

Discovery__dataset

Ariel-ac4391

11/22/2018

Here, I recapitulate the main step related in the research paper with the graphs associated

The first step is data cleansing :

```
training_data=read.csv("data/Data_User_Modeling_training_Dataset.csv")
test_data=read.csv("data/Data_User_Modeling_test_Dataset.csv")
library(gplots)
```

```
##
## Attaching package: 'gplots'
## The following object is masked from 'package:stats':
##
##      lowess
```

```
library(ggplot2)
library(partykit)
```

```
## Loading required package: grid
## Loading required package: libcoin
## Loading required package: mvtnorm
library(rpart) # Popular decision tree algorithm
library(hier.part)
```

```
## Loading required package: gtools
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
## The following objects are masked from 'package:stats':
##
##      filter, lag
## The following objects are masked from 'package:base':
##
##      intersect, setdiff, setequal, union
```

```
library(ipred)
library(randomForest)
```

```
## randomForest 4.6-14
## Type rfNews() to see new features/changes/bug fixes.
##
## Attaching package: 'randomForest'
## The following object is masked from 'package:dplyr':
##
##      combine
```

```
## The following object is masked from 'package:ggplot2':
##
##     margin
#library(rattle) # GUI for building trees and fancy tree plot #Doesn't work
library(rpart.plot) # Enhanced tree plots
library(party) # Alternative decision tree algorithm

## Loading required package: modeltools
## Loading required package: stats4
## Loading required package: strucchange
## Loading required package: zoo
##
## Attaching package: 'zoo'
## The following objects are masked from 'package:base':
##
##     as.Date, as.Date.numeric
## Loading required package: sandwich
##
## Attaching package: 'party'
## The following objects are masked from 'package:partykit':
##
##     cforest, ctree, ctree_control, edge_simple, mob, mob_control,
##     node_barplot, node_bivplot, node_boxplot, node_inner,
##     node_surv, node_terminal, varimp
library(partykit) # Convert rpart object to BinaryTree
#library(RWeka) # Weka decision tree J48.
library(C50) # Original C5.0 implementation.
summary(training_data)

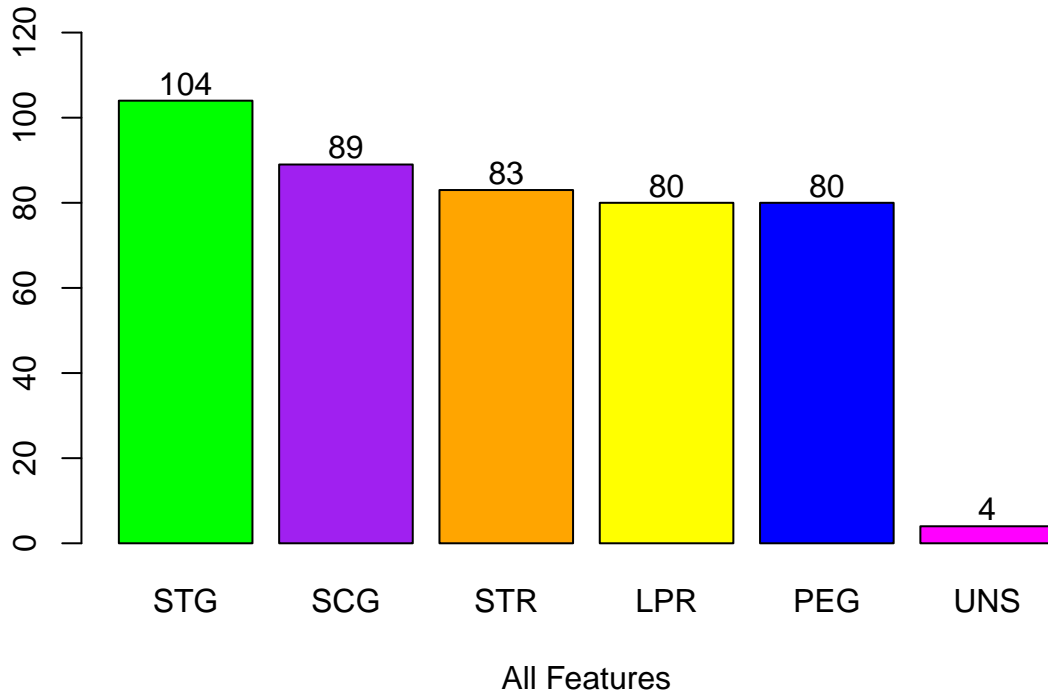
##           STG           SCG           STR           LPR
## Min.      :0.0000   Min.      :0.0000   Min.      :0.0000   Min.      :0.0000
## 1st Qu.:0.2407   1st Qu.:0.2100   1st Qu.:0.2913   1st Qu.:0.2500
## Median :0.3270   Median :0.3025   Median :0.4900   Median :0.3300
## Mean     :0.3711   Mean     :0.3557   Mean     :0.4680   Mean     :0.4327
## 3rd Qu.:0.4950   3rd Qu.:0.4975   3rd Qu.:0.6900   3rd Qu.:0.6475
## Max.     :0.9900   Max.     :0.9000   Max.     :0.9500   Max.     :0.9900
##           PEG           UNS
## Min.      :0.0000   High      :63
## 1st Qu.:0.2500   Low       :83
## Median :0.5000   Middle    :88
## Mean     :0.4585   very_low:24
## 3rd Qu.:0.6600
## Max.     :0.9300

attach(training_data)

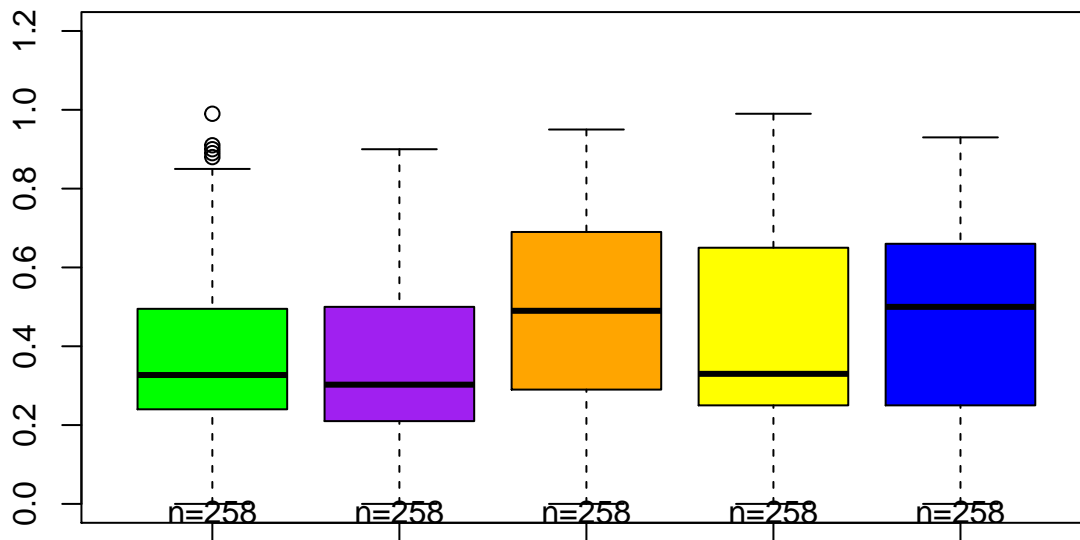
# Number of distinct values in each feture
a = n_distinct(STG)
b = n_distinct(SCG)
c = n_distinct(STR)
```

```
d = n_distinct(LPR)
e = n_distinct(PEG)
f = n_distinct(UNS)
num_distinct = c(a,b,c,d,e,f)
```

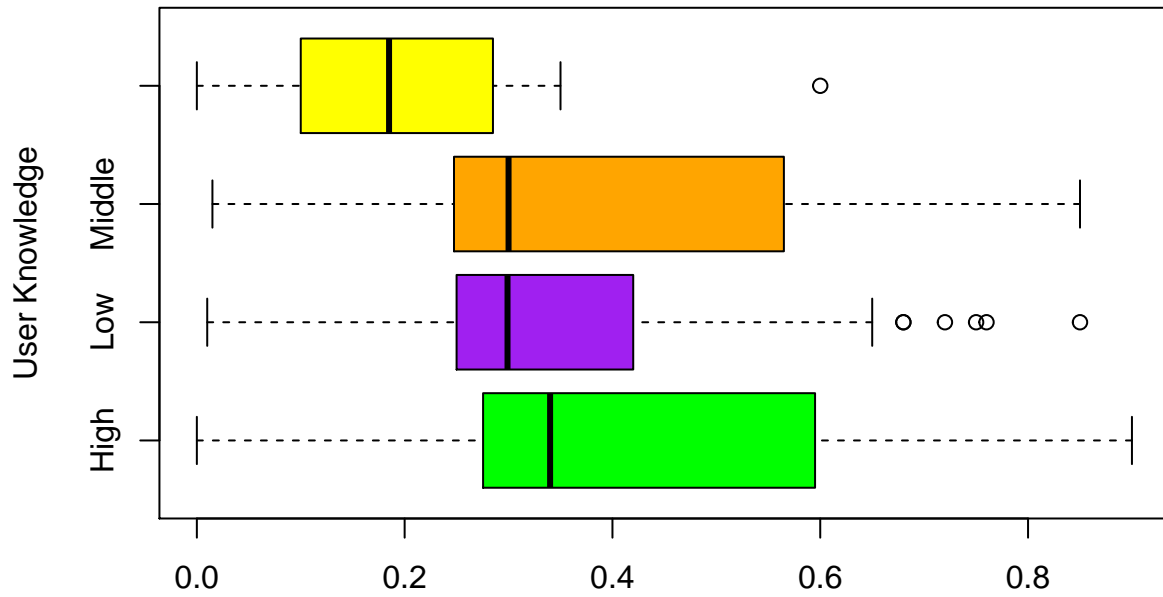
```
plot = barplot(num_distinct, names = c("STG", "SCG", "STR", "LPR", "PEG", "UNS"), ylim=c(0,120), xlab="All Features",
text(plot,num_distinct + 4,labels=as.character(num_distinct))
```



```
# boxplot of all data
boxplot2(STG,SCG,STR,LPR,PEG, col=c("green", "purple", "orange", "yellow", "blue", "magenta"), ylim=c(0,1.2))
```

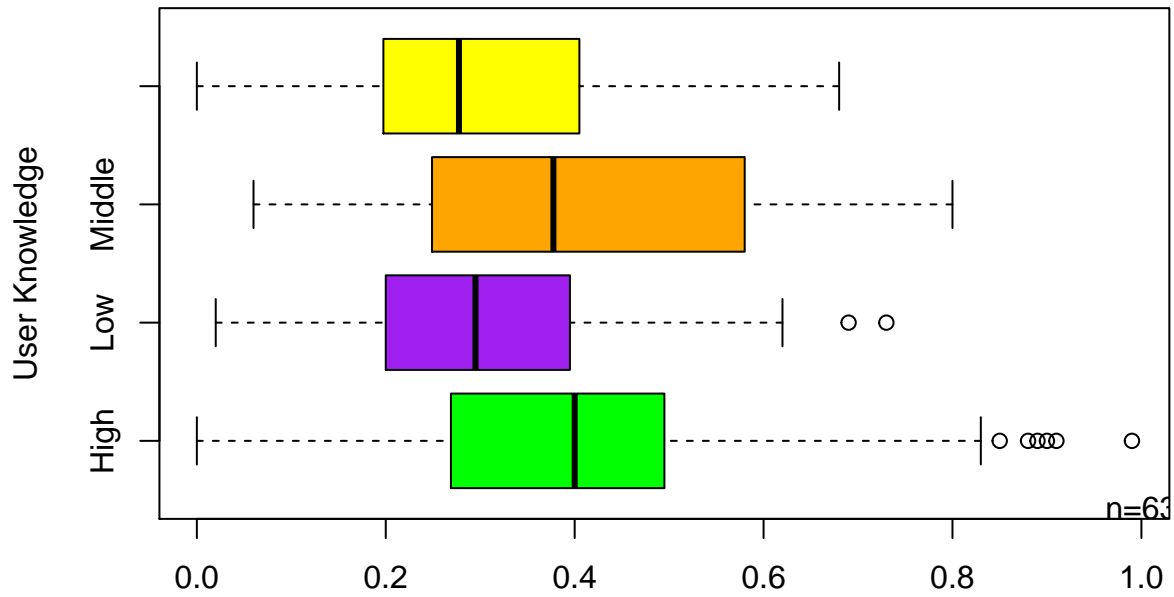


```
# boxplot of SCG divided across UNS values
boxplot2(SCG~UNS,data=training_data, horizontal = TRUE,
xlab="The degree of study time for goal materials", ylab="User Knowledge", col=c("green", "purple", "orange", "yellow", "blue", "magenta"))
```



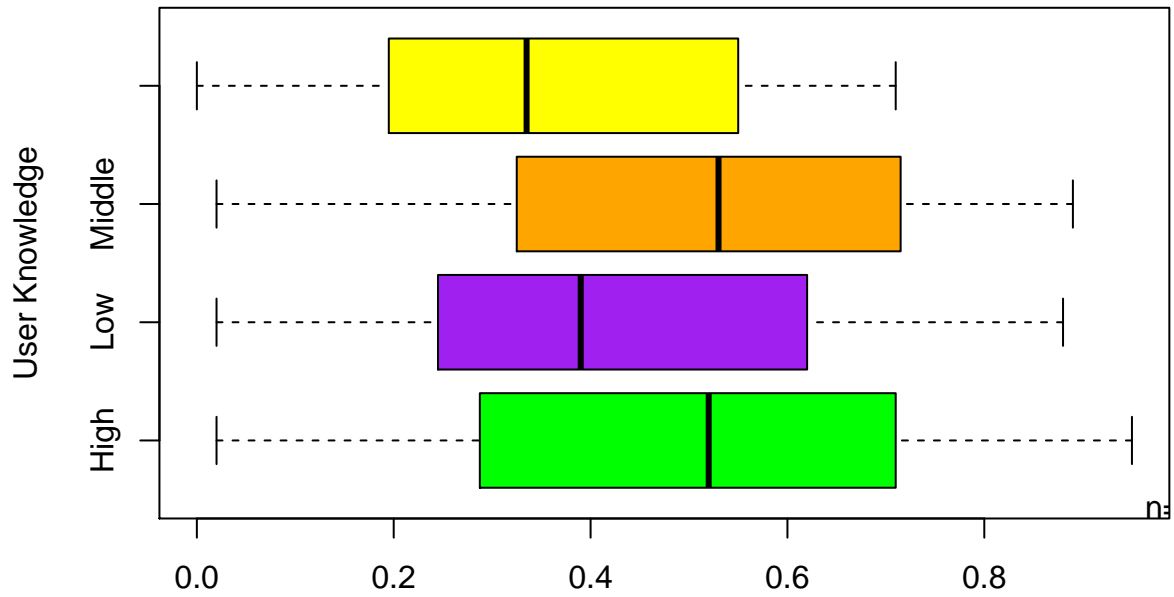
The degree of study time for goal materials

```
# boxplot of STG divided accross UNS values
boxplot2(STG~UNS,data=training_data, horizontal = TRUE,
  xlab="The degree of repetition number of user for goal materials", ylab="User Knowledge", col=c("g
```



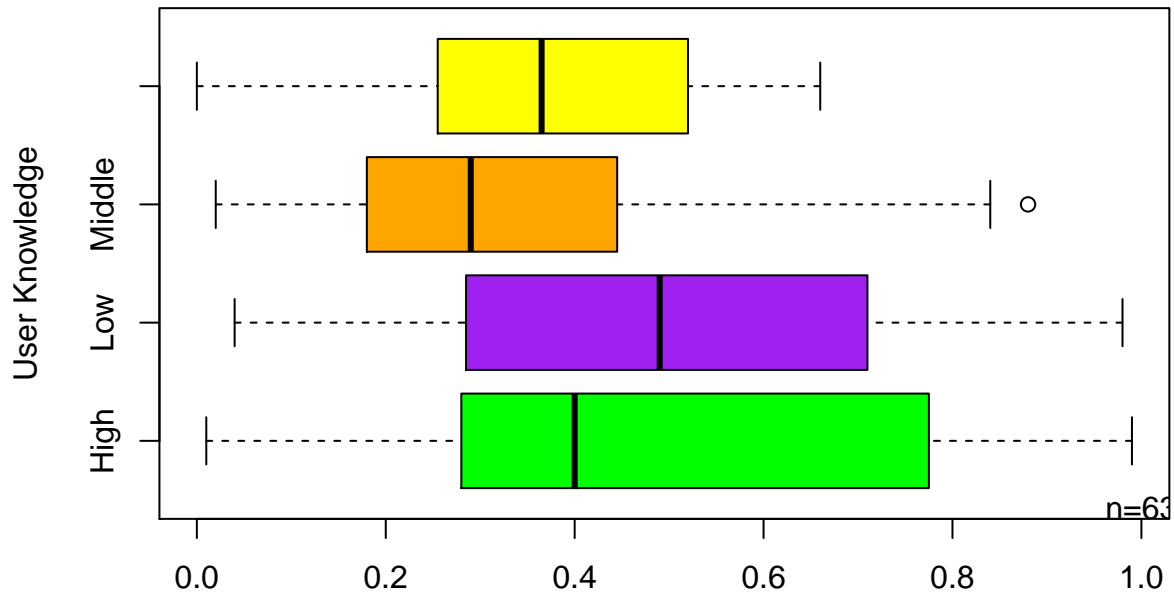
The degree of repetition number of user for goal materials

```
# boxplot of STR divided accross UNS values
boxplot2(STR~UNS,data=training_data, horizontal = TRUE,
  xlab="The degree of study time for related objects with goal materials", ylab="User Knowledge", co
```



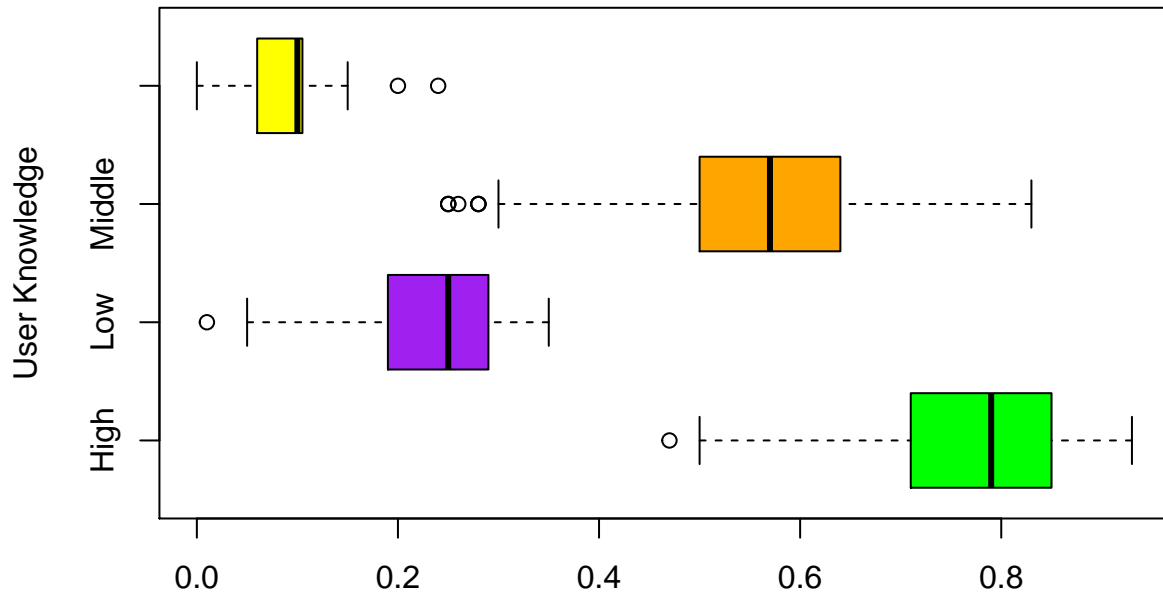
The degree of study time for related objects with goal materials

```
# boxplot of LPR divided accross UNS values
boxplot2(LPR~UNS,data=training_data, horizontal = TRUE,
  xlab="The exam performance of user for related objects with goal materials", ylab="User Knowledge")
```



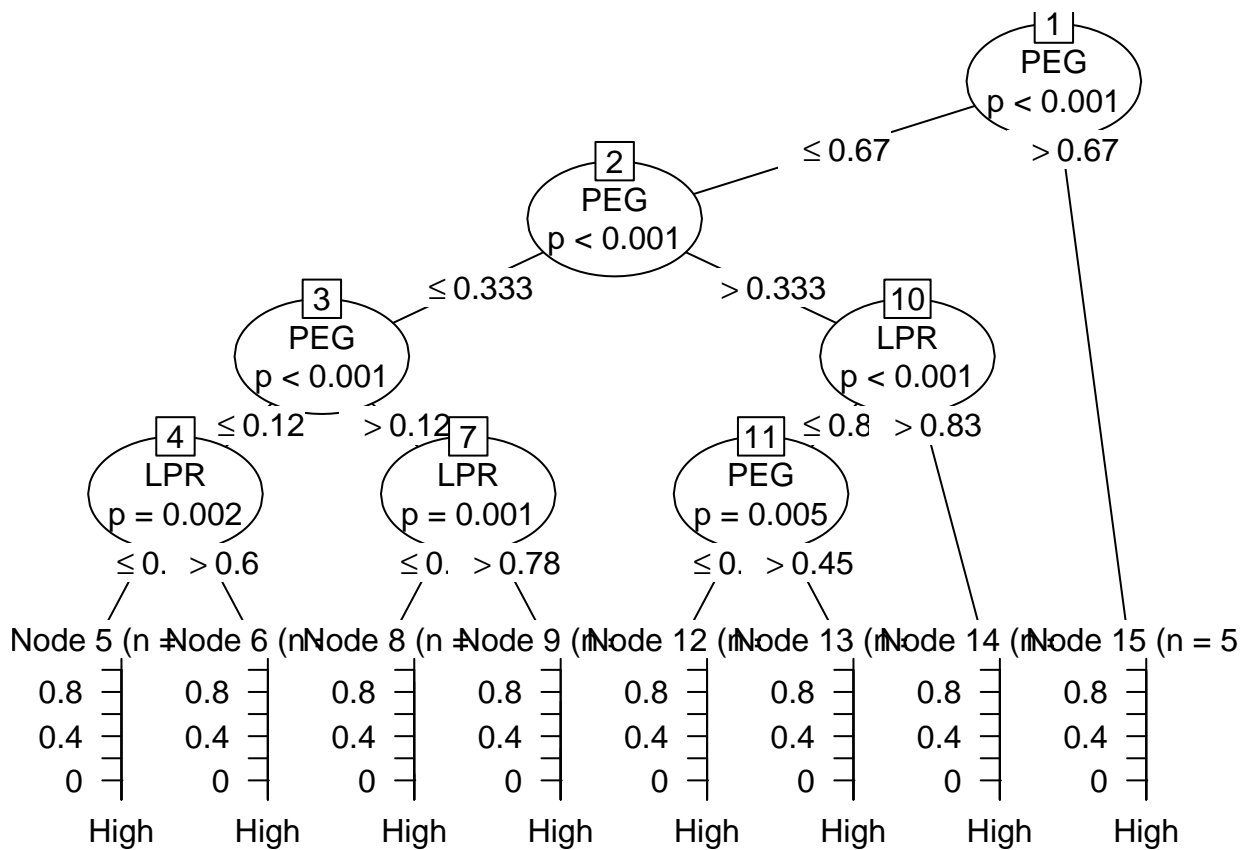
The exam performance of user for related objects with goal materials

```
# boxplot of PEG divided accross UNS values
boxplot2(PEG~UNS,data=training_data, horizontal = TRUE,
  xlab="The exam performance of user for goal materials", ylab="User Knowledge", col=c("green", "purple", "orange", "yellow"))
```



The exam performance of user for goal materials

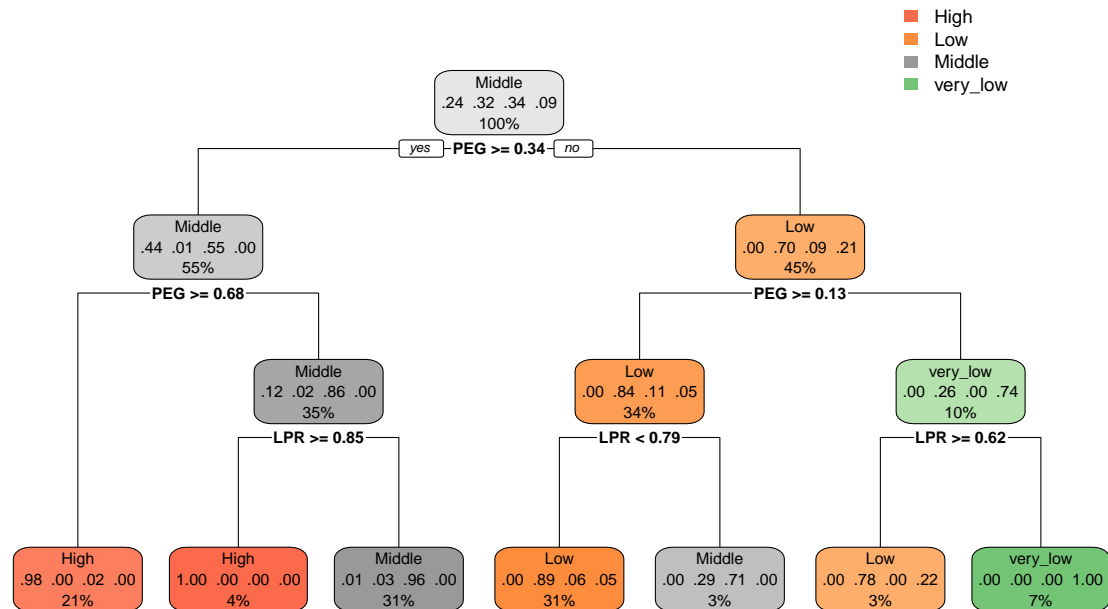
```
# decision tree
tree1 <- ctree(UNS ~ ., data = training_data)
plot(tree1) #Review the design
```



```
fit1 = predict(tree1, test_data)
table(fit1, test_data$UNS)
```

```
##
## fit1      High Low Middle Very Low
## High      39  0      1      0
## Low       0 42      3      5
## Middle    0  4     30      0
## very_low  0  0      0     21

#recursive partition tree
tree2 <- rpart(UNS ~ ., data = training_data)
rpart.plot(tree2)
```



```
rpart.rules(tree2)
```

```
##      UNS  High  Low Midd very
##      High [.98 .00 .02 .00] when PEG >=          0.68
##      High [1.00 .00 .00 .00] when PEG is 0.34 to 0.68 & LPR >= 0.85
##      Low  [.00 .78 .00 .22] when PEG < 0.13          & LPR >= 0.62
##      Low  [.00 .89 .06 .05] when PEG is 0.13 to 0.34 & LPR < 0.79
##      Middle [.00 .29 .71 .00] when PEG is 0.13 to 0.34 & LPR >= 0.79
##      Middle [.01 .03 .96 .00] when PEG is 0.34 to 0.68 & LPR < 0.85
##      very_low [.00 .00 .00 1.00] when PEG < 0.13          & LPR < 0.62
```

```
fit2 = predict(tree2, test_data, type = "class")
table(fit2, test_data$UNS)
```

```
##
## fit2      High Low Middle Very Low
## High      39  0      1      0
## Low       0 42      3      5
## Middle    0  4     30      0
## very_low  0  0      0     21
```

```
# J48 package issues
```

```
# PART package issues
```

```
#hier.part(training_data$UNS, training_data)
```

Bagging tree NOTE: Interesting we did much better than them here, they did something wrong

```
tree3 = bagging(UNS ~., data=training_data, coob=TRUE)
fit3 = predict(tree3, test_data)
table(fit3, test_data$UNS)
```

```
##
## fit3      High Low Middle Very Low
##   High      35  0      1      0
##   Low       0 43      3      4
##   Middle     4  3     30      0
##   very_low   0  0      0     22
```

Random Forest

```
tree4 = randomForest(UNS ~., data=training_data)
fit4 = predict(tree4, test_data)
table(fit4, test_data$UNS)
```

```
##
## fit4      High Low Middle Very Low
##   High      39  0      0      0
##   Low       0 45      3      3
##   Middle     0  1     31      0
##   very_low   0  0      0     23
```

C5.0

```
tree5 <- C5.0(UNS ~., data=training_data)
fit5 = predict(tree5, test_data)
table(fit5, test_data$UNS)
```

```
##
## fit5      High Low Middle Very Low
##   High      39  0      1      0
##   Low       0 39      3      3
##   Middle     0  5     30      0
##   very_low   0  2      0     23
```

Aggregate Data - Add accuracy to each model, compile, add missing algorithms if possible

SVM classification