# FOOTBALL MATCH PREDICTION

# **Submitted by**

RAJA HARIKESH N V: 16BLC1021

PRANAV. N : 16BLC1092

PRAVEEN KUMAR .R: 16BLC1135

J Component - Report

**ECM2002 – Machine Learning Algorithms** 

#### **BACHELOR OF TECHNOLOGY**

in

#### ELECTRONICS AND COMPUTER ENGINEERING



October 2018

# TABLE OF CONTENTS

Chapter	Title	Page
1	ABSTRACT	3
2	SECTION I - Data Set Description 2.1 Data Set Description 2.2 Data Set Pre-Processing	4 6
3	SECTION II - Basic Models  3.1 Logit 3.2 LDA 3.3 QDA 3.4 KNN 3.5 NAÏVE BAYES	7 8 9 10 11
4	SECTION III – Alternate/Advanced Models  4.1 SVM 4.2 NEURAL NETWORKS 4.3 RANDOM FORESTS	12 16 18
5	SECTION IV – Comparison & Conclusion  5.1 Comparison of All models' accuracy 5.2 Conclusion	18 19

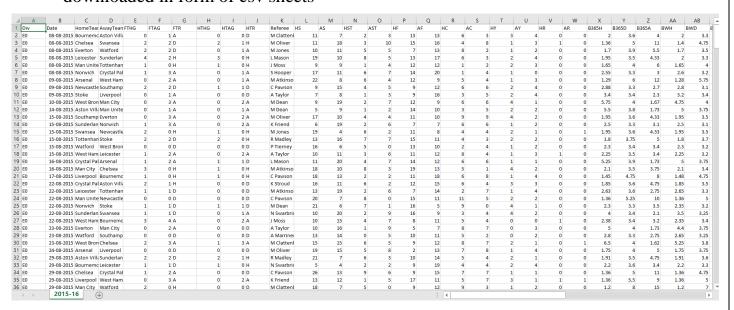
# 1. ABSTRACT:

- 1. Every country in Europe has a club based football league. Each league has 20 teams which play two matches with every other team in their league one at their home stadium and the opponent's home stadium. Each such match has three possible outcomes the home team wins, the match ends in a draw or the visiting team wins.
- 2. English Premier League is the most popular football league, being watched by an estimated figure of 5 billion people across the globe. In a season, since each team plays two games against every other team, there are a total of 380 games.
- 3. Given such a format, it is natural that there are several online fantasy leagues, betting agencies and pundits who try to predict the outcome of each match.
- 4. In this project, an attempt has been made to find out the factors that affect the outcome of a match and also to predict the results of any fixture by using these factors.

# **2.1 DATASET DESCRIPTION:**

Source of Data: <a href="http://www.football-data.co.uk/englandm.php">http://www.football-data.co.uk/englandm.php</a>.

16 seasons of data comprising 380 matches each season were manually downloaded in form of csy sheets



```
All data is in csv format, ready for use within standard spreadsheet
applications.
Div = League Division
Date = Match Date (dd/mm/yy)
HomeTeam = Home Team
AwayTeam = Away Team
FTHG and HG = Full Time Home Team Goals
FTAG and AG = Full Time Away Team Goals
FTR and Res = Full Time Result (H=Home Win, D=Draw, A=Away Win)
HTHG = Half Time Home Team Goals
HTAG = Half Time Away Team Goals
HTR = Half Time Result (H=Home Win, D=Draw, A=Away Win)
Match Statistics (where available)
Attendance = Crowd Attendance
Referee = Match Referee
HS = Home Team Shots
AS = Away Team Shots
HST = Home Team Shots on Target
AST = Away Team Shots on Target
HHW = Home Team Hit Woodwork
AHW = Away Team Hit Woodwork
HC = Home Team Corners
AC = Away Team Corners
HF = Home Team Fouls Committed
AF = Away Team Fouls Committed
HFKC = Home Team Free Kicks Conceded
AFKC = Away Team Free Kicks Conceded
HO = Home Team Offsides
AO = Away Team Offsides
HY = Home Team Yellow Cards
AY = Away Team Yellow Cards
HR = Home Team Red Cards
AR = Away Team Red Cards
HBP = Home Team Bookings Points (10 = yellow, 25 = red)
ABP = Away Team Bookings Points (10 = yellow, 25 = red)
Key to 1X2 (match) betting odds data:
B365H = Bet365 home win odds
B365D = Bet365 draw odds
B365A = Bet365 away win odds
BSH = Blue Square home win odds
BSD = Blue Square draw odds
BSA = Blue Square away win odds
BWH = Bet&Win home win odds
BWD = Bet&Win draw odds
BWA = Bet&Win away win odds
GBH = Gamebookers home win odds
GBD = Gamebookers draw odds
GBA = Gamebookers away win odds
IWH = Interwetten home win odds
IWD = Interwetten draw odds
IWA = Interwetten away win odds
LBH = Ladbrokes home win odds
LBD = Ladbrokes draw odds
LBA = Ladbrokes away win odds
PSH and PH = Pinnacle home win odds
PSD and PD = Pinnacle draw odds
```

```
PSA and PA = Pinnacle away win odds
SOH = Sporting Odds home win odds
SOD = Sporting Odds draw odds
SOA = Sporting Odds away win odds
SBH = Sportingbet home win odds
SBD = Sportingbet draw odds
SBA = Sportingbet away win odds
SJH = Stan James home win odds
SJD = Stan James draw odds
SJA = Stan James away win odds
SYH = Stanleybet home win odds
SYD = Stanleybet draw odds
SYA = Stanleybet away win odds
VCH = VC Bet home win odds
VCD = VC Bet draw odds
VCA = VC Bet away win odds
WHH = William Hill home win odds
WHD = William Hill draw odds
WHA = William Hill away win odds
Bb1X2 = Number of BetBrain bookmakers used to calculate match odds averages
and maximums
BbMxH = Betbrain maximum home win odds
BbAvH = Betbrain average home win odds
BbMxD = Betbrain maximum draw odds
BbAvD = Betbrain average draw win odds
BbMxA = Betbrain maximum away win odds
BbAvA = Betbrain average away win odds
MaxH = Oddsportal maximum home win odds
MaxD = Oddsportal maximum draw win odds
MaxA = Oddsportal maximum away win odds
AvgH = Oddsportal average home win odds
AvgD = Oddsportal average draw win odds
AvgA = Oddsportal average away win odds
Key to total goals betting odds:
BbOU = Number of BetBrain bookmakers used to calculate over/under 2.5 goals
(total goals) averages and maximums
BbMx>2.5 = Betbrain maximum over 2.5 goals
BbAv>2.5 = Betbrain average over 2.5 goals
BbMx<2.5 = Betbrain maximum under 2.5 goals
BbAv<2.5 = Betbrain average under 2.5 goals
GB>2.5 = Gamebookers over 2.5 goals
GB < 2.5 = Gamebookers under 2.5 goals
B365>2.5 = Bet365 \text{ over } 2.5 \text{ goals}
B365 < 2.5 = Bet365 under 2.5 goals
Key to Asian handicap betting odds:
BbAH = Number of BetBrain bookmakers used to Asian handicap averages and
BbAHh = Betbrain size of handicap (home team)
BbMxAHH = Betbrain maximum Asian handicap home team odds
BbAvAHH = Betbrain average Asian handicap home team odds
BbMxAHA = Betbrain maximum Asian handicap away team odds
```

```
BbAvAHA = Betbrain average Asian handicap away team odds
GBAHH = Gamebookers Asian handicap home team odds
GBAHA = Gamebookers Asian handicap away team odds
GBAH = Gamebookers size of handicap (home team)
LBAHH = Ladbrokes Asian handicap home team odds
LBAHA = Ladbrokes Asian handicap away team odds
LBAH = Ladbrokes size of handicap (home team)
B365AHH = Bet365 Asian handicap home team odds
B365AHA = Bet365 Asian handicap away team odds
B365AH = Bet365 size of handicap (home team)

Closing odds (last odds before match starts)

PSCH = Pinnacle closing home win odds
PSCD = Pinnacle closing draw odds
PSCA = Pinnacle closing away win odds
```

# **2.2 PREPROCESSING OF DATA:**

Preprocessing of data was carried using various tools in Python.

#### TOOLS USED:

Pandas: Loading the data, data wrangling and manipulation, feature

engineering.

Scikitlearn: Libraries for classifiers, model evaluation, metrics, cross-

validation

Matplotlib: Data visualization

The main objectives in pre-processing are to extract the meaningful features relating to on-field football and to separate whole data into 2 for model training and testing purposes.

Training data: 15 seasons = 5700 matches

Test data: 1 season = 380 matches.

#### **3.1 LOGIT:**

```
> dataset = read.csv("D:/SEMESTER V/B1_ECM2002/J/Football-Data-Analysis-an
d-Prediction-master/Football-Data-Analysis-and-Prediction-master/Datasets/
final_dataset.csv", header = TRUE)
> dim(dataset)
[1] 6080
> test = read.csv("D:/SEMESTER V/B1_ECM2002/J/Football-Data-Analysis-and-P
rediction-master/Football-Data-Analysis-and-Prediction-master/Datasets/tes
t.csv", header = TRUE)
> dim(test)
[1] 380 43
> library(clusterSim)
> dataset = data.frame(lapply(dataset, function(x) as.numeric(x)))
> dataset = data.Normalization (dataset,type="n4",normalization="column")
> test = data.frame(lapply(test, function(x) as.numeric(x)))
> test = data.Normalization (test,type="n4",normalization="column")
> glm.fits=glm(FTR ~ .,data=dataset,family=binomial)
Warning messages:
1: glm.fit: algorithm did not converge
2: glm.fit: fitted probabilities numerically 0 or 1 occurred
> #glm.fits=glm(FTR ~ HomeTeam + AwayTeam + DiffPts+DiffLP+ DiffFormPts,da
ta=dataset,family=binomial)
> summary(glm.fits)
call:
glm(formula = FTR ~ ., family = binomial, data = dataset)
Deviance Residuals:
                                               Median
                                 10
                                                                      3Q
             Min
-0.0000134258884
                 -0.0000000210734
                                      0.0000000210734
                                                        0.0000074003043
                                                                           0
.0000114591770
Coefficients: (2 not defined because of singularities)
                                                 z value Pr(>|z|)
                      Estimate
                                      Std. Error
                  24.777835096
                                41140.470790097
                                                  0.00060
                                                           0.99952
(Intercept)
                                                           1.00000
                   2.832266444 583824.709939697
                                                  0.00000
Χ
                                                  0.00000
                  -2.740236949 581339.323508215
                                                           1.00000
Date
                  -0.113207924
                                  5940.976712012 -0.00002
                                                           0.99998
HomeTeam
                  -0.101499952
                                  5954.867029370 -0.00002
                                                           0.99999
AwayTeam
                -430.091202898
                                 32670.733741953 -0.01316
                                                           0.98950
FTHG
                 286.241820871
                                 22348.509378651
                                                  0.01281
                                                           0.98978
FTAG
                                                           1.00000
                  -0.134851137
                                 39014.916733706
                                                  0.00000
HTGS
                   0.142451016
                                 37900.521012846
                                                           1.00000
                                                  0.00000
ATGS
                                                  0.00001
                                                           0.99999
HTGC
                   0.255795467
                                 34719.176369211
                   0.246261645
                                                           0.99999
                                 33631.291961249
                                                  0.00001
ATGC
                   0.350616531
                                 29109.485458441
                                                  0.00001
                                                           0.99999
HTP
                   0.187161381
                                 28679.114968246
                                                  0.00001
                                                           0.99999
ATP
                                 29845.014706247
                                                  0.00001
                                                           0.99999
                   0.336088210
HM1
                                                  0.00002
                   0.292520241
                                 12078.595036077
                                                           0.99998
нм2
                   0.139292378
                                  7096.816970097
                                                  0.00002
                                                           0.99998
нм3
                                                  0.00000
                   0.028816661
                                  5963.409513349
                                                           1.00000
нм4
                  -0.040035893
                                  5686.760236158 -0.00001
                                                           0.99999
HM5
                   1.077214689
                                 29490.873608620
                                                  0.00004
                                                           0.99997
AM1
                   0.378240118
                                 11943.243082862
                                                  0.00003
                                                           0.99997
AM2
                   0.098414596
                                  6892.930782557
                                                  0.00001
                                                           0.99999
AM3
                   0.095207671
                                  5984.354733471
                                                  0.00002
                                                           0.99999
AM4
                  -0.007577885
                                  5772.348839886
                                                  0.00000
                                                           1.00000
AM5
                   0.100898393
                                  7202.061698739
                                                  0.00001
                                                           0.99999
HomeTeamLP
                  -0.017877180
                                  7308.578551302
                                                  0.00000
                                                           1.00000
AwayTeamLP
```

```
31448.781734635 -0.00001
                                                            0.99999
MW
                  -0.391787353
HTFormPtsStr
                  -0.890941115
                                 42546.813305330 -0.00002
                                                            0.99998
                                 42384.900918111 -0.00004
ATFormPtsStr
                  -1.813896947
                                                            0.99997
HTFormPts
                  -0.022660524
                                 22486.617562129
                                                  0.00000
                                                            1.00000
ATFormPts
                  -0.209716234
                                 21733.339220800 -0.00001
                                                            0.99999
HTWinStreak3
                  -0.016772026
                                 10169.383971325
                                                  0.00000
                                                            1.00000
HTWinStreak5
                   0.130261856
                                 16493.923311396
                                                  0.00001
                                                            0.99999
HTLossStreak3
                   0.009183538
                                  9406.991368763
                                                  0.00000
                                                            1.00000
HTLossStreak5
                   0.216805324
                                 19411.604241708
                                                  0.00001
                                                            0.99999
ATWinStreak3
                   0.098532600
                                 10023.389944058
                                                  0.00001
                                                            0.99999
ATWinStreak5
                  -0.293865560
                                 17866.714842300 -0.00002
                                                            0.99999
ATLossStreak3
                  -0.067678888
                                  9624.037581689 -0.00001
                                                            0.99999
ATLossStreak5
                   0.095418803
                                 18946.938573252
                                                  0.00001
                                                            1.00000
                  -0.020878120
HTGD
                                 53909.889815444
                                                  0.00000
                                                            1.00000
                  -0.546292650
                                 52997.739892298 -0.00001
ATGD
                                                            0.99999
DiffPts
                             NA
                                              NA
                                                        NA
                                                                 NA
                  -0.569780309
                                 54118.336747709 -0.00001
DiffFormPts
                                                            0.99999
DiffLP
                             NA
                                              NA
                                                        NA
                                                                 NA
(Dispersion parameter for binomial family taken to be 1)
    Null deviance: 8395.6292532105308
                                        on 6079
                                                 degrees of freedom
Residual deviance:
                      0.0000002440788
                                        on 6039
                                                 degrees of freedom
AIC: 82
Number of Fisher Scoring iterations: 25
> alm.probs=predict(alm.fits.test.type="response")
Warning message:
In predict.lm(object, newdata, se.fit, scale = 1, type = ifelse(type ==
  prediction from a rank-deficient fit may be misleading
> glm.probs = ifelse(glm.probs>0.5,1,0)
> table(glm.probs,test[,7])
alm.probs
        0 157
              49
        1
            0 174
> mean(glm.probs==test[,7])
[1] 0.8710526316
3.2 LDA:
> dataset = read.csv("D:/SEMESTER V/B1_ECM2002/J/Football-Data-Analysis-an
d-Prediction-master/Football-Data-Analysis-and-Prediction-master/Datasets/
final_dataset.csv")
> test = read.csv("D:/SEMESTER V/B1_ECM2002/J/Football-Data-Analysis-and-P
rediction-master/Football-Data-Analysis-and-Prediction-master/Datasets/tes
t.csv")
> library(MASS)
 #lda.fit = lda(FTR ~ .,data=dataset)
> lda.fit = lda(FTR ~ HomeTeam + AwayTeam + DiffPts+DiffLP+ DiffFormPts,da
ta=dataset)
> summary(lda.fit)
        Length Class
                      Mode
               -none- numeric
prior
          2
          2
counts
               -none- numeric
        170
means
               -none- numeric
               -none- numeric
scaling
         85
          2
lev
               -none- character
               -none- numeric
svd
          1
          1
               -none- numeric
Ν
call
          3
               -none- call
```

```
3
terms
             terms call
              -none- list
xlevels
         2
> lda.pred=predict(lda.fit, test[-7,])
> names(lda.pred)
[1] "class"
               "posterior" "x"
> lda.class=predict(lda.fit,test)$class
> test[,7]
  NH NH
NH NH NH NH NH NH NH NH NH NH H
 [36] NH NH NH H
                Н
                   Н
                      NH H
                                    NH H NH H
                                               NH H
                                                     Н
                                                        NH NH NH NH
NH H
    NH H
           Н
             NH H
                   Н
                      Н
                         NH NH NH
 [71] NH NH NH H
                Н
                   NH NH H
                            NH NH H
                                    NH NH H
                                            NH NH NH H
                                                        Н
                                                           NH H
                                                                NH H
NH NH H
        NH NH NH H
                   NH NH H
                            NH NH
[106] н
        NH H
             Н
                Н
                   NH H
                            Н
                               Н
                                 NH NH NH NH H
                                                  Н
                                                     NH NH NH NH H
H NH NH NH H
              NH H
                   Н
                      Н
                         NH NH NH
[141] H
        NH NH NH H
                   NH H
                         NH H
                               NH H H
                                       Н
                                          NH NH NH NH NH H
NH NH NH NH H
                Н
                   NH NH NH H
                               Н
[176] NH H
           Н
             Н
                Н
                   Н
                      NH NH NH H
                                 NH H
                                       Н
                                          NH NH H
                                                  NH H
                                                        Н
                                                           Н
                                                             NH H
NH NH H
        NH NH NH NH
                   NH H
                         Н
                            NH NH
[211] NH H
           NH H
                Н
                   Н
                      Н
                         NH NH H
                                 NH H
                                       NH NH NH H
                                                  NH NH NH NH NH H
H NH NH NH H
             Н
                NH
                   Н
                      NH NH H
                              Н
[246] NH NH H
             NH NH NH H
                         NH NH NH H
                                          NH NH H
                                                  NH H
                                                        NH H
     NH NH NH H
                   Н
                      н
                         Н
                            NH NH
[281] H
        NH NH H
                NH NH NH H
                            н
                              NH NH NH H
                                          NH NH NH H
                                                     NH NH NH H
H NH NH H
             NH NH
                   NH H
                         Н
                            NH H
[316] н
        н
           н
             NH NH H
                      NH H
                            NH NH NH H
                                       Н
                                          NH NH NH NH NH NH H
H NH NH NH H
             Н
                NH NH H
                            Н
                              NH
[351] н
       NH NH H
                Н
                   NH NH NH H
                              Н
                                 NH H NH H NH NH H
                                                     NH H
                                                           Н
                                                                NH H
H H H NH NH H
Levels: H NH
> table(lda.class,test[,7])
lda.class
           Н
             NH
      Н
          75
             60
      NH
          82 163
> mean(lda.class==test[,7])
[1] 0.6263157895
```

#### 3.3 QDA:

```
> dataset = read.csv("D:/SEMESTER V/B1_ECM2002/J/Football-Data-Analysis-an
d-Prediction-master/Football-Data-Analysis-and-Prediction-master/Datasets/
final dataset.csv")
> test = read.csv("D:/SEMESTER V/B1_ECM2002/J/Football-Data-Analysis-and-P
rediction-master/Football-Data-Analysis-and-Prediction-master/Datasets/tes
t.csv")
> library(MASS)
> qda.fit = qda(FTR ~ HomeTeam + AwayTeam + DiffPts+DiffLP+ DiffFormPts,da
ta=dataset)
> summary(qda.fit)
        Length Class Mode
prior
              -none- numeric
            2
            2
               -none- numeric
counts
          170
              -none- numeric
means
              -none- numeric
scaling 14450
               -none- numeric
ldet
            2
lev
               -none- character
            2
               -none- numeric
Ν
            1
               -none- call
call
            3
            3 terms call
terms
```

#### 3.4 KNN:

```
> library(class)
> dataset = read.csv("D:/SEMESTER V/B1_ECM2002/J/Football-Data-Analysis-an
d-Prediction-master/Football-Data-Analysis-and-Prediction-master/Datasets/
final_dataset.csv")
> test = read.csv("D:/SEMESTER V/B1_ECM2002/J/Football-Data-Analysis-and-P
rediction-master/Football-Data-Analysis-and-Prediction-master/Datasets/tes
t.csv")
> attach(dataset)
> library(clusterSim)
> dataset = data.frame(lapply(dataset, function(x) as.numeric(x)))
> dataset = data.Normalization (dataset,type="n4",normalization="column")
> test = data.frame(lapply(test, function(x) as.numeric(x)))
> test = data.Normalization (test,type="n4",normalization="column")
> train.X=cbind(dataset[])
> attach(test)
> test.X=cbind(test[])
> train.FTR=dataset[,7]
> set.seed(1)
> knn.pred=knn(train.X,test.X,train.FTR,k=1)
> table(knn.pred,test[,7])
knn.pred
       0 157
           0 223
       1
> mean(knn.pred==test[,7])
[1] 1
> knn.pred=knn(train.X,test.X,train.FTR,k=100)
> mean(knn.pred==test[,7])
[1] 0.9947368421
> knn.pred=knn(train.X,test.X,train.FTR,k=200)
> mean(knn.pred==test[,7])
[1] 0.9973684211
> knn.pred=knn(train.X,test.X,train.FTR,k=250)
> mean(knn.pred==test[,7])
[1] 1
> knn.pred=knn(train.X,test.X,train.FTR,k=300)
> mean(knn.pred==test[,7])
[1] 1
```

# 3.5 NAÏVE BAYES:

```
> dataset = read.csv("D:/SEMESTER V/B1_ECM2002/J/Football-Data-Analysis-an
d-Prediction-master/Football-Data-Analysis-and-Prediction-master/Datasets/
final_dataset.csv")
> test = read.csv("D:/SEMESTER V/B1_ECM2002/J/Football-Data-Analysis-and-P
rediction-master/Football-Data-Analysis-and-Prediction-master/Datasets/tes
t.csv")
> library('varhandle')
> library('e1071',warn.conflicts=FALSE)
> naive_bayes_model<-naiveBayes(FTR ~ ., data = dataset)</pre>
> naive_bayes_predictions<-predict(naive_bayes_model, newdata=test)</pre>
> naive_bayes_predictions
  Н
                                                                Н
                                                                      Н
NH H
 [36] NH NH NH H
                                   NH NH H
                                                  NH H
                                                              NH NH NH H
                 NH H
                       NH H
                                Н
                                            NH H
                                                        NH H
NH H
           Н
              NH H
                    NH NH H
     NH H
                             NH NH
[71] NH NH NH H
                    NH NH H
                                      NH NH H
                                                              NH H
                 Н
                             NH NH H
                                                  Н
                                                     NH H
                                                           Н
                                                                   NH H
NH H
        NH NH H
                 NH NH NH H
     Н
                             NH NH
[106] NH NH H
                    NH H
                             NH H
                                   NH NH H
              Н
                 Н
                          Н
                                            NH H
                                                  Н
                                                    Н
                                                        NH NH H
                                                                   NH NH
н н
     NH NH H
                    NH H
              NH H
                          NH H
                                NH
[141] H
        NH H
                 NH NH H
                          NH NH NH H
                                         Н
                                            NH NH NH H
                                                           Н
                                                              Н
                                                                Н
                                                                   NH H
              Н
H NH NH NH H
                    NH NH H
                 Н
                             NH H
[176] NH NH H
                          NH NH H
              NH H
                    Н
                       Н
                                      Н
                                         Н
                                            NH NH H
                                                        Н
                                                          NH NH NH H
                                   Н
                                                     Н
NH NH NH NH H
                       Н
                 NH H
                          Н
                             NH H
[211] NH H
                 NH
                       Н
              Н
                    Н
                          NH NH H
                                   NH H
                                         Н
                                            NH NH H
                                                        NH H
                                                              NH H
                                                                   NH H
           Н
                                                     Н
H NH NH NH H
                       Н
              Н
                 NH H
                          NH NH H
[246] NH H
              NH NH NH H
                          NH H
                                NH NH NH H
                                            NH NH H
                                                     NH H
                                                              Н
           Н
                                                           Н
                                                                      Н
                    Н
NH NH H
        NH NH H
                 Н
                       NH NH H
                                NH
[281] н
        NH H
              Н
                 Н
                    NH NH NH H
                                NH NH NH H
                                               NH NH H
                                                              NH NH NH H
                                            Н
                                                        NH H
H NH NH H
              Н
                 NH H
                       Н
           Н
                          Н
                             NH H
                                NH NH H
[316] н
           NH NH H
                       NH H
                             Н
        Н
                    Н
                                         Н
                                            NH H
                                                  Н
                                                     NH H
                                                           NH
                                                             NH H
                                                                      Н
H NH NH H
           H H NH H H
                            NH H
[351] H NH NH NH NH NH NH NH H
                                   NH NH NH NH NH NH H
                                Н
                                                          Н
                                                              Н
                                                                   NH H
H H NH NH NH H
Levels: H NH
> table(naive_bayes_predictions,test[,7])
naive_bayes_predictions
                            NH
                         Н
                       121
                            56
                    Н
                        36 167
                    NH
> mean(naive_bayes_predictions==test[,7])
[1] 0.7578947368
```

# 4.1 SVM:

```
> dataset = read.csv("D:/SEMESTER V/B1_ECM2002/J/Football-Data-Analysis-an
d-Prediction-master/Football-Data-Analysis-and-Prediction-master/Datasets/
final_dataset.csv")
> test = read.csv("D:/SEMESTER V/B1_ECM2002/J/Football-Data-Analysis-and-P
rediction-master/Football-Data-Analysis-and-Prediction-master/Datasets/tes
t.csv")
> library("e1071")
> library(caret)
> options(repos = c(CRAN = "http://cran.rstudio.com"))
> set.seed(319)
> ind <- sample(2, nrow(dataset), replace = TRUE, prob = c(0.90, 0.10))
> training = dataset[ind==1,]
> testing = dataset[ind==2,]
> dim(training)
[1] 5449
          43
> dim(testing)
[1] 631 43
> svm_model1 <- svm(FTR ~ .,kernel = "radial", data=training)</pre>
> summary(svm_model1)
call:
svm(formula = FTR ~ ., data = training, kernel = "radial")
Parameters:
   SVM-Type: C-classification
 SVM-Kernel:
              radial
       cost:
              0.0004103406
      gamma:
Number of Support Vectors: 3265
 ( 1632 1633 )
Number of Classes: 2
Levels:
 H NH
> pred1 <- predict(svm_model1, newdata = testing)</pre>
> confusionMatrix(pred1, testing$FTR )
Confusion Matrix and Statistics
          Reference
Prediction H NH
        н 269
        NH 25 329
               Accuracy : 0.9477
                 95% CI: (0.9273, 0.9637)
    No Information Rate: 0.5341
    P-Value [Acc > NIR] : < 2.2e-16
                  Kappa: 0.8945
 Mcnemar's Test P-Value: 0.005349
```

```
Sensitivity: 0.9150
            Specificity: 0.9763
         Pos Pred Value : 0.9711
         Neg Pred Value: 0.9294
            Prevalence: 0.4659
         Detection Rate: 0.4263
   Detection Prevalence: 0.4390
      Balanced Accuracy: 0.9456
       'Positive' Class: H
> svm_model2 <- svm(FTR ~ .,kernel = "linear", data=training)</pre>
> summary(svm_model2)
svm(formula = FTR ~ ., data = training, kernel = "linear")
Parameters:
   SVM-Type: C-classification
 SVM-Kernel:
             linear
       cost:
             0.0004103406
      gamma:
Number of Support Vectors: 1512
 (724 788)
Number of Classes: 2
Levels:
H NH
> pred2 <- predict(svm_model2, newdata = testing)</pre>
> confusionMatrix(pred2, testing$FTR )
Confusion Matrix and Statistics
          Reference
Prediction H NH
        н 294
        NH 0 337
              Accuracy : 1
                95% CI: (0.9942, 1)
    No Information Rate : 0.5341
    P-Value [Acc > NIR] : < 2.2e-16
                  Kappa: 1
 Mcnemar's Test P-Value: NA
           Sensitivity: 1.0000
           Specificity: 1.0000
         Pos Pred Value : 1.0000
         Neg Pred Value: 1.0000
             Prevalence: 0.4659
         Detection Rate: 0.4659
   Detection Prevalence: 0.4659
      Balanced Accuracy: 1.0000
```

```
'Positive' Class: H
> svm_model3 <- svm(FTR ~ .,kernel = "polynomial",degree = 3, data=trainin</pre>
g)
> summary(svm_model3)
svm(formula = FTR ~ ., data = training, kernel = "polynomial", degree = 3)
Parameters:
   SVM-Type: C-classification
 SVM-Kernel:
              polynomial
       cost:
     degree:
      gamma:
             0.0004103406
     coef.0:
Number of Support Vectors: 5044
 ( 2522 2522 )
Number of Classes: 2
Levels:
H NH
> pred3 <- predict(svm_model3, newdata = testing)</pre>
> confusionMatrix(pred3, testing$FTR )
Confusion Matrix and Statistics
          Reference
Prediction H NH
             0
        NH 294 337
               Accuracy : 0.5341
                 95% CI: (0.4943, 0.5735)
    No Information Rate : 0.5341
    P-Value [Acc > NIR] : 0.5163
                  карра: 0
 Mcnemar's Test P-Value : <2e-16
            Sensitivity: 0.0000
            Specificity: 1.0000
         Pos Pred Value:
         Neg Pred Value: 0.5341
             Prevalence: 0.4659
         Detection Rate: 0.0000
   Detection Prevalence: 0.0000
      Balanced Accuracy: 0.5000
       'Positive' Class: H
 svm_model4 <- svm(FTR ~ .,kernel ="polynomial",degree = 5, data=training</pre>
> summary(svm_model4)
```

```
call:
svm(formula = FTR ~ ., data = training, kernel = "polynomial", degree = 5)
Parameters:
  SVM-Type: C-classification
 SVM-Kernel:
             polynomial
       cost:
     degree:
     gamma: 0.0004103406
     coef.0: 0
Number of Support Vectors: 5044
 ( 2522 2522 )
Number of Classes: 2
Levels:
H NH
> pred4 <- predict(svm_model4, newdata = testing)</pre>
> confusionMatrix(pred4, testing$FTR )
Confusion Matrix and Statistics
         Reference
Prediction H NH
            0
        NH 294 337
              Accuracy : 0.5341
                95% CI: (0.4943, 0.5735)
    No Information Rate: 0.5341
    P-Value [Acc > NIR] : 0.5163
                 карра: 0
 Mcnemar's Test P-Value : <2e-16
           Sensitivity: 0.0000
           Specificity: 1.0000
         Pos Pred Value:
         Neg Pred Value : 0.5341
            Prevalence: 0.4659
         Detection Rate: 0.0000
   Detection Prevalence: 0.0000
      Balanced Accuracy: 0.5000
       'Positive' Class : H
```

#### **4.2 NEURAL NETWORKS:**

```
> dataset = read.csv("D:/SEMESTER V/B1_ECM2002/J/Football-Data-Analysis-an
d-Prediction-master/Football-Data-Analysis-and-Prediction-master/Datasets/
final dataset.csv")
> test = read.csv("D:/SEMESTER V/B1_ECM2002/J/Football-Data-Analysis-and-P
rediction-master/Football-Data-Analysis-and-Prediction-master/Datasets/tes
t.csv")
> #Pre-Processing
> str(dataset)
               6080 obs. of 43 variables:
'data.frame':
               : int 0 1 2 3 4 5 6 7 8 9 ...
$ X
                : Factor w/ 1592 levels "2000-08-19", "2000-08-20", ...: 1 1
$ Date
1 1 1 1 1 1 1 2 ...
$ HomeTeam
                : Factor w/ 42 levels "Arsenal", "Aston Villa", ...: 11 12 13
15 20 21 22 35 37 24 ...
$ AwayTeam
               : Factor w/ 42 levels "Arsenal", "Aston Villa", ...: 23 40 26
33 16 2 8 1 19 27 ...
                      4 4 1 2 2 0 1 1 3 2 ...
$ FTHG
               : int
$ FTAG
                : int 0 2 3 2 0 0 0 0 1 0 ...
                : Factor w/ 2 levels "H", "NH": 1 1 2 2 1 2 1 1 1 1 ...
$ FTR
                : int 0000000000...
$ HTGS
                      0 0 0 0 0 0 0 0 0 0 ...
$ ATGS
               : int
                      0 0 0 0 0 0 0 0 0 0 ...
$ HTGC
               : int
                      00000000000...
               : int
$ ATGC
                      0000000000...
$ HTP
               : num
               : num 0 0 0 0 0 0 0 0 0
$ ATP
               : Factor w/ 4 levels "D", "L", "M", "W": 3 3 3 3 3 3 3 3 3 3
$ HM1
               : Factor w/ 4 levels "D", "L", "M", "W": 3 3 3 3 3 3 3 3 3 3
$ HM2
               : Factor w/ 4 levels "D", "L", "M", "W": 3 3 3 3 3 3 3 3 3 3
 $ HM3
                : Factor w/ 4 levels "D", "L", "M", "W": 3 3 3 3 3 3 3 3 3 3
 $ HM4
                : Factor w/ 4 levels "D", "L", "M", "W": 3 3 3 3 3 3 3 3 3 3
 $ HM5
                : Factor w/ 4 levels "D", "L", "M", "W": 3 3 3 3 3 3 3 3 3 3
 $ AM1
                : Factor w/ 4 levels "D", "L", "M", "W": 3 3 3 3 3 3 3 3 3 3
 $ AM2
                : Factor w/ 4 levels "D", "L", "M", "W": 3 3 3 3 3 3 3 3 3 3
 $ AM3
               : Factor w/ 4 levels "D", "L", "M", "W": 3 3 3 3 3 3 3 3 3 3
 $ AM4
                : Factor w/ 4 levels "D", "L", "M", "W": 3 3 3 3 3 3 3 3 3 3
 $ AM5
 $ HomeTeamLP
               : num
                      18 5 14 16 3 8 4 7 10 1
                      18 9 12 15 13 6 17 2 18 11 ...
$ AwayTeamLP
                : num
                      1 1 1 1 1 1 1 1 1 1 . .
                : num
$ HTFormPtsStr : Factor w/ 349 levels "DDDDD", "DDDDL", ...: 230 230 230
230 230 230 230 230 ...
 $ ATFormPtsStr : Factor w/ 359 levels "DDDDD", "DDDDL",..: 241 241 241
241 241 241 241 241 ...
               : int 0000000000...
 $ HTFormPts
                : int
 $ ATFormPts
                      0 0 0 0 0 0 0 0 0 0 ...
                      00000000000...
$ HTWinStreak3 : int
 $ HTWinStreak5 : int
                      00000000000...
 $ HTLossStreak3: int
                      00000000000...
 $ HTLossStreak5: int
                      00000000000...
$ ATWinStreak3 : int 0000000000...
```

```
$ ATWinStreak5 : int 0 0 0 0 0 0 0 0 0 ...
  ATLossStreak3: int
                       0 0 0 0 0 0 0 0 0 0 ...
  ATLossStreak5: int
                       0 0 0 0 0 0 0 0 0
 $ HTGD
                  num
                       0 0 0 0 0 0 0 0 0 0 ...
 $
  ATGD
                  num
                       0000000000...
  DiffPts
                : num
                       0 0 0 0 0 0 0 0 0 0 ...
                       0 0 0 0 0 0 0 0 0 0 ...
 $ DiffFormPts
                : num
                       0 -4 2 1 -10 2 -13 5 -8 -10 ...
 $ DiffLP
                : num
> library(clusterSim)
> dataset = data.frame(lapply(dataset, function(x) as.numeric(x)))
> dataset = data.Normalization (dataset,type="n4",normalization="column")
> test = data.frame(lapply(test, function(x) as.numeric(x)))
> test = data.Normalization (test,type="n4",normalization="column")
 # Neural Networks
> library(neuralnet)
> set.seed(319)
> # 5 neurons hidden layer
> n = neuralnet(FTR ~ X + Date + HomeTeam + AwayTeam + FTHG + FTAG + HTGS
+ ATGS +
                  \mathsf{HTGC} + \mathsf{ATGC} + \mathsf{HTP} + \mathsf{ATP} + \mathsf{HM1} + \mathsf{HM2} + \mathsf{HM3} + \mathsf{HM4} + \mathsf{HM5} +
AM1 +
                  AM2 + AM3 + AM4 + AM5 + HomeTeamLP + AwayTeamLP + MW + H
TFormPtsStr +
                  ATFORMPTSStr + HTFORMPTS + ATFORMPTS + HTWinStreak3 + HT
WinStreak5 +
                  HTLossStreak3 + HTLossStreak5 + ATWinStreak3 + ATWinStre
ak5 +
                  ATLossStreak3 + ATLossStreak5 + HTGD + ATGD + DiffPts +
DiffFormPts +
                  DiffLP, data = dataset, hidden = 5, err.fct = "ce", linear.o
utput = FALSE)
> # Prediction
> output <- compute(n, dataset[,-7])</pre>
> head(output$net.result)
1 0.00000000005496020571
2 0.00000000094983300103
4 0.999999843402808674675
5 0.00000000065638662296
6 0.999999795743205188714
> head(dataset[1,])
 X Date
            HomeTeam
                         AwayTeam
                                           FTHG FTAG FTR HTGS ATGS HTGC ATG
C HTP ATP
                   HM1
                                 нм2
                                              нм3
1 0
       0 0.243902439 0.5365853659 0.444444444
                                                        0
                                                             0
                                                                  0
                                                                        0
        0 0.666666667 0.666666667 0.6666666667
                        HM5
                                                    AM2
                                                                 AM3
             AM5 HomeTeamLP AwayTeamLP MW
AM4
1 0.666666667 0.666666667 0.666666667 0.666666667 0.666666667 0.66666
66667 0.6666666667
                             1
                                        1
  HTFormPtsStr ATFormPts ATFormPts HTWinStreak3 HTWinStreak5
HTLossStreak3 HTLossStreak5 ATWinStreak3
  0.658045977 0.6703910615
                                                                           0
  ATWinStreak5 ATLossStreak3 ATLossStreak5
                                                     HTGD
                                                                 ATGD
                                                                            D
iffPts DiffFormPts DiffLP
                                          0 0.4285714286 0.487804878 0.5123
152709
               0.5
                      0.5
> # Confusion Matrix & Misclassification Error - training data
> output <- compute(n, dataset[,-7])</pre>
> p1 <- output$net.result</pre>
> pred1 <- ifelse(p1>0.5, 1, 0)
```

```
> tab1 <- table(pred1, dataset$FTR)</pre>
> tab1
pred1
          0
                1
    0 2816
                0
          0 3264
    1
> sum(diag(tab1))/sum(tab1)
> # Confusion Matrix & Misclassification Error - testing data
> output <- compute(n, test[,-7])</pre>
> p2 <- output$net.result</pre>
> pred2 <- ifelse(p2>0.5, 1, 0)
> tab2 <- table(pred2, test$FTR)</pre>
pred2
       0
           51
    0 157
        0 172
    1
> sum(diag(tab2))/sum(tab2)
[1] 0.8657894737
```

### **5.1: COMPARISON OF ALL MODELS:**

MODEL	TEST
	ACCURACY
LOGISTIC REGRESSION	87.10%
LDA	62.63%
QDA	59.21%
KNN	100%
NAÏVE BAYES	75.78%
SVM	100%
NEURAL NETWOKRS	86.57%

#### **5.2: CONCLUSION:**

- Home teams have a definite advantage over Away teams. On aggregate, home team win 46.65% matches compared to 27.72% matches won by the away teams.
- But Football is an unpredictable affair. The English Premier League is known for its famous upsets by bogey teams i.e. a consistently lower ranked team beating a higher ranked team.
- Despite all the facts, an attempt has been made by us to use statistical models to predict whether a home team will win or not given any match.
- Models like K Nearest Neighbours and Support Vector Machines with linear kernel have been able to predict correctly the 380 matches in test data, whereas other models have their test accuracy in the range of 60-90%.