# Configuration Manual

MSc Research Project
Data Analytics

## Anamika Chavan
Student ID: x18199950

School of Computing
National College of Ireland

Supervisor:     Pierpaolo Dondio

# National College of Ireland
## Project Submission Sheet
## School of Computing

| | |
|---|---|
| **Student Name:** | Anamika Chavan |
| **Student ID:** | x18199950 |
| **Programme:** | Data Analytics |
| **Year:** | 2019 |
| **Module:** | MSc Research Project |
| **Supervisor:** | Pierpaolo Dondio |
| **Submission Due Date:** | 12/12/2019 |
| **Project Title:** | Recruitment of Suitable Football Player by using Machine Learning Techniques. |
| **Word Count:** | 547 |
| **Page Count:** | 8 |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

**ALL** internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

| | |
|---|---|
| **Signature:** | |
| **Date:** | 11th December 2019 |

## PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies). | ☐ |
| **Attach a Moodle submission receipt of the online project submission**, to each project (including multiple copies). | ☐ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | ☐ |

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Configuration Manual

Anamika Chavan
x18199950

# 1 Introduction

A complete guidelines for the implementation of the research "Recruitment of Suitable Football Player By using Machine Learning Techniques" is given in this document.This research have been developed in R-studio. All the libraries and packages used to develop this project are mentioned in this document.

# 2 Hardware Requirements

This research study was done on "DELL Inspiron 13" laptop. Hardware configuration of this laptop is as follows:
**Operating System**: Windows 10
**RAM** : 8GB
**Processor** : Core i5
**Storage**: 256GBSSD
　　All above configuration are sufficient to run this project.

# 3 Software Requirements

For this research study R Studio and PowerBI was used. So below steps will explain installation process for R and PowerBI.

## 3.1 Download and Install R Studio

- Download the R-studio server for Windows 10 from the below link.
  https://rstudio.com/products/rstudio/download/

- After downloading the R-studio, next step is to install the R-studio. To install R-studio, refer the following link which contains clear instruction about installation.
  http://rprogramming.net/download-and-install-rstudio/

## 3.2 Download and Install PowerBI

- Download the PowerBI for Windows 10 from the below link.
  https://powerbi.microsoft.com/en-us/downloads/

- After downloading the powerBI, next step is to install the powerBI. To install powerBI, refer the following link which contains clear instruction about installation. https://www.knowledgehut.com/blog/business-intelligence-and-visualization/how-to-install-power-bi

# 4 Implementation of the Models

After installation of all software, implementation of the project can be done as follows.

## 4.1 Download the dataset

For this research, data is collected from kaggle which contains information about various players. Download the dataset from below link website.
https://www.knowledgehut.com/blog/business-intelligence-and-visualization/how-to-install-power-bi

## 4.2 Import the libraries

For this research, we have used R language to develop machine learning models. To do so, we need to firstly clean the data. Hence, to clean and develop the model, we need to import some R libraries. Below are the R libraries used for this research.

```
1
2   library(dplyr)
3   library(Hmisc)
4   library(corrplot)
5   library(caTools)
6   library(tidyr)
7   library(e1071)
8   library(caret)
9   library(FNN)
10  library(MASS)
11  library(rpart)
12  library(xgboost)
13
```

Figure 1: Imported Libraries

## 4.3 Pre-processing of the data

After downloading the dataset,we preprocessed the data before applying it to the machine learning models. Preprocessing steps includes checking missing values,feature engineering, data encoding and data scaling. The screenshot of code snippet is given below.

```r
#checking for null values
players_data[players_data==""] <- NA
which(is.na(players_data) == TRUE, arr.ind=TRUE)
#remove unwanted column
players_data[ ,c('Name','Club','Contract_Expiry',
                 'Nationality','Club_Joining','Birth_Date','National_Position',
                 'National_Kit','Club_Position','Club_Kit')] <- list(NULL)
#data Encoding
players_data$Work_Rate_1 <- as.factor(players_data$Work_Rate_1)
players_data$Work_Rate_1 <- factor(players_data$Work_Rate_1,
                                   levels = c("Low ", "Medium ", "High "),ordered = TRUE)

# convert columns character to numeric
cols.num <- c("Height","Weight")
players_data[cols.num] <- sapply(players_data[cols.num],as.numeric)

#convert rating into categories
players_data$Rating <- cut(players_data$Rating, breaks=c(40,45,50,55,60,65,70,75,80,90,100),
                           labels = c("1","2","3","4","5","6","7","8","9","10"))

#separate dataframe for each category
ForwardPlayers<- players_data%>% filter(players_data$Preffered_Position =="Forward")
ForwardPlayers$Preffered_Position <- NULL

MidfielderPlayers<- players_data%>% filter(players_data$Preffered_Position =="Midfielder")

DefenderPlayers<- players_data%>% filter(players_data$Preffered_Position =="Defender")

GoalkeeperPlayers<- players_data%>% filter(players_data$Preffered_Position =="Goalkeeper")
# Feature Scaling
training_set[c(2,3,5,7:41)] <- scale(training_set[c(2,3,5,7:41)])
test_set[c(2,3,5,7:41)] <- scale(test_set[c(2,3,5,7:41)])
```

Figure 2: Preprocessing of the Data

## 4.4 SVM Model

After pre-processing of the data, now we will apply all the models one by one. We will start with Support Vector Machine(SVM) model.Atter applying the model on the dataset, we have calculated evaluation metrics. The screenshot of all the above process is given below.

```
35  #implementation of SVM
36
37  library(e1071)
38  library(caret)
39  folds <- createFolds(data$Rating, k = 10)
40
41▾ cv <- lapply(folds, function(x) {
42    training_fold = data[-x, ] # training fold =  training set minus (-) it's sub test fold
43    test_fold = data[x, ] # here we describe the test fold individually
44    # now apply (train) the classifer on the training_fold
45    classifier = svm(formula = Rating ~ .,
46                     data = training_fold,
47                     type = 'C-classification',
48                     kernel = 'radial')
49    y_pred = predict(classifier, newdata = test_fold[-1])
50
51    cm = table(test_fold$Rating, y_pred)
52
53    return(cm)
54  })
55
56  cm <- Reduce('+', cv)
57  accuracy_svm <- sum(diag(cm))/sum(cm)
58  precision_svm <- diag(cm)/colSums(cm)
59  avg_precision_svm <- mean(as.numeric(precision_svm), na.rm=TRUE)
60  recall_svm <- diag(cm)/rowSums(cm)
61  avg_recall_svm <- mean(recall_svm, na.rm=TRUE)
62  Fmeasure_SVM <- 2 * avg_precision_svm * avg_recall_svm / (avg_precision_svm + avg_recall_svm)
63
```

Figure 3: Implementation of SVM model

## 4.5   LDA model

Now we will apply Linear Discriminant Analysis(LDA) model on the data. The screenshot
of the process is given below.

```
66  #applying LDA model
67
68  library(MASS)
69
70▾ cv_lda <- lapply(folds, function(x) { # start of function
71    # in the next two lines we will separate the Training set into it's 10 pieces
72    training_fold = data[-x, ] # training fold =  training set minus (-) it's sub test fold
73    test_fold = data[x, ] # here we describe the test fold individually
74    # now apply (train) the classifer on the training_fold
75    classifier = lda(Rating ~ .,training_fold)
76    # next step in the loop, we calculate the predictions and cm and we equate the accuracy
77    # note we are training on training_fold and testing its accuracy on the test_fold
78    y_pred = predict(classifier, newdata = test_fold[-1])
79
80    cm = table(test_fold$Rating, y_pred$class)
81    return(cm)
82  })
83  cm <- Reduce('+', cv_lda)
84  accuracy_lda <- sum(diag(cm))/sum(cm)
85  precision_lda <- diag(cm)/colSums(cm)
86  avg_precision_lda <- mean(precision_lda, na.rm=TRUE)
87  recall_lda <- diag(cm)/rowSums(cm)
88  avg_recall_lda <- mean(recall_lda, na.rm=TRUE)
89  Fmeasure_lda <- 2 * avg_precision_lda * avg_recall_lda / (avg_precision_lda + avg_recall_lda)
90
```

Figure 4: Implementation of LDA model

## 4.6 Naive Bayes model

Developing the Naive Bayes model is given below.

```
254  #naive_bayes
255
256 ▾ cv_naive_bayes <- lapply(folds, function(x) { # start of function
257      # in the next two lines we will separate the Training set into it's 10 pieces
258      training_fold = ForwardPlayers[-x, ] # training fold =  training set minus (-) it's sub test fold
259      test_fold = ForwardPlayers[x, ] # here we describe the test fold individually
260      # now apply (train) the classifer on the training_fold
261      classifier = naiveBayes(formula = Rating ~ .,
262                              data = training_fold,
263                              probability = TRUE)
264      # next step in the loop, we calculate the predictions and cm and we equate the accuracy
265      # note we are training on training_fold and testing its accuracy on the test_fold
266      y_pred = predict(classifier, newdata = test_fold[-1])
267
268      cm = table(test_fold$Rating, y_pred)
269      accuracy <- sum(diag(cm))/sum(cm)
270      return(cm)
271  })
272  cm <- Reduce('+', cv_naive_bayes)
273  accuracy_naive_bayes <- sum(diag(cm))/sum(cm)
274  precision_naive_bayes <- diag(cm)/colSums(cm)
275  avg_precision_naive_bayes <- mean(precision_naive_bayes, na.rm=TRUE)
276  recall_naive_bayes <- diag(cm)/rowSums(cm)
277  avg_recall_naive_bayes <- mean(recall_naive_bayes, na.rm=TRUE)
278  Fmeasure_naive_bayes <- 2 * avg_precision_naive_bayes * avg_recall_naive_bayes /
279                          (avg_precision_naive_bayes + avg_recall_naive_bayes)
280
281
282  accuracy_naive_bayes <- mean(as.numeric(cv_naive_bayes))
283
284
```

Figure 5: Implementation of Naive Bayes model

## 4.7 Decision Tree model

Code snippet of building decision tree model is given below.

```
92   #decision tree
93
94   library(rpart)
95
96 - cv_rpart <- lapply(folds, function(x) { # start of function
97       # in the next two lines we will separate the Training set into it's 10 pieces
98       training_fold = ForwardPlayers[-x, ]
99       # training fold =  training set minus (-) it's sub test fold
100      test_fold = ForwardPlayers[x, ]
101      # here we describe the test fold individually
102      # now apply (train) the classifer on the training_fold
103      classifier = rpart(Rating~., data = training_fold, method = 'class')
104      # next step in the loop, we calculate the predictions and cm and we equate the accuracy
105      # note we are training on training_fold and testing its accuracy on the test_fold
106      y_pred = predict(classifier, test_fold[-1], type = 'class')
107
108      cm = table(test_fold$Rating, y_pred)
109      accuracy <- sum(diag(cm))/sum(cm)
110      return(cm)
111  })
112
113  cm <- Reduce('+', cv_rpart)
114  conf_matrix <- confusionMatrix(cm)
115  accuracy_rpart <- sum(diag(cm))/sum(cm)
116  precision_rpart <- diag(cm)/colSums(cm)
117  avg_precision_rpart <- mean(precision_rpart, na.rm=TRUE)
118  recall_rpart <- diag(cm)/rowSums(cm)
119  avg_recall_rpart <- mean(recall_rpart, na.rm=TRUE)
120  Fmeasure_rpart <- 2 * avg_precision_rpart * avg_recall_rpart /
121              (avg_precision_rpart + avg_recall_rpart)
```

Figure 6: Implementation of Decision Tree model

## 4.8   XGBoost model

Code snippet of building XGBoost model is given below.

```
154  #implementation of XGBoost
155  folds <- createFolds(training_set$Rating, k = 10)
156  library(xgboost)
157 - cv_xgboost <- lapply(folds, function(x) { # start of function
158      # in the next two lines we will separate the Training set into it's 10 pieces
159      ForwardPlayers$Preffered_Foot <- as.numeric(ForwardPlayers$Preffered_Foot)
160      training_fold = training_set[-x, ] # training fold =  training set minus (-) it's sub test fold
161      test_fold = training_set[x, ] # here we describe the test fold individually
162      # now apply (train) the classifer on the training_fold
163
164      XTrain <- data.frame(lapply(training_fold[,-1],as.numeric))
165      XTrain <- as.matrix(XTrain)
166      yTrain <- unclass(training_fold$Rating)-1
167      m.xg.def <- xgboost(data=XTrain,label=yTrain,objective="multi:softmax",num_class=10,nrounds = 1)
168
169      XTest <- data.frame(lapply(test_set[,-1],as.numeric))
170      XTest <- as.matrix(XTest)
171
172      y.xg.def <- predict(m.xg.def,newdata=XTest)+1
173      y.gbm.default1 <- factor(y.xg.def, levels = c(1,2,3,4,5,6,7,8,9,10), ordered = TRUE)
174      cm <- table(test_set$Rating,y.gbm.default1)
175      return(cm)
176  })
177  cm <- Reduce('+',cv_xgboost)
178  accuracy_xgboost <- sum(diag(cm))/sum(cm)
179  precision_xgboost <- diag(cm)/colSums(cm)
180  avg_precision_xgboost <- mean(precision_xgboost, na.rm=TRUE)
181  recall_xgboost <- diag(cm)/rowSums(cm)
182  avg_recall_xgboost <- mean(recall_xgboost, na.rm=TRUE)
183  Fmeasure_xgboost <- 2 * avg_precision_xgboost * avg_recall_xgboost / (avg_precision_xgboost + avg_recall_xgboost)
184
```

Figure 7: Implementation of XGBoost model

## 4.9   KNN model

Code snippet of building KNN model is given below.

```
127   #knnn
128   library(FNN)
129   ForwardPlayers$Rating <- as.numeric(ForwardPlayers$Rating)
130   ForwardPlayers$Work_Rate_1 <- as.numeric(ForwardPlayers$Work_Rate_1)
131   ForwardPlayers$Work_Rate_2 <- as.numeric(ForwardPlayers$Work_Rate_2)
132   ForwardPlayers$Preffered_Foot <- as.numeric(ForwardPlayers$Preffered_Foot)
133 ▾ cv_knn <- lapply(folds, function(x) { # start of function
134     # in the next two lines we will separate the Training set into it's 10 pieces
135     training_fold = ForwardPlayers[-x, ] # training fold =  training set minus (-) it's sub test fold
136     test_fold = ForwardPlayers[x, ] # here we describe the test fold individually
137     # now apply (train) the classifer on the training_fold
138     training_fold <- ForwardPlayers[-x,]
139     ##extract testing set
140     test_fold <- ForwardPlayers[x,]
141     ##extract 5th column of train dataset because it will be used as 'cl' argument in knn function.
142     target_category <- ForwardPlayers[-x,1]
143     ##extract 5th column if test dataset to measure the accuracy
144     test_category <- ForwardPlayers[x,1]
145     test_category <- factor(test_category, levels = c(1,2,3,4,5,6,7,8,9,10), ordered = TRUE)
146     k <- knn(training_fold, test_fold, training_fold$Rating, k = 8)
147
148     k <- factor(k, levels = c(1,2,3,4,5,6,7,8,9,10), ordered = TRUE)
149     cm <- table(k,test_category)
150     return(cm)
151   })
152   cm <- Reduce('+',cv_knn)
153   accuracy_knn <- sum(diag(cm))/sum(cm)
154   precision_knn <- diag(cm)/colSums(cm)
155   avg_precision_knn <- mean(precision_knn, na.rm=TRUE)
156   recall_knn <- diag(cm)/rowSums(cm)
157   avg_recall_knn <- mean(recall_knn, na.rm=TRUE)
158   Fmeasure_knn <- 2 * avg_precision_knn * avg_recall_knn / (avg_precision_knn + avg_recall_knn)
```

Figure 8: Implementation of KNN model

## 4.10   Finding the closest match

To find the closest match for the replaced player, we have used knn model.Code snippet of finding the closest match is given below.

```
385
386   k <- knn(training_set[,-1], test_set[,-1], labels, k = 10)
387   indices = attr(k, "nn.index")
388   print(indices[156, ])
389
```

Figure 9: Implementation of closest match

## 4.11 Evaluation Results by using powerBI

We have compared the two evaluation metrics accuracy and F-measure for different models. Following screenshot shows the comparison between different models in terms of accuracy and F-measure.
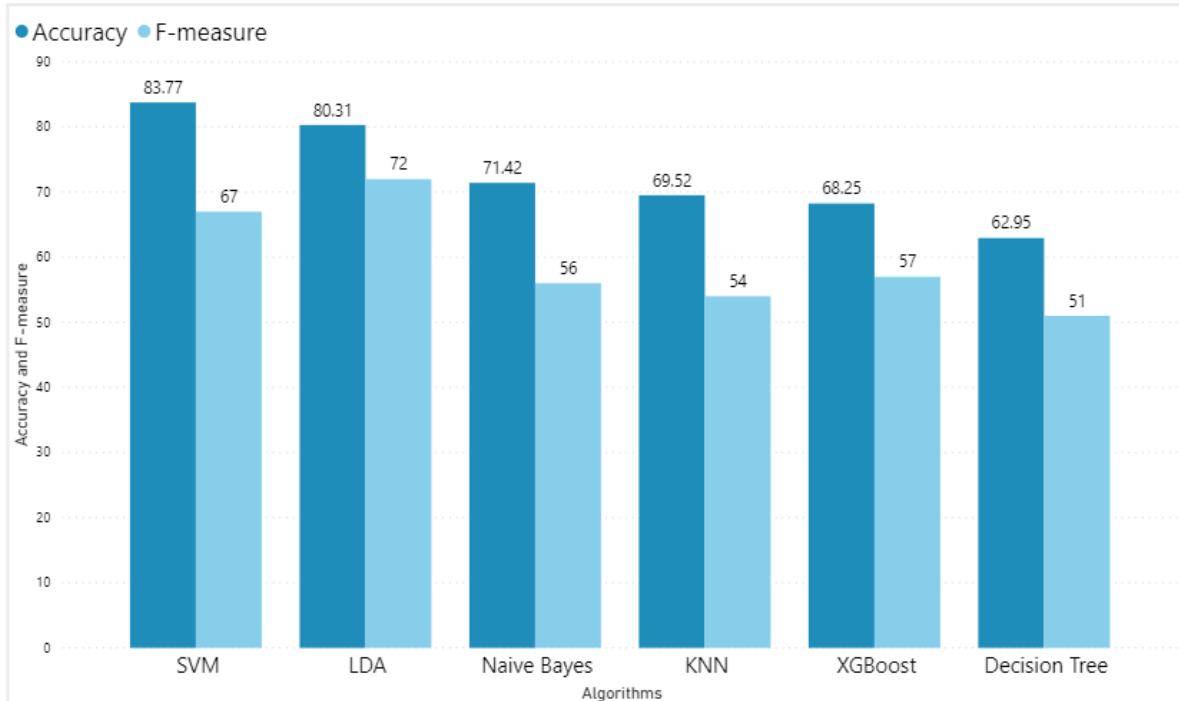


Figure 10: Accuracy and F-measure score by %