

Coursework  
Advanced Data Science (CMM536)

Andrew Tait, [1504693@rgu.ac.uk](mailto:1504693@rgu.ac.uk)

April 25, 2018

# 1 Research

[your text goes here] The paper that was chosen for this work is [1]. The authors provided a full review on different streaming algorithms and methods that handle concept drift. Below is a review of this paper that includes problem statement, related work and methods applied.

*Notice how I cited the paper, and how does it appear in the document, to do so, you need to have a file in your working director (references.bib), this file simply contains the bibtext items for the papers you chose. These BibTex items are often available to download from publishers website, see Figure 1*

## 1.1 Problem Statement

What is this paper about? your text goes here, your text goes here your text goes here your text goes here  
your text goes here your text goes here your text goes here

## 1.2 Relevant Work

[your text goes here your text goes here your text goes here your text goes here your text goes here your  
text goes here your text goes here your text goes here your text goes here your text goes here your text goes  
here your text goes here your text goes here your text goes here]

### 1.3 Methods

[your text goes here your text goes here your text goes here your text goes here your text goes here your  
text goes here your text goes here your text goes here your text goes here your text goes here your text goes  
here your text goes here your text goes here your text goes here]

## 1.4 Results

[your text goes here your text goes here your text goes here your text goes here your text goes here your  
text goes here your text goes here your text goes here your text goes here your text goes here your text goes  
here your text goes here your text goes here your text goes here]

## 1.5 Conclusion

[your text goes here your text goes here your text goes here your text goes here your text goes here your  
text goes here your text goes here your text goes here your text goes here your text goes here your text goes  
here your text goes here your text goes here your text goes here]

## 2 Data Streams

### 2.1 Dataset Choice

The dataset that has been chosen for this part of the coursework is the 'adult' dataset. This is available on the UCI repository.

<https://archive.ics.uci.edu/ml/datasets/adult>

The dataset that has been chosen for this part of the course work is IRIS. This is available on the UCI repository. The set was chosen because of .... It proves to be a good set for evaluating 'x' methods ....

### 2.2 Data Exploration

The main purpose of the adult dataset is to find out which characteristics of the us population affect if their income is either  $\leq 50k$  or  $> 50k$

To start off I will clear the RStudio environment and import the required libraries.

```
#Clean RStudio Environment
rm(list = ls())

#Import libraries
library(caret)
library(partykit)
library(mlbench)
library(RWeka)
library(C50)
library(datasets)
library(rpart)
library(ggplot2)
library(jsonlite)
library(data.table)
library(RMOA)
library(ROCR)
library(stream)
library(mlbench)
library(doParallel)
require(tm)
require(wordcloud)
library(wordcloud,quietly=TRUE)
library(RColorBrewer,quietly=TRUE)
```

Then set the working directory to the coursework folder

```
#Set working directory
setwd("~/CMM536 Advanced Data Science/Coursework/CMM536_Coursework")
```

In order to import the adult dataset, the feature names first need to be defined.

```
#Feature names
adultNames <- c("age", "workclass", "fblwgt",
               "education", "education-num",
               "marital-status",
               "occupation",
               "relationship",
               "race",
               "sex",
               "captial-gain",
               "captain-loss",
               "hours-per-week",
               "native-country",
               "class")
```

The adult dataset is then imported from adult.data file:

```
#Import datasets
adult <- read.table("data/adult.data", header = FALSE, sep = ",",
                  strip.white = TRUE, col.names = adultNames,
                  na.strings = "?", stringsAsFactors = TRUE)
```

Now the adult dataset is imported, it's time for some basic exploration  
Number of rows (instances) in the dataset

```
nrow(adult)
```

```
## [1] 32561
```

The number of columns (features)

```
ncol(adult)
```

```
## [1] 15
```

Summary of the adult dataset:

```
#inspect dataset
str(adult)
```

```
## 'data.frame': 32561 obs. of 15 variables:
## $ age      : int  39 50 38 53 28 37 49 52 31 42 ...
## $ workclass : Factor w/ 8 levels "Federal-gov",...: 7 6 4 4 4 4 4 6 4 4 ...
## $ fblwgt    : int  77516 83311 215646 234721 338409 284582 160187 209642 45781 159449 ...
## $ education : Factor w/ 16 levels "10th","11th",...: 10 10 12 2 10 13 7 12 13 10 ...
## $ education.num : int  13 13 9 7 13 14 5 9 14 13 ...
## $ marital.status: Factor w/ 7 levels "Divorced","Married-AF-spouse",...: 5 3 1 3 3 3 4 3 5 3 ...
## $ occupation  : Factor w/ 14 levels "Adm-clerical",...: 1 4 6 6 10 4 8 4 10 4 ...
## $ relationship : Factor w/ 6 levels "Husband","Not-in-family",...: 2 1 2 1 6 6 2 1 2 1 ...
## $ race        : Factor w/ 5 levels "Amer-Indian-Eskimo",...: 5 5 5 3 3 5 3 5 5 5 ...
## $ sex         : Factor w/ 2 levels "Female","Male": 2 2 2 2 1 1 1 2 1 2 ...
```

```
## $ captial.gain : int 2174 0 0 0 0 0 0 0 14084 5178 ...
## $ captain.loss : int 0 0 0 0 0 0 0 0 0 0 ...
## $ hours.per.week: int 40 13 40 40 40 40 16 45 50 40 ...
## $ native.country: Factor w/ 41 levels "Cambodia","Canada",...: 39 39 39 39 5 39 23 39 39 39 ...
## $ class          : Factor w/ 2 levels "<=50K", ">50K": 1 1 1 1 1 1 1 2 2 2 ...
```

Now that some basic data exploration is covered, next to inspect the dataset a bit further. Starting with the class distribution in the adult dataset, see (Figure 1)

```
#Class Distribution
barplot(table(adult$class))
```

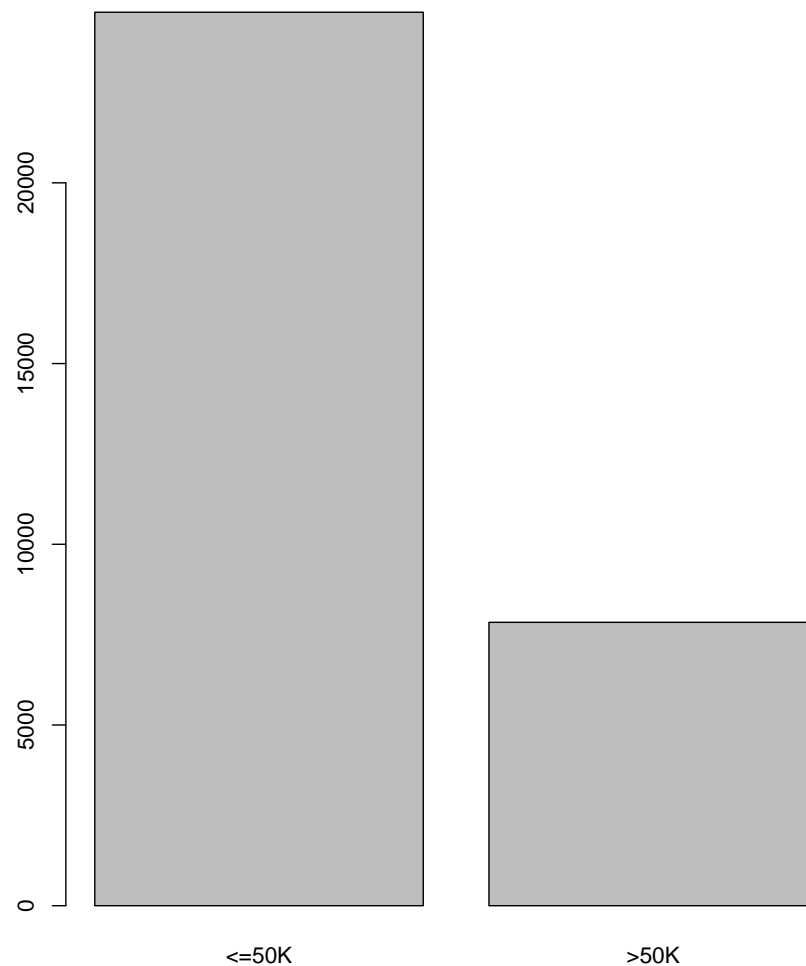


Figure 1: Barplot of Class Distribution

## 2.3 Build Classifier

### 2.3.1 Pre-Processing

Before the adult dataset can be classified, some pre-processing is required first.

The dataset is checked for missing values ('?')

```
#Check for missing values ('?')
table(complete.cases (adult))
```

```
##
## FALSE  TRUE
## 2399 30162
```

As shown above there are missing values in the workclass, occupation and class columns. so these are removed.

```
cleanadult = adult[!is.na(adult$workclass)& !is.na(adult$occupation) & !is.na(adult$class),]
```

The flwgt feature is also removed as it's not required.

```
#Remove flwgt feature
cleanadult$fblwgt = NULL
```

Lets inspect the cleanadult dataframe before going further:

```
str(cleanadult)
```

```
## 'data.frame': 30718 obs. of 14 variables:
## $ age : int 39 50 38 53 28 37 49 52 31 42 ...
## $ workclass : Factor w/ 8 levels "Federal-gov",...: 7 6 4 4 4 4 4 6 4 4 ...
## $ education : Factor w/ 16 levels "10th","11th",...: 10 10 12 2 10 13 7 12 13 10 ...
## $ education.num : int 13 13 9 7 13 14 5 9 14 13 ...
## $ marital.status: Factor w/ 7 levels "Divorced","Married-AF-spouse",...: 5 3 1 3 3 3 4 3 5 3 ...
## $ occupation : Factor w/ 14 levels "Adm-clerical",...: 1 4 6 6 10 4 8 4 10 4 ...
## $ relationship : Factor w/ 6 levels "Husband","Not-in-family",...: 2 1 2 1 6 6 2 1 2 1 ...
## $ race : Factor w/ 5 levels "Amer-Indian-Eskimo",...: 5 5 5 3 3 5 3 5 5 5 ...
## $ sex : Factor w/ 2 levels "Female","Male": 2 2 2 2 1 1 1 2 1 2 ...
## $ captial.gain : int 2174 0 0 0 0 0 0 0 14084 5178 ...
## $ captain.loss : int 0 0 0 0 0 0 0 0 0 0 ...
## $ hours.per.week: int 40 13 40 40 40 40 16 45 50 40 ...
## $ native.country: Factor w/ 41 levels "Cambodia","Canada",...: 39 39 39 39 5 39 23 39 39 39 ...
## $ class : Factor w/ 2 levels "<=50K", ">50K": 1 1 1 1 1 1 1 2 2 2 ...
```

As the adult dataset is a mixture of factors, integers and characters, I decided that transforming the dataset into binary. This will be create features based on every possible value in the dataset.  
The cleanadult dataframe is first copied and the class is removed (it's not being transforme dinto binary)

```
#Copy dataset  
noClass <- cleanadult  
#Remove class as it is not being transformed to binary  
noClass$class <- NULL
```

Then the noClass dataframe is transformed into binary

```
binaryVars <- caret::dummyVars(~ ., data = noClass)  
newAdult <- predict(binaryVars, newdata = noClass)
```

The class feature is then added to the binarised dataset

```
#add class to binarised dataset  
binAdult <- cbind(newAdult, cleanadult[14])
```

Any rows with NA values after being binary transformed.

```
#remove any rows with NA values  
row.has.na <- apply(binAdult, 1, function(x){any(is.na(x))})  
sum(row.has.na)  
binAdult <- binAdult[!row.has.na,]
```

Number of NA rows removed.

```
## [1] 556
```

### 2.3.2 Classification

When it came to dividing the mushroom dataset into training and testing subsets, I decided to go with 80 percent training and 20 percent testing split as a starting point/baseline. Since the class distribution is unbalanced, I thought this split would cover the majority of cases.

```
#split 80% training and 20% testing datasets
inTrain <- createDataPartition(y=binAdult$class, p=0.8, list=FALSE)

#Assign indexes to split the binAdult dataset into training and testing
training <- binAdult[inTrain,]
testing <- binAdult[!inTrain,]
```

For the initial classifier I decided to go with the kNN Classifier as it has proven to be a good baseline in previous labs and exercises in R. Before the classification begins, parallel processing is enabled to speed up this process.

```
#Setup Parallel processing to speed up classification modelling
cl <- makeCluster(detectCores(), type='PSOCK')
registerDoParallel(cl)
```

The train control is set to repeated-cross-validation with 10 folds and 3 repeats

```
ctrl <- trainControl(method = "repeatedcv",
                     repeats = 3,
                     number = 10,
                     verboseIter=TRUE)
```

Next the seed is set to 1, in order to make the model reproducible and the kNN model is set up with the train control from above and k value set to 3.

```
# ensure reproducibility of results by setting the seed to a known value
set.seed(1)
#use knn
mod21.knn<- train(class~., data=training,
                  method="knn", tuneGrid=expand.grid(.k=3),trControl=ctrl)
```

```
## Aggregating results
## Fitting final model on full training set
```

Once the knn Model is complete, it's time to analyse the results, first with a print of the kNNModel as shown below.

```
print(mod21.knn)
```

```
## k-Nearest Neighbors
##
## 24131 samples
##   104 predictor
##     2 classes: '<=50K', '>50K'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold, repeated 3 times)
## Summary of sample sizes: 21717, 21718, 21719, 21718, 21718, 21718, ...
## Resampling results:
##
##   Accuracy   Kappa
##   0.8371529  0.5550044
##
## Tuning parameter 'k' was held constant at a value of 3
```



### 2.3.3 Evaluation

To evaluate the module, a confusion matrix is produced by predicting the against the testing (20 percent of the total dataset ) subset.

```
#Evaluation
predictkNN <- predict(mod21.knn,testing)
confusionMatrix(predictkNN, testing$class)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction <=50K >50K
##    <=50K 17070  1504
##    >50K  1054  4503
##
##           Accuracy : 0.894
##           95% CI : (0.89, 0.8979)
##    No Information Rate : 0.7511
##    P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.7092
##  Mcnemar's Test P-Value : < 2.2e-16
##
##           Sensitivity : 0.9418
##           Specificity : 0.7496
##           Pos Pred Value : 0.9190
##           Neg Pred Value : 0.8103
##           Prevalence : 0.7511
##           Detection Rate : 0.7074
##   Detection Prevalence : 0.7697
##           Balanced Accuracy : 0.8457
##
##           'Positive' Class : <=50K
##
```

## 2.4 Build Stream Classifier

### 3 Text Classification

For this task a csv of leave/remain tweets of the brexit campaign was provided. The was first imported:

```
leaveRemainTweets <- read.csv("data/leaveRemainTweets_CW.csv", header=TRUE)
```

#### 3.1 Preprocessing

In order to pre-process the leaveRemain tweets, I built a custom corpus function. Which can be found at (Figure ??)

```
#Custom Corpus Function for preprocessing

buildCorpus <- function (tweets, wholeDataSet, wordAssociation){
  corp <- Corpus(VectorSource(tweets$text))

  corp <- tm_map(corp,
                 content_transformer(function(x) iconv(x, to='ASCII',
                                                         sub='byte')))

  # remove stop words and other preprocessing

  corp <- tm_map(corp, content_transformer(tolower))
  corp <- tm_map(corp, removeNumbers)

  toSpace = content_transformer( function(x, pattern) gsub(pattern, " ",x) )

  ##Tweet cleaning

  #credit to https://stackoverflow.com/a/31352005/8816204
  corp <-tm_map(corp, toSpace, "(RT|via)((?:\\b\\W*@[\\w+)+)")
  corp <-tm_map(corp, toSpace, "@\\w+")
  corp <-tm_map(corp, toSpace, "&amp;")

  ##Remove URLs from tweets
  corp <-tm_map(corp, toSpace, "httpsw+")
  corp <-tm_map(corp, toSpace, "http:\\w+")
  corp <-tm_map(corp, toSpace, "https:\\w+")

  corp <-tm_map(corp, toSpace, "[ \\t]{2,}")
  corp <-tm_map(corp, toSpace, "^\\s+|\\s+$")

  # ##Remove obvious words /stopwords
  if(missing(wholeDataSet)){
    corp <- tm_map(corp,
                   function(x)removeWords(x,c(
                     stopwords("english"),"amp", "will", "\\%Û_", "https", "http", "httpsdb")))
  } else if(wholeDataSet == TRUE){
    corp <- tm_map(corp,
                   function(x)removeWords(x,c(
```

```

        stopwords("english"), "amp", "will", "%Ű_", "https", "http",
        "httpsdb", "eu", "brexit", "rt", "leave", "remain", "vote")))
    }

    #remove punctuation last so urls are removed correctly
    corp <- tm_map(corp, removePunctuation)

    if(missing(wordAssociation)){
        tdm <- TermDocumentMatrix(corp)

        m <- as.matrix(tdm)
        v <- sort(rowSums(m), decreasing = TRUE)
        d <- data.frame(word = names(v), freq = v)
        d$word <- gsub("~", " ", d$word) ## Edit 2

        tweets <- d$word
    }else if(wordAssociation == TRUE){
        dtm <- DocumentTermMatrix(corp)
    }
}

```

## 3.2 Text Analysis

### 3.3 Text Classification

Use one of the ‘R’ packages to build a classifier that classifies the tweets as leave tweet or remain tweets. Complete this part as required by the coursework sheet. Again, be clear, visuals always helps in communicating results. Justify your choices and explain your methods.

## 4 Appendix

### 4.1 Custom buildCorpus Function

For the pre-processing in the text-classification task, I created a custom corpus function that will build a corpus of the leave/remain tweets and change the word filtering based on the optional function parameters.

## 5 References

- [1] P. B. Dongre and L. G. Malik. “A review on real time data stream classification and adapting to various concept drift scenarios”. In: *2014 IEEE International Advance Computing Conference (IACC)*. 2014, pp. 533–537. DOI: [10.1109/IAAdCC.2014.6779381](https://doi.org/10.1109/IAAdCC.2014.6779381).