

Homework 5

Marc Mendez & Joel Cantero

5th May, 2019

First of all, we are going to install and load all the libraries we need for this exercise.

```
packages <- c("rpart", "ROCR", "readxl", "randomForest", "mice")
for (package in packages) {
  if(!require(package, character.only=TRUE)){
    install.packages(package, repos="http://cran.rstudio.com")
    library(package, character.only=TRUE)
  }
}
```

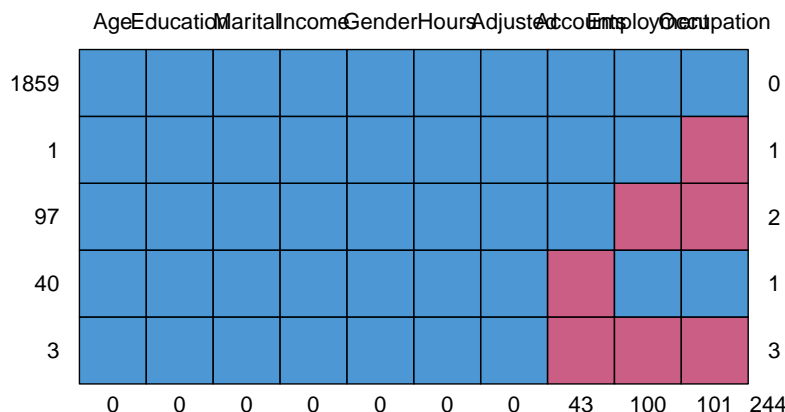
1. Read the Audit.xlsx file and convert it to the csv extension.

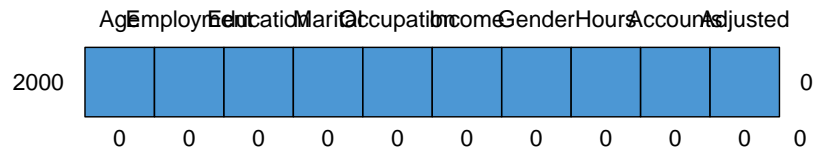
Once we have loaded and installed all the packages, we are going to load the excel file. Also, we are going to convert Employment, Education, Marital, Occupation, Gender, Accounts and Adjusted attributes to factors. Then, we will convert it to CSV file thanks to write.csv2 function.

```
audit <- read_excel("audit.xlsx", na = "NA")
factors <- c("Employment", "Education", "Marital", "Occupation", "Gender",
            "Accounts", "Adjusted")
for (fac in factors) {
  audit[fac] <- lapply(audit[fac], factor)
}
write.csv2(audit, "audit.csv")
```

2. The goal is to use a decision tree to predict the binary “Adjusted” variable, whether the individuals had made a correct financial statement or not. Decide which predictors you would use and eventually preprocess these variables.

First of all, we have to remove these attributes that will not help us for splitting the tree. We can observe that “ID”, “Deductions” and “Adjusted” are related to statement made and we do not need them. After that, we will impute missing values with MICE function. We have found 244 missing values and we have decided to use MICE because it is a small percentage of all our instances.





So we can see that with MICE we solved all the missing values, and now we can proceed with our study.

3. Select the 1/3 of the last observations as test data.

We have to select the last 33% of data instances as test, and the first 66% instances as training data.

```
test <- imputedAudit[seq(0.66*nrow(imputedAudit), nrow(imputedAudit)), ]
training <- imputedAudit[-seq(0.66*nrow(imputedAudit), nrow(imputedAudit)), ]
```

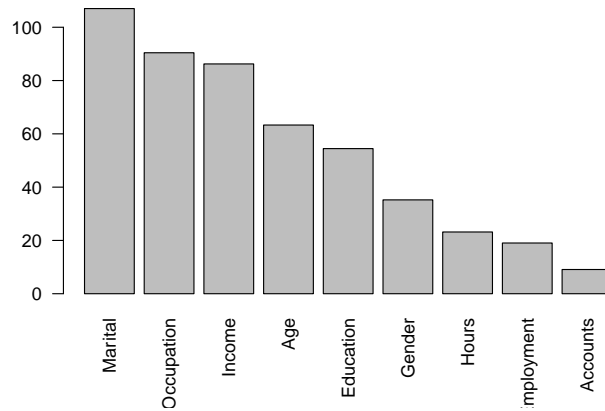
4. Obtain the decision tree to predict whether the variable “Adjusted” on the training data. Decide the cutoff value for taking the decision.

```
##
## Classification tree:
## rpart(formula = Adjusted ~ ., data = training, method = "class",
##       control = rpart.control(xval = 10, cp = 0.001))
##
## Variables actually used in tree construction:
## [1] Accounts   Age          Education  Employment  Gender      Hours
## [7] Income     Marital      Occupation
##
## Root node error: 317/1319 = 0.24033
##
## n= 1319
##
##      CP nsplit rel error  xerror    xstd
## 1 0.1388013      0  1.00000 1.00000 0.048953
## 2 0.0378549      2  0.72240 0.77918 0.044696
## 3 0.0189274      4  0.64669 0.75394 0.044130
## 4 0.0105152      7  0.58991 0.72240 0.043396
## 5 0.0094637     10  0.55836 0.70347 0.042941
## 6 0.0063091     12  0.53943 0.72555 0.043471
## 7 0.0042061     15  0.52050 0.77287 0.044556
## 8 0.0037855     18  0.50789 0.78549 0.044834
## 9 0.0031546     23  0.48896 0.77603 0.044626
## 10 0.0021030     25  0.48265 0.78864 0.044903
## 11 0.0010515     28  0.47634 0.80442 0.045244
## 12 0.0010000     31  0.47319 0.80442 0.045244
```

If we calculate the cutoff and with the cutoff we look at the table we obtain the following:

```
## Cutoff idx: 5 With min value: 0.7464114
## CP: 0.009463722 nsplit: 10 rel error: 0.5583596 xerror: 0.70347 xstd: 0.04294139
```

In our case the cutoff will be on 5 splits. ## 5. Plot the importance of variables in the prediction.



As we can see in this plot, the three most important variables are: marital, occupation and income.

6. Compute the accuracy, precision, recall and AUC on the test individuals.

Before calculating the accuracy, precision, recall and AUC, we need to calculate the confusion matrix. After it we can start calculating the variables said earlier.

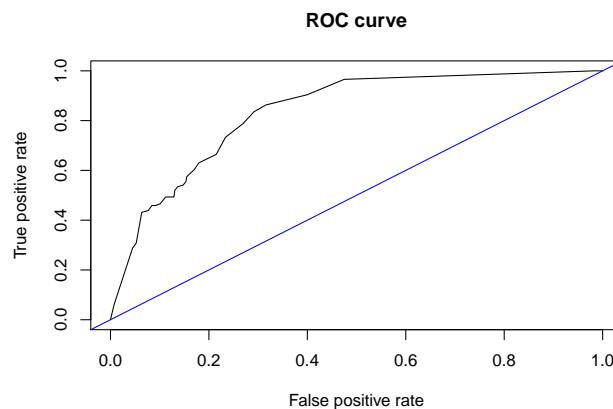
```
prediction <- predict(rt, test, type="class")
(results <- table(test$Adjusted, prediction))
```

```
##      prediction
##      0      1
## 0 486  49
## 1  79  67
```

```
# Accuracy
error <- (results[1,2] + results[2,1])/nrow(test)
accuracy <- 1 - error

# Precision
precision_p <- results[1,1]/(results[1,1] + results[2,1])
precision_n <- results[2,2]/(results[1,2] + results[2,2])
precision <- (precision_n + precision_p)/2

# Recall
recall <- results[1,1]/(results[1,1] + results[1,2])
```



We obtain a very high accuracy. If we look at the amount of results predicted we see that our precision is also high, and last but not least, the recall value which is also very high.

```
## Accuracy: 0.8120411
```

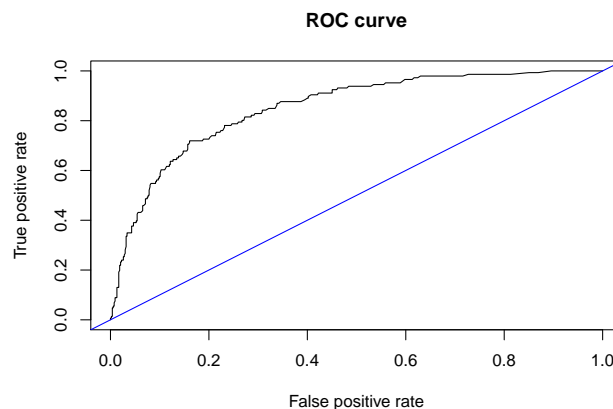
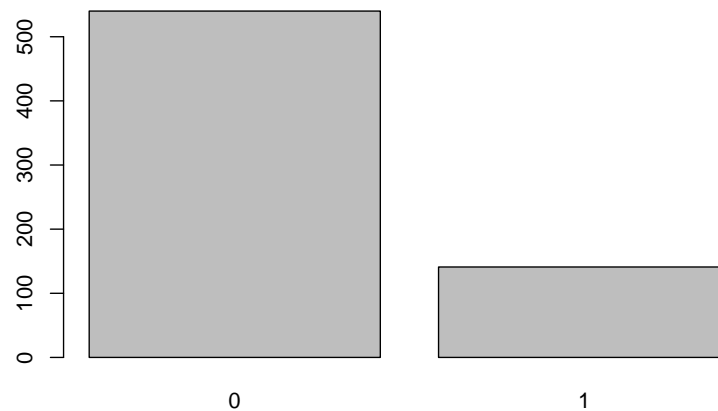
```
## Precision: 0.7188816
## Recall: 0.9084112
## AUC: 0.8348739
```

7. Perform a Random Forest on the same data.

As we did previously, we have to build the model and then test the results with the test data (the same split as we have used before, if we want to compare it).

The goal of this exercise is to compare the previous results with a random forest. For this reason, we will calculate again the accuracy, precision, recall and AUC.

```
##
## Call:
## randomForest(formula = Adjusted ~ ., data = training, importance = TRUE)
##               Type of random forest: classification
##               Number of trees: 500
## No. of variables tried at each split: 3
##
##               OOB estimate of error rate: 16.53%
## Confusion matrix:
##      0   1 class.error
## 0 906  96 0.09580838
## 1 122 195 0.38485804
```



```
## Accuracy: 0.8311307
## Precision: 0.749409
```

```
## Recall: 0.8971963
```

```
## AUC: 0.8503009
```

Conclusions

To conclude, we can say that all the metrics have been improved using a random forest just using 500 trees (we can see that if we print `randomForest` variable). If we just use one decision tree (the previous exercise) against a random forest, we can observe that the results are not good as random forest ones.