# Term Project Report

DA5020 Collect/Store/Retrieve Data

Harsh Desai
MS in Bioinformatics (Spring 2017)

# Premier League 2016-17 Player Database

Motivation: I am a Bioinformatics student, so ideally I should select any biological database for the final term project. However, biological data is big responsibility to work with and I have worked on biological data during my other courses. So, my final choice was to work with sports data and I have decided to work on Premier League 2016-17 player and stats database.

Introduction:

Premier League or English Premier League is a major soccer league contested between top 20 teams in England every year. It is the one of the most followed and major soccer league in the world. For the final term project, I have decided to collect data for the Premier League 2016-17 season, which is also current running season. I have tried to collect most of the possible and significant data for the season.

Work Approach:

1)Collecting 2) Storing 3) Retrieval 4) Graphical Representation

## 1)Collecting:

For the collection of the data I will be scraping directly from the www.premierleague.com and www.nbcsports.com

I am using import.io as a scraping tools and URL get query to collect data team wise using team id as URL query.

For example: http://scores.nbcsports.com/epl/teamstats.asp?team=21, here team=21 is a team query for the team called Arsenal. I have extracted data using import.io tool using URL as input.



Import.io allow us to store collected data into .csv files. So, I have stored 20 .csv file for each team separately.

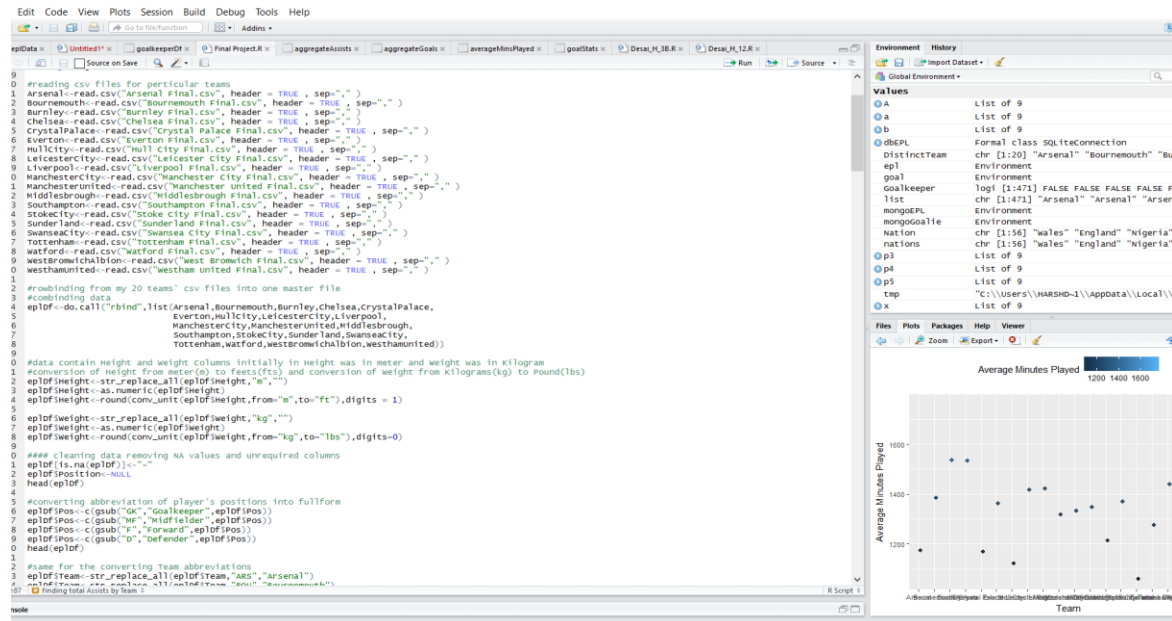Reading data into R using read.csv() function and combining all 20 files into master file using R functions. I might have combined data manually but I wanted to use R functions that's the main I choose to extract file per teams and not a single master file.

Packages used: stringr ,measurements

Loading .csv files, combining files and Cleaning data

```
#reading csv files for perticular teams
Arsenal<-read.csv("Arsenal Final.csv", header = TRUE , sep=",")
Bournemouth<-read.csv("Bournemouth Final.csv", header = TRUE , sep=",")
Burnley<-read.csv("Burnley Final.csv", header = TRUE , sep=",")
Chelsea<-read.csv("Chelsea Final.csv", header = TRUE , sep=",")
CrystalPalace<-read.csv("Crystal Palace Final.csv", header = TRUE , sep=",")
Everton<-read.csv("Everton Final.csv", header = TRUE , sep=",")
HullCity<-read.csv("Hull City Final.csv", header = TRUE , sep=",")
LeicesterCity<-read.csv("Leicester City Final.csv", header = TRUE , sep=",")
Liverpool<-read.csv("Liverpool Final.csv", header = TRUE , sep=",")
ManchesterCity<-read.csv("Manchester City Final.csv", header = TRUE , sep=",")
Manchesterunited<-read.csv("Manchester United Final.csv", header = TRUE , sep=",")
Middlesbrough<-read.csv("Middlesbrough Final.csv", header = TRUE , sep=",")
Southampton<-read.csv("Southampton Final.csv", header = TRUE , sep=",")
StokeCity<-read.csv("Stoke City Final.csv", header = TRUE , sep=",")
Sunderland<-read.csv("Sunderland Final.csv", header = TRUE , sep=",")
SwanseaCity<-read.csv("Swansea City Final.csv", header = TRUE , sep=",")
Tottenham<-read.csv("Tottenham Final.csv", header = TRUE , sep=",")
Watford<-read.csv("Watford Final.csv", header = TRUE , sep=",")
WestBromwichAlbion<-read.csv("West Bromwich Final.csv", header = TRUE , sep=",")
westhamunited<-read.csv("westham united Final.csv", header = TRUE , sep=",")

#rowbinding from my 20 teams' csv files into one master file
#combinding data
eplDf<-do.call("rbind",list(Arsenal,Bournemouth,Burnley,chelsea,CrystalPalace,
                Everton,HullCity,LeicesterCity,Liverpool,
                ManchesterCity,Manchesterunited,Middlesbrough,
                Southampton,StokeCity,Sunderland,SwanseaCity,
                Tottenham,Watford,WestBromwichAlbion,westhamunited))

#data contain Height and weight columns initially in Height was in meter and weight was in kilogram
#conversion of Height from meter(m) to feets(fts) and conversion of weight from Kilograms(kg) to Pound(lbs)
eplDf$Height<-str_replace_all(eplDf$Height,"m","")
eplDf$Height<-as.numeric(eplDf$Height)
eplDf$Height<-round(conv_unit(eplDf$Height,from="m",to="ft"),digits = 1)

eplDf$weight<-str_replace_all(eplDf$weight,"kg","")
eplDf$weight<-as.numeric(eplDf$weight)
eplDf$weight<-round(conv_unit(eplDf$weight,from="kg",to="lbs"),digits=0)

#### cleaning data removing NA values and unrequired columns
eplDf[is.na(eplDf)]<-"-"
eplDf$Position<-NULL
head(eplDf)

#converting abbreviation of player's positions into fullform
eplDf$Pos<-c(gsub("GK","Goalkeeper",eplDf$Pos))
eplDf$Pos<-c(gsub("MF","Midfielder",eplDf$Pos))
eplDf$Pos<-c(gsub("F","Forward",eplDf$Pos))
eplDf$Pos<-c(gsub("D","Defender",eplDf$Pos))
head(eplDf)

#same for the converting Team abbreviations
eplDf$Team<-str_replace_all(eplDf$Team,"ARS","Arsenal")
eplDf$Team<-str_replace_all(eplDf$Team,"BOU","Bournemouth")
```

Combined final data frame :

| Name | Position | Matches Played | Minutes Played | Goals | Assists | Penalties Scored | Shots | Shots on Target | Shots on Target % | Yellow Cards | Red Cards | Fouls | Fouls Suffered | Crosses | Offsi |
|------|----------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 1 Aaron Ramsey | Midfielder | 16 | 772 | 0 | 2 | 0 | 29 | 6 | 20.7 | 3 | 0 | 11 | 15 | 24 | |
| 2 Ainsley Maitland-Niles | Midfielder | 1 | 1 | 0 | 0 | 0 | 0 | 0 | - | 0 | 0 | 0 | 0 | 0 | |
| 3 Alex Iwobi | Forward | 24 | 1419 | 3 | 4 | 0 | 35 | 9 | 25.7 | 1 | 0 | 6 | 12 | 33 | |
| 4 Alex Oxlade-Chamberlain | Midfielder | 26 | 1352 | 2 | 5 | 0 | 26 | 7 | 26.9 | 1 | 0 | 23 | 16 | 74 | |
| 5 Alexis Sánchez | Forward | 31 | 2628 | 19 | 9 | 2 | 104 | 41 | 39.4 | 5 | 0 | 40 | 52 | 136 | |
| 6 Carl Jenkinson | Defender | 1 | 83 | 0 | 0 | 0 | 0 | 0 | - | 0 | 0 | 1 | 1 | 3 | |
| 7 Damián Martínez | Goalkeeper | 2 | 180 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | |
| 8 Danny Welbeck | Forward | 9 | 417 | 1 | 0 | 0 | 9 | 5 | 55.6 | 0 | 0 | 5 | 5 | 2 | |
| 9 David Ospina | Goalkeeper | 2 | 142 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | |
| 10 Francis Coquelin | Midfielder | 24 | 1577 | 0 | 0 | 0 | 12 | 2 | 16.7 | 5 | 0 | 33 | 34 | 9 | |
| 11 Gabriel | Defender | 16 | 1298 | 0 | 0 | 0 | 7 | 2 | 28.6 | 4 | 0 | 14 | 8 | 14 | |
| 12 Granit Xhaka | Midfielder | 25 | 1923 | 1 | 2 | 0 | 24 | 4 | 16.7 | 4 | 2 | 30 | 14 | 35 | |
| 13 Héctor Bellerín | Defender | 26 | 2064 | 0 | 2 | 0 | 17 | 4 | 23.5 | 2 | 0 | 17 | 11 | 75 | |
| 14 Kieran Gibbs | Defender | 6 | 295 | 0 | 1 | 0 | 3 | 2 | 66.7 | 3 | 0 | 7 | 1 | 19 | |
| 15 Laurent Koscielny | Defender | 28 | 2449 | 2 | 0 | 0 | 8 | 3 | 37.5 | 3 | 0 | 14 | 24 | 4 | |
| 16 Lucas Pérez | Forward | 11 | 265 | 1 | 0 | 0 | 9 | 4 | 44.4 | 0 | 0 | 3 | 2 | 11 | |
| 17 Mathieu Debuchy | Defender | 1 | 16 | 0 | 0 | 0 | 0 | 0 | - | 0 | 0 | 0 | 0 | 0 | |
| 18 Mesut Özil | Midfielder | 26 | 2225 | 7 | 7 | 0 | 36 | 15 | 41.7 | 0 | 0 | 8 | 20 | 194 | |
| 19 Mohamed Elneny | Midfielder | 14 | 694 | 0 | 1 | 0 | 14 | 2 | 14.3 | 1 | 0 | 10 | 5 | 10 | |
| 20 Nacho Monreal | Defender | 29 | 2523 | 0 | 2 | 0 | 9 | 2 | 22.2 | 3 | 0 | 26 | 33 | 78 | |
| 21 Olivier Giroud | Forward | 23 | 913 | 9 | 3 | 0 | 30 | 14 | 46.7 | 1 | 0 | 16 | 10 | 3 | |
| 22 Petr Cech | Goalkeeper | 28 | 2468 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 2 | 1 | 0 | |

Collected Variables:

Basic Player Info: Name, Position, Team, Matches Played, Minutes Played, Nationality etc

Standard Stats: Goals, Assists, Shots, Yellow Cards, Red Cards, Crosses, Offside etc.

Goalkeeper Stats: Saves, Saves, Penalty Kicks Saved, Goals Allowed, Shots faced etc.

## 2)Storing:

Initially, I was planning to store collected data using Relational Database. After working on relational database schema, I realized that it will be difficult to store data into relational database due to complexity between player positions and their statistics. For example, Goalkeeper statistics attributes are totally different than other player position's attribute for soccer. So, I decided to go with non-relational (NoSQL) database system to store collected data. It is simpler than relational database system, other advantages of NoSQL database system are no unique constrain, massive parallel processing and tolerance for failure. Non-relational system I have used for project is MongoDB.

Packages used: mongolite

Steps for running MongoDB server
1.installing MongoDB from mongodb.com
2.Setting up MongoDB environment
mkdir \data\db
"C:\Program Files\MongoDB\Server\3.4\bin\mongod.exe" --dbpath d:\test\mongodb\data
3. C:\MongoDB\Server\3.2\bin\mongod.exe   Code to start MongoDB server
3. Connect using R

After connecting MongoDB server with R studio using mongolite R package, next step is to insert collected and cleaned data into MongoDB system using insert() function provided my mongolite. Mongolite also provides option of exporting JSON structure to file using export() function.

Setting MongoDB connections, inserting files and exporting JSON structure to designated files.

## 3.Retrieval:

Functions used : sort(),find(),distinct()

Packages used: dplyr,ggplot2,mongolite

find() function is used to find specific value from the stored data.

Firstly using find() function to search all the forward position players from the database.



Using find() function to fetch multiple columns from the database and to make data frame

The other function used for retrieval was distinct(),which will find distinct values from particular columns.

Use of distinct() function to find distinct teams and players from different nations.
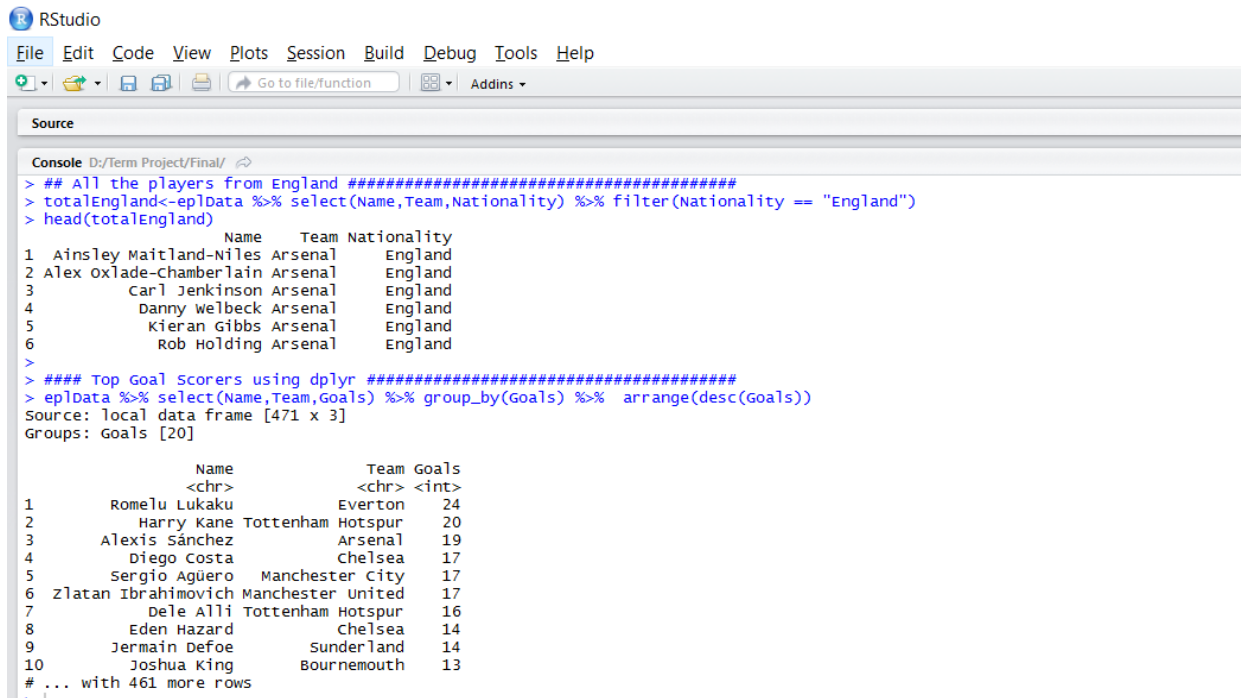


We can see that there are 20 distinct teams in a database and players are from 56 distinct countries.

Next function I have used for the retrieval is sort()function. I have used sort()function to find top 10 Goal scorers and top 10 Assists by players.



I have also used **dplyr package** to work on my retrieved data. I have used filter, group_by, select, arrange etc.

I have used dplyr package to find all the player from England and Top Goal Scorer for the season.
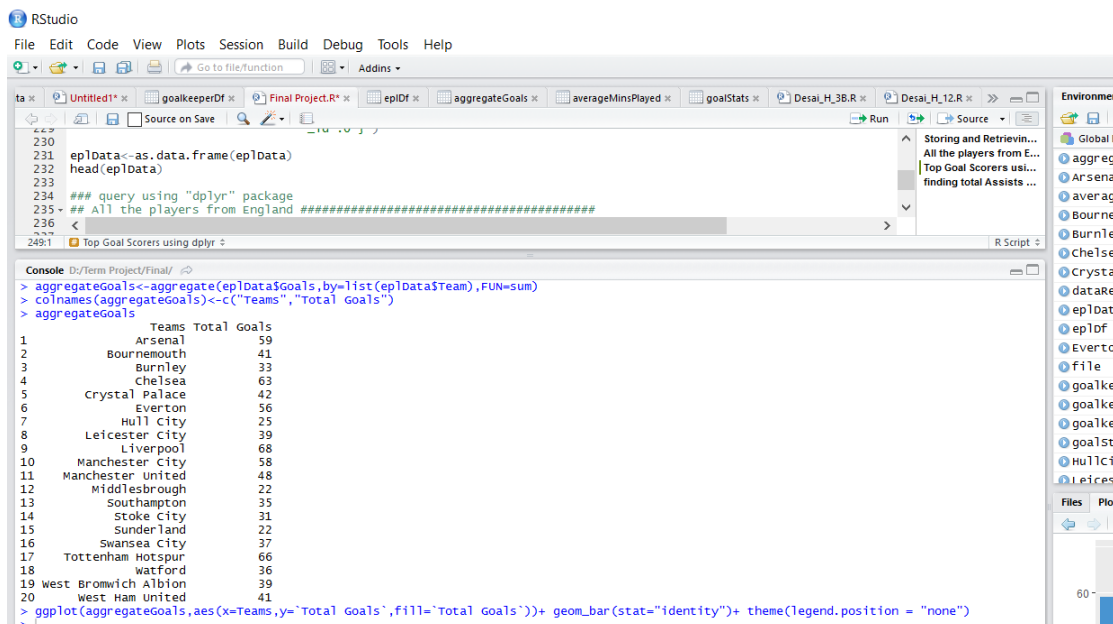
```
R RStudio

File  Edit  Code  View  Plots  Session  Build  Debug  Tools  Help

Source

Console D:/Term Project/Final/
> ## All the players from England ######################################
> totalEngland<-eplData %>% select(Name,Team,Nationality) %>% filter(Nationality == "England")
> head(totalEngland)
                   Name    Team Nationality
1  Ainsley Maitland-Niles Arsenal      England
2 Alex Oxlade-Chamberlain Arsenal      England
3          Carl Jenkinson Arsenal      England
4           Danny Welbeck Arsenal      England
5            Kieran Gibbs Arsenal      England
6            Rob Holding Arsenal      England
>
> #### Top Goal Scorers using dplyr #####################################
> eplData %>% select(Name,Team,Goals) %>% group_by(Goals) %>%  arrange(desc(Goals))
Source: local data frame [471 x 3]
Groups: Goals [20]

                  Name               Team Goals
                 <chr>              <chr> <int>
1         Romelu Lukaku            Everton    24
2           Harry Kane Tottenham Hotspur    20
3        Alexis Sánchez            Arsenal    19
4          Diego Costa            Chelsea    17
5        Sergio Agüero    Manchester City    17
6  Zlatan Ibrahimovich Manchester United    17
7            Dele Alli Tottenham Hotspur    16
8          Eden Hazard            Chelsea    14
9        Jermain Defoe          Sunderland    14
10          Joshua King        Bournemouth    13
# ... with 461 more rows
```

## 4.Graphical Representation:

Last but not the list, I have used ggplot2 package to present my data and results graphically.

Use of aggregate function to calculate total number of goals by each team and present data graphically using ggplot.
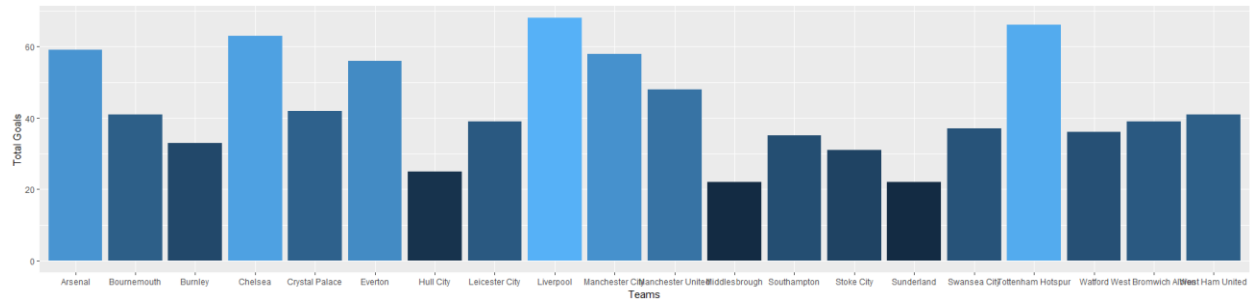
```
R RStudio

File  Edit  Code  View  Plots  Session  Build  Debug  Tools  Help

ta × | Untitled1* × | goalkeeperDf × | Final Project.R* × | eplDf × | aggregateGoals × | averageMinsPlayed × | goalStats × | Desai_H_3B.R × | Desai_H_12.R × | >>       Environmen

230
231  eplData<-as.data.frame(eplData)
232  head(eplData)
233
234  ### query using "dplyr" package
235 ## All the players from England ######################################
236

249:1  Top Goal Scorers using dplyr

Console D:/Term Project/Final/
> aggregateGoals<-aggregate(eplData$Goals,by=list(eplData$Team),FUN=sum)
> colnames(aggregateGoals)<-c("Teams","Total Goals")
> aggregateGoals
              Teams Total Goals
1            Arsenal          59
2        Bournemouth          41
3            Burnley          33
4            Chelsea          63
5      Crystal Palace          42
6            Everton          56
7          Hull City          25
8      Leicester City          39
9          Liverpool          68
10    Manchester City          58
11  Manchester United          48
12      Middlesbrough          22
13        Southampton          35
14         Stoke City          31
15         Sunderland          22
16        Swansea City          37
17  Tottenham Hotspur          66
18            Watford          36
19 West Bromwich Albion          39
20      West Ham United          41
> ggplot(aggregateGoals,aes(x=Teams,y=`Total Goals`,fill=`Total Goals`))+ geom_bar(stat="identity")+ theme(legend.position = "none")
>
```

Then I plotted abundance graph for goals and assists for by the team.





Presented two points in graph shows highest goals and highest assists and color for the team matches with our data. Top goal scorer Romelu Lukaku is from Everton and Top Assists player Kevin De Bruyne is from Manchester City.
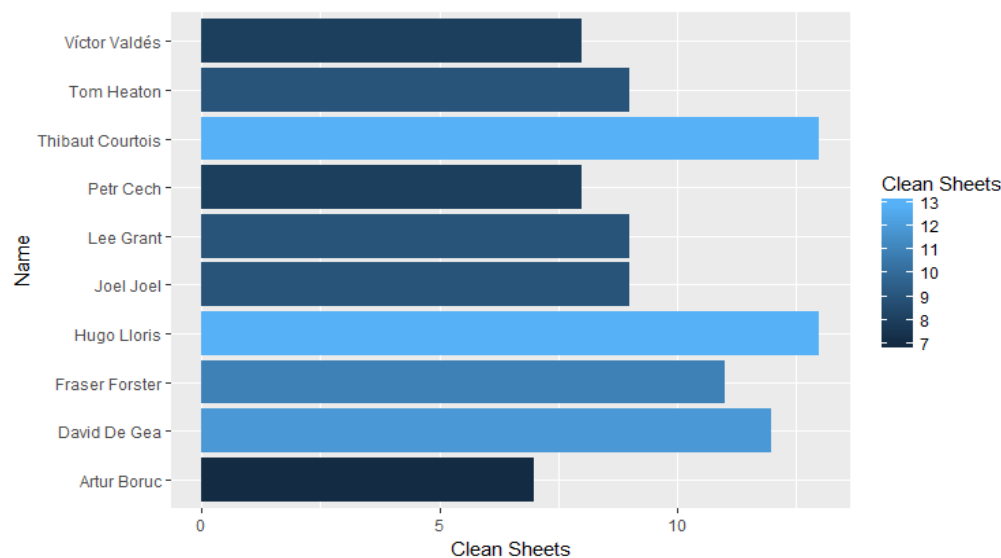
Other graph would be to find Top 10 Cleansheets by Goalkeeper and to plot it.



```r
> ## to find higest CleanSheet by GoalKeepers and Plot a graph
> TopCleanSheets<-GOALIE$find('{}',
+                             fields = '{"Name":1,
+                             "Team":1,
+                             "Clean Sheets":1,
+                             "_id":0}',
+                             sort='{"clean Sheets":-1}')
> TopCleanSheets<-TopCleanSheets[1:10,]
> TopCleanSheets$`Clean Sheets`<-as.numeric(TopCleanSheets$`Clean Sheets`)
> head(TopCleanSheets)
           Name              Team Clean Sheets
1 Thibaut Courtois         Chelsea           13
2    Hugo Lloris Tottenham Hotspur           13
3   David De Gea Manchester United           12
4  Fraser Forster      Southampton           11
5     Tom Heaton          Burnley            9
6      Joel Joel          Everton            9
>
> ggplot() +
+   layer(
+     data = TopCleanSheets, mapping = aes(x = Name, y = `Clean Sheets`,fill=`clean Sheets`),
+     geom = "bar", stat = "identity", position = "identity"
+   ) +
+   scale_y_continuous() +
+   coord_flip()
> |
```
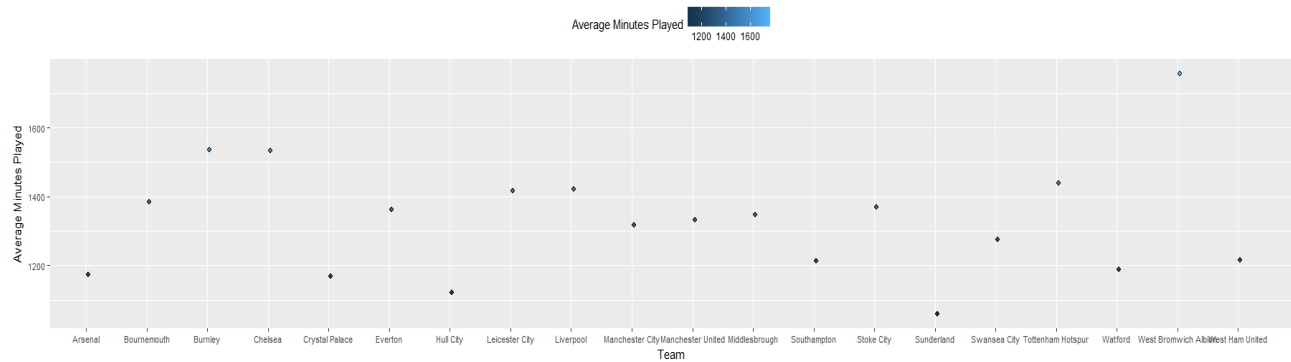


In the end I calculated average minutes played by every team and plotted a graph for it.



```r
> ## To find average minutes played by team and plot a graph to represent it graphically
>
> eplData$`Minutes Played`<-as.numeric(eplData$`Minutes Played`)
>
> averageMinsPlayed<-aggregate(eplData$`Minutes Played`,by=list(eplData$Team),FUN=mean)
> colnames(averageMinsPlayed)<-c("Team","Average Minutes Played")
> averageMinsPlayed
                Team Average Minutes Played
1            Arsenal               1173.769
2        Bournemouth               1383.217
3            Burnley               1536.524
4            Chelsea               1532.000
5     Crystal Palace               1167.259
6            Everton               1362.261
7           Hull City               1121.385
8     Leicester City               1416.000
9          Liverpool               1420.435
10   Manchester City               1316.708
11 Manchester United               1332.043
12     Middlesbrough               1347.348
13       Southampton               1213.583
14        Stoke City               1369.435
15        Sunderland               1059.214
16      Swansea City               1274.958
17 Tottenham Hotspur               1439.955
18           Watford               1187.577
19 West Bromwich Albion              1756.056
20     West Ham United               1215.520
>
> x<-ggplot(averageMinsPlayed,aes(x=Team,y=`Average Minutes Played`,fill=`Average Minutes Played`))
> x+geom_dotplot(binaxis = "y",binwidth = 15)+  theme(legend.position="top",
+                             axis.text=element_text(size = 8))
> |
```

Future Work:

- To collect more stats for the same seasons and previous seasons
- To use better and dynamic scraping technique
- To collect data for different soccer leagues like La Liga, Champions League etc.
- After collecting other Premier League Seasons data, produce comparative statistics.
- Better graphical representation and will try to create soccer statistics application.

References:

- Collecting, Storing and Retrieving Data by Yatish Jain and Martin Schedlbauer
- Data Manipulation with R by Jaynal Abedin
- Seven Database in Seven Weeks by Eric Redmond and Jim R. Wilson
- www.premierleague.com
- www.nbcsports.com
- http://ggplot2.org/

Acknowledgment:

I would like to thank ,Teaching Assistants for the course DA5020 Jianchao (Jesse) Yang and Japan Mehta for their continuous support throughout the course. My special thanks for the Dr. Durant for continuous guidance and knowledge during course.