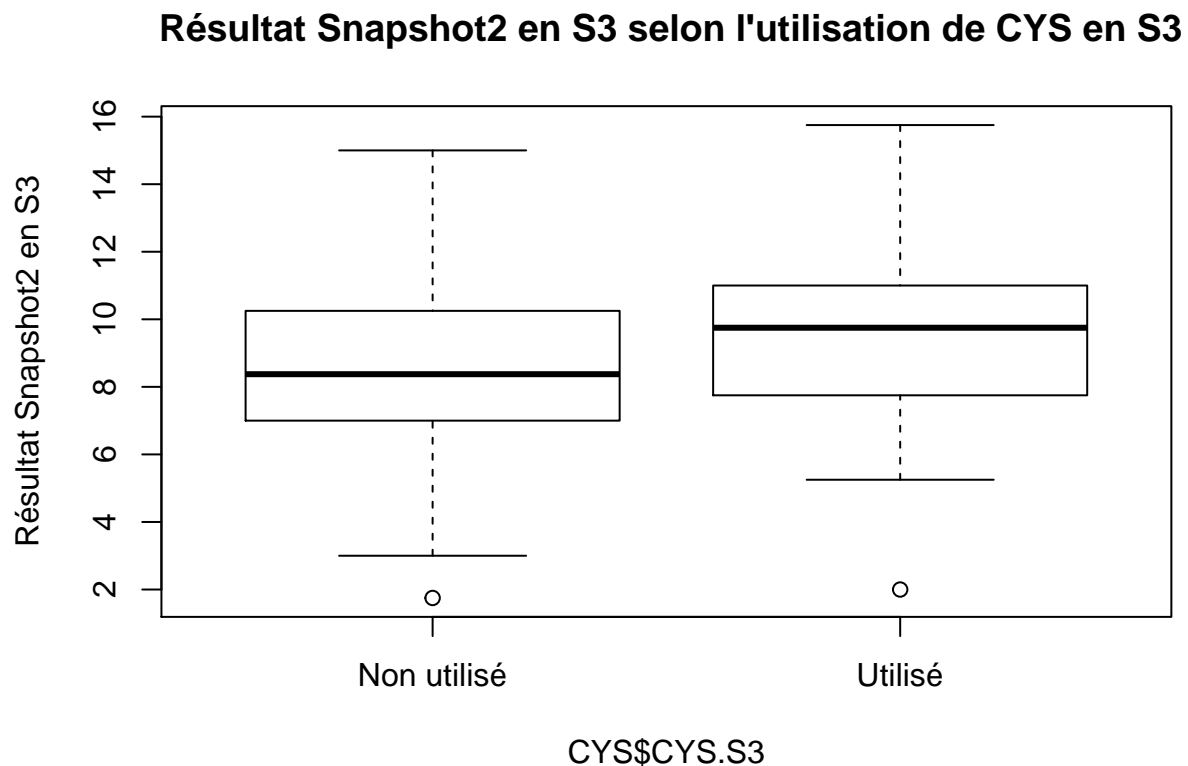
```
# Stat. descriptives à compléter
ks.test(CYS$snapshot.1,CYS$snapshot.2,alternative = "greater")
```

```
## Warning in ks.test(CYS$snapshot.1, CYS$snapshot.2, alternative = "greater"): p-
## value will be approximate in the presence of ties
##
```

```
## Two-sample Kolmogorov-Smirnov test
##
## data: CYS$snapshot.1 and CYS$snapshot.2
## D+ = 0.35912, p-value = 7.286e-11
## alternative hypothesis: the CDF of x lies above that of y
```

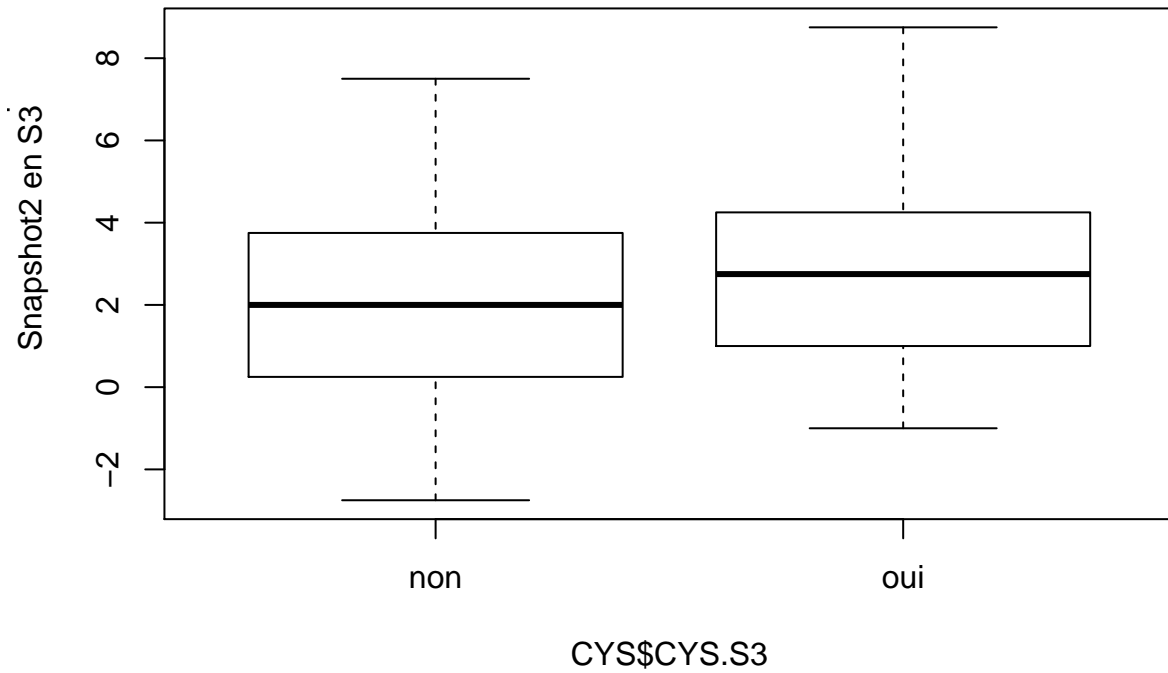
La p\_valeur associée au test Kolmogorov est inférieure à 0,05 donc on accepte que les étudiants ont des notes de Snapshot2 plus élevé que Snapshot1.

```
# Stat. descriptives à compléter
boxplot(CYS$snapshot.2~CYS$CYS.S3,names=c("Non utilisé","Utilisé"),
        ylab="Résultat Snapshot2 en S3",
        main=" Résultat Snapshot2 en S3 selon l'utilisation de CYS en S3")
```



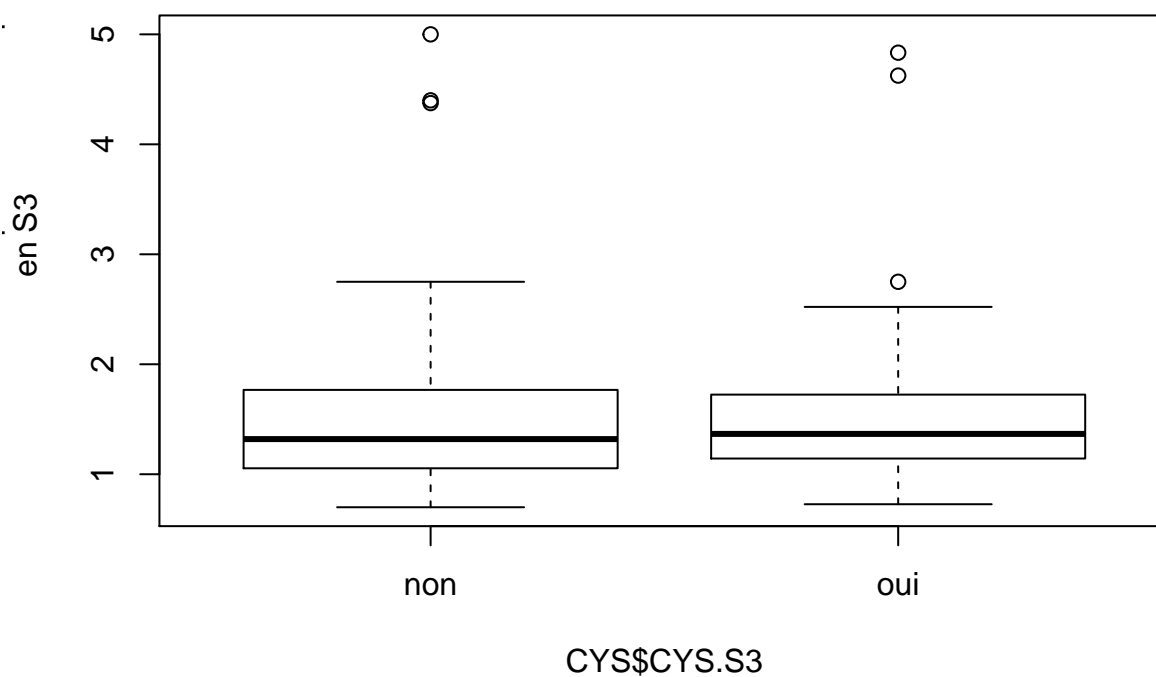
```
# Stat. descriptives à compléter
dif_snap=CYS$snapshot.2-CYS$snapshot.1
ratio_snap=CYS$snapshot.2/CYS$snapshot.1
boxplot(dif_snap~CYS$CYS.S3, ylab="Différence de résultat entre Snapshot1 et
        Snapshot2 en S3",
        main=" Différence de résultat entre Snapshot1 et
        Snapshot2 en S3 selon l'utilisation de CYS en S3")
```

### Différence de résultat entre Snapshot1 et Snapshot2 en S3 selon l'utilisation de CYS en S3



```
boxplot(ratio_snap~CYS$CYS.S3, ylab="Ratio de résultat entre Snapshot1 et Snapshot2  
en S3",  
main=" Ratio de résultat entre Snapshot1 et Snapshot2  
en S3 selon l'utilisation de CYS en S3")
```

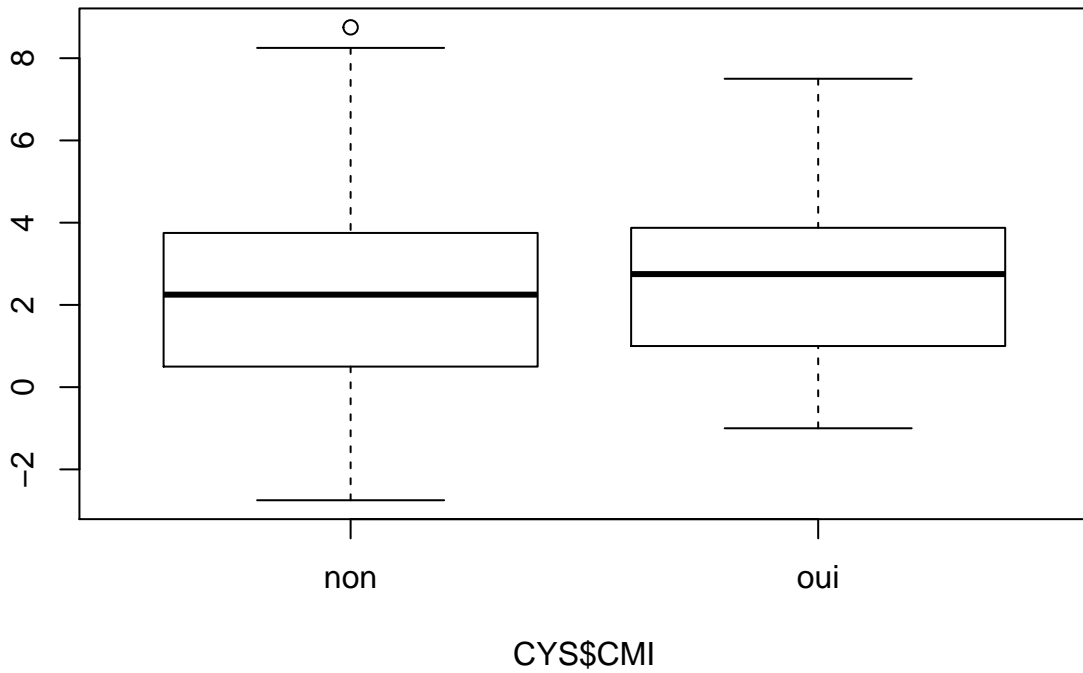
### Ratio de résultat entre Snapshot1 et Snapshot2 en S3 selon l'utilisation de CYS en S3



```
# Stat. descriptives à compléter  
boxplot(dif_snap~CYS$CMI, ylab="Différence de résultat entre Snapshot1 et Snapshot2 en S3",  
        main=" Différence de résultat entre Snapshot1 et Snapshot2 en S3 selon CMI en S3")
```

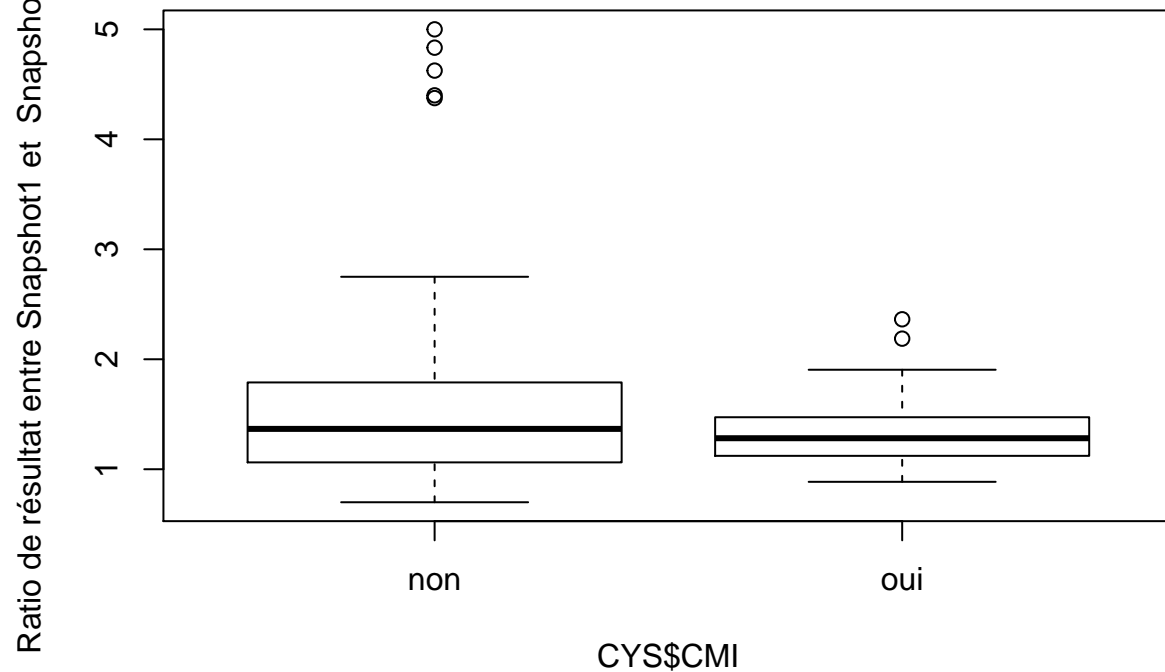
Différence de résultat entre Snapshot1 et Snapshot2 en S3

Différence de résultat entre Snapshot1 et Snapshot2 en S3 selon CMI



```
# Stat. descriptives à compléter
boxplot(ratio_snap~CYS$CMI, ylab="Ratio de résultat entre Snapshot1 et Snapshot2 en S3",
        main=" Ratio de résultat entre Snapshot1 et Snapshot2 en S3 selon CMI en S3")
```

## Ratio de résultat entre Snapshot1 et Snapshot2 en S3 selon CMI en



Remarque: On observe peu d'évolution entre les résultats de Snapshot1 et Snapshot2. Il vaut donc mieux considérer un modèle avec plusieurs variables.

## Test d'un modèle ANOVA de 3 facteurs(CYS S3, CMI, TP S3)

```
# A compléter
mod1=lm(dif_snap~(CYS$CYS.S3+CYS$TP.S3+ CYS$CMI)^2,data=CYS)
summary(mod1)
```

```
##
## Call:
## lm(formula = dif_snap ~ (CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI)^2,
##     data = CYS)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.875 -1.739  0.000   1.255   5.005
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      1.9954     0.2063   9.675  <2e-16 ***
## CYS$CYS.S3oui       0.7437     0.4941   1.505   0.134
## CYS$TP.S3GB        -0.9954     2.1632  -0.460   0.646
## CYS$CMIoui         3.1250     2.2584   1.384   0.168
## CYS$CYS.S3oui:CYS$TP.S3GB  2.4328     2.2702   1.072   0.285
```



```
## CYS$CYS.S3oui:CYS$CMIoui    -5.5515      2.3652   -2.347    0.020 *
## CYS$TP.S3GB:CYS$CMIoui      NA          NA      NA      NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.153 on 175 degrees of freedom
## Multiple R-squared:  0.1238, Adjusted R-squared:  0.09876
## F-statistic: 4.945 on 5 and 175 DF,  p-value: 0.0002929
# A completer - fonction lm
step.backward = step(mod1)
```

```
## Start:  AIC=283.56
## dif_snap ~ (CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI)^2
##
##
## Step:  AIC=283.56
## dif_snap ~ CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI + CYS$CYS.S3:CYS$TP.S3 +
##           CYS$CYS.S3:CYS$CMI
##
##           Df Sum of Sq    RSS    AIC
## - CYS$CYS.S3:CYS$TP.S3  1     5.3247 816.76 282.74
## <none>                        811.43 283.56
## - CYS$CYS.S3:CYS$CMI    1    25.5452 836.98 287.17
##
## Step:  AIC=282.74
## dif_snap ~ CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI + CYS$CYS.S3:CYS$CMI
##
##           Df Sum of Sq    RSS    AIC
## <none>                        816.76 282.74
## - CYS$TP.S3              1    15.852 832.61 284.22
## - CYS$CYS.S3:CYS$CMI    1    52.880 869.64 292.09
```

Selon le test d'AIC, on trouve le meilleur modèle modAIC1:

$$dif\_snap \sim CYS\$CYS.S3 + CYS\$TP.S3 + CYS\$CMI + CYS\$CYS.S3 : CYS\$CMI$$

```
# A completer
modAIC1=lm(dif_snap ~ CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI + CYS$CYS.S3:CYS$CMI,data=CYS)
summary(modAIC1)
```

```
##
## Call:
## lm(formula = dif_snap ~ CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI + CYS$CYS.S3:CYS$CMI,
##     data = CYS)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.875 -1.725  0.000  1.275  5.025
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      1.9753     0.2055   9.613 < 2e-16 ***
## CYS$CYS.S3oui      0.8590     0.4825   1.780 0.076740 .
## CYS$TP.S3GB       1.2134     0.6565   1.848 0.066252 .
```

```
## CYS$CMIoui          0.9362      0.9641   0.971 0.332834
## CYS$CYS.S3oui:CYS$CMIoui -3.2340      0.9580  -3.376 0.000906 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.154 on 176 degrees of freedom
## Multiple R-squared:  0.118, Adjusted R-squared:  0.098
## F-statistic: 5.889 on 4 and 176 DF,  p-value: 0.0001802
anova(modAIC1,mod1)

## Analysis of Variance Table
##
## Model 1: dif_snap ~ CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI + CYS$CYS.S3:CYS$CMI
## Model 2: dif_snap ~ (CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI)^2
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1      176 816.76
## 2      175 811.43   1    5.3247 1.1484 0.2854
```

La p\_valeur de Test Fisher est 0,2854 supérieure que 0,05 donc on accepte le modèle modAIC1.

```
step.backward = step(mod1,direction="backward",k=log(nrow(CYS)))
```

```
## Start:  AIC=302.75
## dif_snap ~ (CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI)^2
##
##
## Step:  AIC=302.75
## dif_snap ~ CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI + CYS$CYS.S3:CYS$TP.S3 +
##   CYS$CYS.S3:CYS$CMI
##
##           Df Sum of Sq    RSS    AIC
## - CYS$CYS.S3:CYS$TP.S3   1    5.3247 816.76 298.73
## <none>                      811.43 302.75
## - CYS$CYS.S3:CYS$CMI     1   25.5452 836.98 303.16
##
## Step:  AIC=298.73
## dif_snap ~ CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI + CYS$CYS.S3:CYS$CMI
##
##           Df Sum of Sq    RSS    AIC
## - CYS$TP.S3             1   15.852 832.61 297.01
## <none>                   816.76 298.73
## - CYS$CYS.S3:CYS$CMI    1   52.880 869.64 304.89
##
## Step:  AIC=297.01
## dif_snap ~ CYS$CYS.S3 + CYS$CMI + CYS$CYS.S3:CYS$CMI
##
##           Df Sum of Sq    RSS    AIC
## <none>                   832.61 297.01
## - CYS$CYS.S3:CYS$CMI    1   76.921 909.53 307.81
```

Selon le test d'BIC, on trouve le meilleur modèle:

*dif\_snap ~ CYS\$CYS.S3 + CYS\$CMI + CYS\$CYS.S3 : CYS\$CMI*

```
# A completer
modBIC1=lm(dif_snap ~ CYS$CYS.S3 + CYS$CMI + CYS$CYS.S3:CYS$CMI,data=CYS)
```

```
anova(modBIC1,mod1)
```

```
## Analysis of Variance Table
##
## Model 1: dif_snap ~ CYS$CYS.S3 + CYS$CMI + CYS$CYS.S3:CYS$CMI
## Model 2: dif_snap ~ (CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI)^2
##   Res.Df    RSS Df Sum of Sq      F Pr(>F)
## 1      177 832.61
## 2      175 811.43  2    21.177 2.2835  0.105
```

La p\_valeur de Test Fisher est 0,105 supérieure que 0,05 donc on accepte le modèle modBIC1.

```
# A completer
```

```
anova(modBIC1,modAIC1)
```

```
## Analysis of Variance Table
##
## Model 1: dif_snap ~ CYS$CYS.S3 + CYS$CMI + CYS$CYS.S3:CYS$CMI
## Model 2: dif_snap ~ CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI + CYS$CYS.S3:CYS$CMI
##   Res.Df    RSS Df Sum of Sq      F Pr(>F)
## 1      177 832.61
## 2      176 816.76  1    15.852 3.4158 0.06625 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

La p\_valeur de Test Fisher est 0,105 supérieure que 0,05 donc on accepte le modèle modBIC1.

```
# A completer
```

```
mod2=lm(ratio_snap~(CYS$CYS.S3+CYS$TP.S3+ CYS$CMI)^2,data=CYS)
summary(mod1)
```

```
##
## Call:
## lm(formula = dif_snap ~ (CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI)^2,
##     data = CYS)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.875 -1.739  0.000  1.255  5.005
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      1.9954     0.2063   9.675  <2e-16 ***
## CYS$CYS.S3oui      0.7437     0.4941   1.505   0.134
## CYS$TP.S3GB      -0.9954     2.1632  -0.460   0.646
## CYS$CMIoui       3.1250     2.2584   1.384   0.168
## CYS$CYS.S3oui:CYS$TP.S3GB  2.4328     2.2702   1.072   0.285
## CYS$CYS.S3oui:CYS$CMIoui -5.5515     2.3652  -2.347   0.020 *
## CYS$TP.S3GB:CYS$CMIoui      NA         NA      NA      NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.153 on 175 degrees of freedom
## Multiple R-squared:  0.1238, Adjusted R-squared:  0.09876
## F-statistic: 4.945 on 5 and 175 DF, p-value: 0.0002929
```

```
step.backward = step(mod2,direction="backward",k=log(nrow(CYS)))
```

```
## Start:  AIC=-101.45
## ratio_snap ~ (CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI)^2
##
##
## Step:  AIC=-101.45
## ratio_snap ~ CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI + CYS$CYS.S3:CYS$TP.S3 +
##      CYS$CYS.S3:CYS$CMI
##
##              Df Sum of Sq    RSS      AIC
## - CYS$CYS.S3:CYS$TP.S3  1   0.05397 87.032 -106.54
## - CYS$CYS.S3:CYS$CMI    1   0.52871 87.507 -105.56
## <none>                                86.978 -101.45
##
## Step:  AIC=-106.54
## ratio_snap ~ CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI + CYS$CYS.S3:CYS$CMI
##
##              Df Sum of Sq    RSS      AIC
## - CYS$TP.S3          1   0.18258 87.214 -111.36
## - CYS$CYS.S3:CYS$CMI  1   1.61594 88.648 -108.41
## <none>                                87.032 -106.54
##
## Step:  AIC=-111.36
## ratio_snap ~ CYS$CYS.S3 + CYS$CMI + CYS$CYS.S3:CYS$CMI
##
##              Df Sum of Sq    RSS      AIC
## - CYS$CYS.S3:CYS$CMI  1     1.438 88.652 -113.60
## <none>                                87.214 -111.36
##
## Step:  AIC=-113.6
## ratio_snap ~ CYS$CYS.S3 + CYS$CMI
##
##              Df Sum of Sq    RSS      AIC
## - CYS$CYS.S3  1   0.72326 89.376 -117.33
## - CYS$CMI     1   1.57982 90.232 -115.60
## <none>                                88.652 -113.60
##
## Step:  AIC=-117.33
## ratio_snap ~ CYS$CMI
##
##              Df Sum of Sq    RSS      AIC
## - CYS$CMI  1   1.0724 90.448 -120.36
## <none>                                89.376 -117.33
##
## Step:  AIC=-120.37
## ratio_snap ~ 1
```

Selon le test d'BIC, on trouve le meilleur modèle modBIC2:

$ratio\_snap \sim CYS\$CMI$

De même façon, on trouve le meilleur modèle pour modéliser le ratio\_snap:

```

# A completer
modBIC2=lm(ratio_snap ~ CYS$CMI,data=CYS)
summary(modBIC2)

##
## Call:
## lm(formula = ratio_snap ~ CYS$CMI, data = CYS)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.8666 -0.4505 -0.1435  0.1981  3.4334
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.56657    0.05769  27.153  <2e-16 ***
## CYS$CMIoui   -0.20431    0.13941  -1.466   0.145
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7066 on 179 degrees of freedom
## Multiple R-squared:  0.01186,    Adjusted R-squared:  0.006336
## F-statistic: 2.148 on 1 and 179 DF,  p-value: 0.1445
anova(modBIC2,mod2)

```

```

## Analysis of Variance Table
##
## Model 1: ratio_snap ~ CYS$CMI
## Model 2: ratio_snap ~ (CYS$CYS.S3 + CYS$TP.S3 + CYS$CMI)^2
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1      179 89.376
## 2      175 86.978  4     2.3978 1.2061  0.31

```

La p\_valeur de Test Fisher est 0,31 supérieure que 0,05 donc on accepte le modèle modBIC1.

Remarque: Grâce au test ANOVA on trouve que:

- Si on considère la différence entre les deux snapshots on obtient le modèle

$$dif\_snap \sim CYS\$CYS.S3 + CYS\$CMI + CYS\$CYS.S3 : CYS\$CMI$$

Cela montre l'impact de CMI et CYS S3 sur l'évolution de résultat.

- SI on considère la ratio entre les deux snapshots on obtient le modèle

$$ratio\_snap \sim CYS\$CMI$$

Cela montre l'impact de CMI S3 sur l'évolution de résultat.

Dans ces deux cas on ne trouve pas l'effet de la variable TP S3.

## Test d'un modèle ANCOVA de 3 facteurs qualitatives (CYS S3, CMI, TP S3) et 1 facteur quantitative(Snapshot1) pour modéliser le Snapshot2

On souhaite expliquer le Snapshot2 en fonction du Snapshot1, des choix (de la langue de TP, l'utilisation de CMI et celle de CYS) en S3. On met en place un modèle d'analyse de la covariance.

```

# A completer
modR=lm(CYS$snapshot.2 ~ (CYS$snapshot.1+CYS$TP.S3+CYS$CYS.S3+ CYS$CMI)^2,data=CYS)
summary(modR)

##
## Call:
## lm(formula = CYS$snapshot.2 ~ (CYS$snapshot.1 + CYS$TP.S3 + CYS$CYS.S3 +
##     CYS$CMI)^2, data = CYS)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.9027 -1.3388 -0.0402  1.1452  5.7622
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      4.52851    0.50089   9.041 3.29e-16 ***
## CYS$snapshot.1    0.58867    0.07554   7.792 6.10e-13 ***
## CYS$TP.S3GB       0.36310    2.94340   0.123  0.90197
## CYS$CYS.S3oui     0.21415    1.12806   0.190  0.84966
## CYS$CMIoui        3.96204    3.15382   1.256  0.21073
## CYS$snapshot.1:CYS$TP.S3GB -0.18738    0.33981  -0.551  0.58207
## CYS$snapshot.1:CYS$CYS.S3oui  0.01915    0.19788   0.097  0.92300
## CYS$snapshot.1:CYS$CMIoui   0.04565    0.34210   0.133  0.89400
## CYS$TP.S3GB:CYS$CYS.S3oui   3.03385    2.06231   1.471  0.14310
## CYS$TP.S3GB:CYS$CMIoui      NA         NA      NA      NA
## CYS$CYS.S3oui:CYS$CMIoui   -6.05498    2.16507  -2.797  0.00575 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.938 on 171 degrees of freedom
## Multiple R-squared:  0.4853, Adjusted R-squared:  0.4582
## F-statistic: 17.91 on 9 and 171 DF,  p-value: < 2.2e-16

# A completer
step.backward = step(modR,direction="backward",k=log(nrow(CYS)))

## Start:  AIC=281.22
## CYS$snapshot.2 ~ (CYS$snapshot.1 + CYS$TP.S3 + CYS$CYS.S3 + CYS$CMI)^2
##
##
## Step:  AIC=281.22
## CYS$snapshot.2 ~ CYS$snapshot.1 + CYS$TP.S3 + CYS$CYS.S3 + CYS$CMI +
##     CYS$snapshot.1:CYS$TP.S3 + CYS$snapshot.1:CYS$CYS.S3 + CYS$snapshot.1:CYS$CMI +
##     CYS$TP.S3:CYS$CYS.S3 + CYS$CYS.S3:CYS$CMI
##
##              Df Sum of Sq    RSS    AIC
## - CYS$snapshot.1:CYS$CYS.S3  1    0.0352 642.30 276.03
## - CYS$snapshot.1:CYS$CMI      1    0.0669 642.33 276.04
## - CYS$snapshot.1:CYS$TP.S3    1    1.1420 643.41 276.35
## - CYS$TP.S3:CYS$CYS.S3        1    8.1283 650.39 278.30
## <none>                        642.27 281.22
## - CYS$CYS.S3:CYS$CMI          1   29.3765 671.64 284.12
##
## Step:  AIC=276.03
## CYS$snapshot.2 ~ CYS$snapshot.1 + CYS$TP.S3 + CYS$CYS.S3 + CYS$CMI +

```

```

##      CYS$snapshot.1:CYS$TP.S3 + CYS$snapshot.1:CYS$CMI + CYS$TP.S3:CYS$CYS.S3 +
##      CYS$CYS.S3:CYS$CMI
##
##              Df Sum of Sq    RSS    AIC
## - CYS$snapshot.1:CYS$CMI      1    0.0514 642.35 270.85
## - CYS$snapshot.1:CYS$TP.S3    1    1.2895 643.59 271.20
## - CYS$TP.S3:CYS$CYS.S3        1    8.3852 650.69 273.18
## <none>                        642.30 276.03
## - CYS$CYS.S3:CYS$CMI          1   29.9130 672.21 279.07
##
## Step:  AIC=270.85
## CYS$snapshot.2 ~ CYS$snapshot.1 + CYS$TP.S3 + CYS$CYS.S3 + CYS$CMI +
##      CYS$snapshot.1:CYS$TP.S3 + CYS$TP.S3:CYS$CYS.S3 + CYS$CYS.S3:CYS$CMI
##
##              Df Sum of Sq    RSS    AIC
## - CYS$snapshot.1:CYS$TP.S3    1    2.6647 645.02 266.40
## - CYS$TP.S3:CYS$CYS.S3        1    8.3467 650.70 267.99
## <none>                        642.35 270.85
## - CYS$CYS.S3:CYS$CMI          1   29.8840 672.24 273.88
##
## Step:  AIC=266.4
## CYS$snapshot.2 ~ CYS$snapshot.1 + CYS$TP.S3 + CYS$CYS.S3 + CYS$CMI +
##      CYS$TP.S3:CYS$CYS.S3 + CYS$CYS.S3:CYS$CMI
##
##              Df Sum of Sq    RSS    AIC
## - CYS$TP.S3:CYS$CYS.S3      1     8.280 653.30 263.51
## <none>                        645.02 266.40
## - CYS$CYS.S3:CYS$CMI        1    29.003 674.02 269.16
## - CYS$snapshot.1            1   280.332 925.35 326.52
##
## Step:  AIC=263.51
## CYS$snapshot.2 ~ CYS$snapshot.1 + CYS$TP.S3 + CYS$CYS.S3 + CYS$CMI +
##      CYS$CYS.S3:CYS$CMI
##
##              Df Sum of Sq    RSS    AIC
## <none>                        653.30 263.51
## - CYS$TP.S3                  1    37.826 691.12 268.50
## - CYS$CYS.S3:CYS$CMI        1    46.252 699.55 270.69
## - CYS$snapshot.1            1   285.184 938.48 323.88
##
## # A completer
modbest=lm(CYS$snapshot.2 ~ CYS$snapshot.1 + CYS$TP.S3 + CYS$CYS.S3 + CYS$CMI +
      CYS$CYS.S3:CYS$CMI,data=CYS)
summary(modbest)

##
## Call:
## lm(formula = CYS$snapshot.2 ~ CYS$snapshot.1 + CYS$TP.S3 + CYS$CYS.S3 +
##      CYS$CMI + CYS$CYS.S3:CYS$CMI, data = CYS)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.8965 -1.2725 -0.0205  1.1728  5.8232
##
## Coefficients:

```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      4.62384    0.44064  10.493 < 2e-16 ***
## CYS$snapshot.1    0.56912    0.06511   8.740 1.85e-15 ***
## CYS$TP.S3GB       1.90361    0.59802   3.183 0.00172 **
## CYS$CYS.S3oui     0.43506    0.43744   0.995 0.32132
## CYS$CMIoui        1.28154    0.86629   1.479 0.14084
## CYS$CYS.S3oui:CYS$CMIoui -3.02652    0.85983  -3.520 0.00055 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.932 on 175 degrees of freedom
## Multiple R-squared:  0.4764, Adjusted R-squared:  0.4615
## F-statistic: 31.85 on 5 and 175 DF, p-value: < 2.2e-16
```

```
anova(modbest,modR)
```

```
## Analysis of Variance Table
##
## Model 1: CYS$snapshot.2 ~ CYS$snapshot.1 + CYS$TP.S3 + CYS$CYS.S3 + CYS$CMI +
##          CYS$CYS.S3:CYS$CMI
## Model 2: CYS$snapshot.2 ~ (CYS$snapshot.1 + CYS$TP.S3 + CYS$CYS.S3 + CYS$CMI)^2
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1      175 653.30
## 2      171 642.27  4    11.031 0.7343 0.5698
```

Selon le test d'BIC, on trouve le meilleur modèle modBIC2:

$CYS\$snapshot.2 \sim CYS\$snapshot.1 + CYS\$TP.S3 + CYS\$CYS.S3 + CYS\$CMI + CYS\$CYS.S3 : CYS\$CMI$

## Train+ Test par le modèle ANCOVA

### Train

```
n <- nrow(CYS)
trainIndices <- sample(n, size = 2 * n / 3)
train <- CYS[trainIndices, ]
test <- CYS[-trainIndices, ]
```

### Test

```
m0=lm(snapshot.2 ~ snapshot.1 + TP.S3 + CYS.S3 + CMI +
      CYS.S3:CMI,data=train)
p <- predict(m0, newdata = test, type = "response")
table(p)
```

```
## p
## 5.91388375547579 6.16982615982907 6.19377473858344 6.44971714293672
##           1           1           1           1
## 6.75355670479873 6.86955361759819 6.89350219635255 7.03344768790638
##           3           1           2           2
## 7.1733931794602 7.31333867101402 7.59322965412167 7.73317514567549
##           1           2           3           1
## 7.84917205847495 7.87312063722931 7.98911755002877 8.01306612878314
##           1           6           1           1
## 8.26900853313642 8.29295711189078 8.43290260344461 8.54889951624406
```



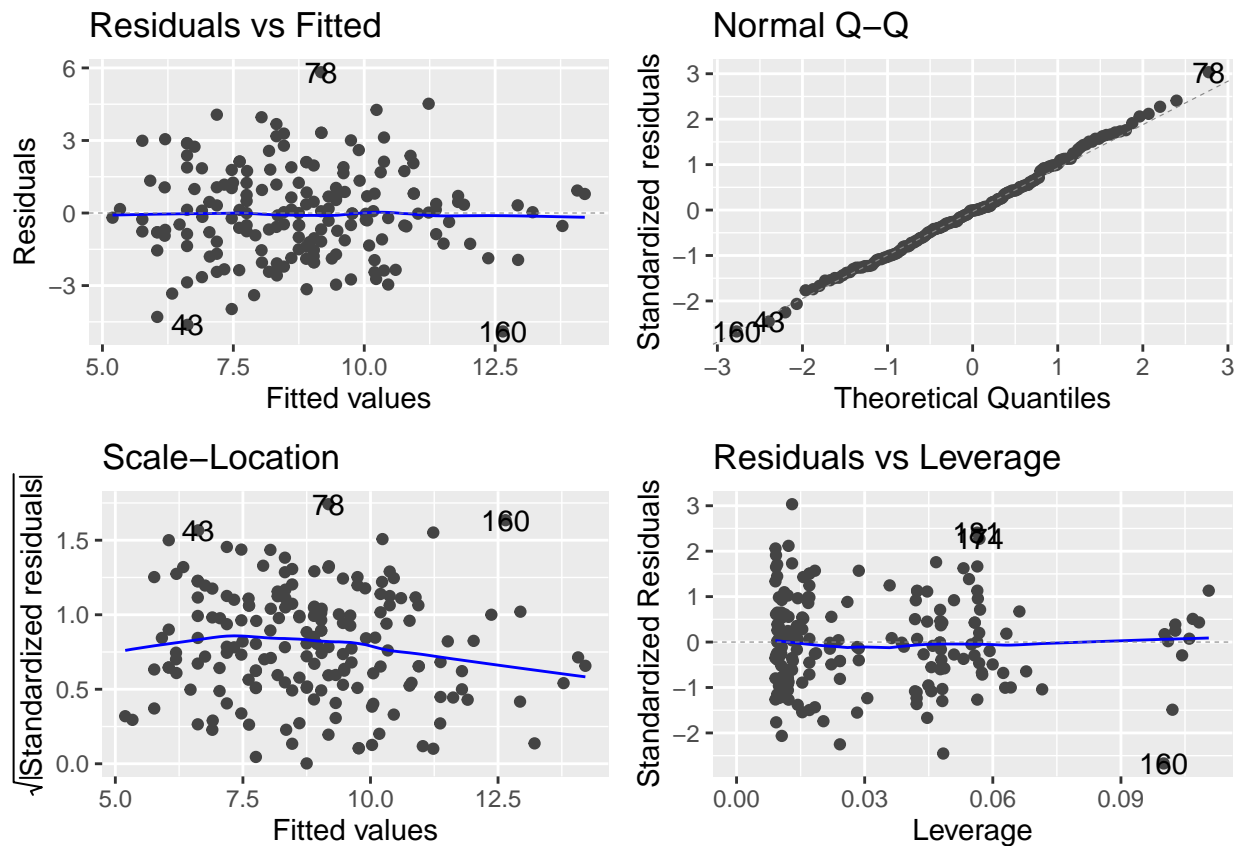
```
##           1           1           1           1
## 8.57284809499843 8.83047080849156 8.85273907810608 8.9926845696599
##           3           1           2           1
## 9.13263006121372 9.24862697401318 9.27257555276754 9.41252104432137
##           1           1           2           2
## 9.62286492977803 9.69241202742901 9.76281042133185 9.90275591288567
##           1           1           1           1
## 10.0427014044395 10.1122485020905 10.1826468959933 10.3225923875471
##           1           1           1           1
## 10.342393445949 10.7424288622086 11.0695986733527 11.182066395272
##           1           1           1           1
## 11.2095441649066 12.071485383844 12.3016303277026 13.4211942601331
##           1           1           1           1
```

```
#p0<-(CYS$snapshot.2)[-trainIndices]
p0<-test$snapshot.2
SCR=norm(p-p0,type="2")^2
SCT=norm(p0-mean(p0),type="2")^2
Rajustee<-1-SCR/SCT
Rajustee
```

```
## [1] 0.5245512
```

```
#norm(p,type="2")
```

```
autoplot(modbest)
```



Vérifions les 4 hypothèses pour ce modèle:

Hypothèse H1 : Les erreurs sont centrées

Hypothèse H2 : La variance des erreurs est constante

Hypothèse H3 : Les résidus sont indépendantes

Hypothèse H4 : Les données suivent des lois gaussiennes

- Les résidus sont situés aléatoirement sur les deux côtés de 0. Or, on ne trouve pas la forme bannane ni la forme trompette, donc on peut déduire que H1 et H2 sont vérifiées
- On observe dans le graphe Residuals vs Fitted aucun phénomène de paquets sur les côtés de 0 donc H3 est vérifiée.
- Au regard du graphe Normal Q-Q, on constate que les résidus suivent une loi gaussienne. Ainsi, H4 est vérifiée.

## Conclusion:

On a décidé de modéliser la note de Snapshot2 en fonction de Snapshot1, l'utilisation de l'outil CheckYourSmile, la langue de TP et le fait que l'étudiant est en CMI ou pas.

Alors, sous le modèle ANCOVA on a trouvé que les trois facteurs qualitatives ont des impacts sur le résultat Snapshot2. Or, le Snapshot1 a un gros effet sur le résultat de Snapshot2. On y trouve aussi un terme d'interaction entre la variable CYS.S3 et la variable CMI et celui dernier a un effet important négatif sur le Snapshot2.

ie, un étudiant en CMI utilisant l'outil Check Your Smile a tendance de dégrader environ 2,6 points (-3,02652+0,43506) et un étudiant non CMI utilisant l'outil Check Your Smile a tendance de progresser environ 0,4 points (0,43506)

On voit que le modèle ANCOVA nous donne un R-ajusté bien meilleur que le modèle ANOVA. (0,4615»0,006 et 0,4615»0,09)

Pourtant, l'analyse des données peut s'améliorer car on n'a gardé que 181 données en semestre 3 et 9 en semestre 4 sur 459 étudiants.

S3	Snapshot1	Snapshot2	CYS S3	TP S3	CMI
Manquant	38	27	1	23	21
Total	242				
Restant	181				
Nb négligé	61				
S4	Snapshot1	Snapshot2	CYS S4	TP S4	CMI
Manquant	126	51	28	86	91
Total	217				
Restant	9				
Nb négligé	208				

Figure 1: Données Manquantes

En semestre 4, les données ne sont pas assez bonnes:

- 126 individus n'ont pas de note de Snapshot1
- 51 individus n'ont pas de note de Snapshot2
- 28 individus ont des réponses invalidées si ils utilisent ou pas l'outil CYS (la valeur de réponse est soit BLANK soit ?)

- 86 individus ont des réponses invalidées de la langue pour TP (la valeur de réponse est soit non soit BLANK)
- 91 individus ont des réponses invalidées si ils sont CMI (la valeur de réponse est soit BLANK soit ? soit N/A)