

# CYO PROJECT: FIFA player overall 2022

Isaac Cyrman Casafont

3/31/2021

## Contents

I) Introduction .....	2
II) Data cleaning, filtering and selecting process .....	3
III) Data analysis process .....	5
IV) Create FIFA set (Training set) and validation set (test set) .....	24
V) Modeling process .....	26
5.1 Linear regression model .....	26
5.2 Random forest model .....	29
5.3 K-nearest neighbors model .....	32
5.4 Ridge regression model .....	35
5.5 Gradient boosting machine .....	37
VI) Results and interpretation for the models .....	39
VII) Conclusions .....	39
VIII) Bibliography .....	40

```
## Warning: package 'pdftools' was built under R version 4.0.4  
## Warning: package 'dslabs' was built under R version 4.0.4  
## Warning: package 'corrplot' was built under R version 4.0.4
```

### Introduction:

The video game **FIFA** is one of the most played and recognized games around the world, currently this game has 30.4 million active players and much of this support for the game is due to the fact that many of its players are fans of one of the biggest sports in the world, which is soccer. In this game you can play soccer matches with any team, in any stadium, in any competition and most importantly with anyone who wants it. In addition to that, this game has the option of being able to start a career as a technical director or as a player of a team. Also, you can play online with friends or others players, with a predetermined or created team, etc.

In this work, we will focus on being technical directors and we will manage the database with all existing players and their information in the FIFA 2021 game. Much of the information that will be used is the following: **Personal, club and player statistics data**. Within this information, we can find data such as: **Player overall, height, weight, club, team position, statistics, ID, date of birth, market value, wages, among many others..** Before we focus on the purpose of the project, we will make a brief description of the data. The initial data set has **18944 obs. and 106 variables**, where we will only focus on the players who play on the field (**NON-GOALKEPEERS IN THE DATA SET**) and with this initial data set, elaborate three data sets with specific selected information, in consequence, our final data sets has **PlayerInfo: (16861 obs. of 11 variables)**, **PlayerClubInfo: (15959 obs. of 11 variables)** and **PlayerStats: (16861 obs. of 43 variables)**.

The purpose of this project is to create predictions for the player overalls in order to know and anticipate their overalls for the next FIFA game, which will be the “FIFA 2022”. These predictions will be build, based on the current statistics and personal information of the player such as: age, international reputation, weak foot, etc,. Continuing with the purpose of the project, the following processes can be found in order to achieve the proposed objective: **DOWNLOAD, UNZIP AND READ DATA PROCESS, DATA CLEANING, FILTERING AND SELECTING PROCESS, DATA ANALYSIS PROCESS and MODELING PROCESS**. Finally, as a general objective in this project, it is to demonstrate our abilities when **creating, analyzing, shaping, visualizing, presenting** and above all, being precise when **constructing and evaluating** the data for predictive models.

**IMPORTANT:** In this report, code, messages and some results were hidden.

## DATA CLEANING, FILTERING AND SELECTING PROCESS:

In this process, we created three data sets (**PlayerInfo**, **PlayerCLubInfo** and **PlayerStats**) with specific selected columns for future data analysis and modeling work. Also, for this process, we filtered the data by player positions in order to only have field players in the data sets and cleaned them omitting and replacing NA values.

**PlayerInfo data set:** This data set was create to save the information of the player.

```
## [1] 0

## # A tibble: 16,861 x 11
##   sofifa_id short_name      age dob    nationality height_cm weight_kg
##   <dbl> <chr>        <dbl> <date> <chr>          <dbl>     <dbl>
## 1 158023 L. Messi       33 1987-06-24 Argentina      170      72
## 2 20801 Cristiano Ronaldo 35 1985-02-05 Portugal      187      83
## 3 188545 R. Lewandowski 31 1988-08-21 Poland       184      80
## 4 190871 Neymar Jr      28 1992-02-05 Brazil       175      68
## 5 192985 K. De Bruyne    29 1991-06-28 Belgium      181      70
## 6 231747 K. Mbappé       21 1998-12-20 France      178      73
## 7 203376 V. van Dijk     28 1991-07-08 Netherlands   193      92
## 8 208722 S. Mané         28 1992-04-10 Senegal      175      69
## 9 209331 M. Salah        28 1992-06-15 Egypt       175      71
## 10 153079 S. Agüero      32 1988-06-02 Argentina    173      70
## # ... with 16,851 more rows, and 4 more variables: player_positions <chr>,
## #   wage_eur <dbl>, value_eur <dbl>, overall <dbl>
```

**PlayerCLubInfo data set:** This data set was create to save the information of the player and club.

```
## [1] 0

## # A tibble: 15,959 x 11
##   sofifa_id short_name   club_name league_name league_rank team_position
##   <dbl> <chr>        <chr>      <chr>          <dbl> <chr>
## 1 158023 L. Messi    FC Barcelona Spain Primera~      1 CAM
## 2 20801 Cristiano Ro~ Juventus   Italian Serie~      1 LS
## 3 188545 R. Lewandows~ FC Bayern M~ German 1. Bun~      1 ST
## 4 190871 Neymar Jr    Paris Saint~ French Ligue 1      1 LW
## 5 192985 K. De Bruyne Manchester ~ English Premi~      1 RCM
## 6 231747 K. Mbappé     Paris Saint~ French Ligue 1      1 LS
## 7 203376 V. van Dijk    Liverpool   English Premi~      1 LCB
## 8 208722 S. Mané       Liverpool   English Premi~      1 LW
## 9 209331 M. Salah      Liverpool   English Premi~      1 RW
## 10 153079 S. Agüero     Manchester ~ English Premi~      1 SUB
## # ... with 15,949 more rows, and 5 more variables: value_eur <dbl>,
## #   wage_eur <dbl>, release_clause_eur <dbl>, player_positions <chr>,
## #   overall <dbl>
```

**PlayerStats data set:** This data set was created to save the general information of the player and player stats.

```
## [1] 0

## # A tibble: 16,861 x 43
##   sofifa_id short_name    potential work_rate international_reputation overall    age
##   <dbl> <chr>          <dbl> <chr>           <dbl> <dbl> <dbl>
## 1     158023 L. Messi      93 Medium/Low       5     93    33
## 2     20801 Cristiano Ronaldo 92 High/Low        5     92    35
## 3     188545 R. Lewandowski 91 High/Medium      4     91    31
## 4     190871 Neymar Jr      91 High/Medium      5     91    28
## 5     192985 K. De Bruyne   91 High/High        4     91    29
## 6     231747 K. Mbappé      95 High/Low        3     90    21
## 7     203376 V. van Dijk    91 Medium/Medium    3     90    28
## 8     208722 S. Mané        90 High/Medium      3     90    28
## 9     209331 M. Salah       90 High/Medium      3     90    28
## 10    153079 S. Agüero      89 High/Medium      4     89    32
## # ... with 16,851 more rows, and 36 more variables: weak_foot <dbl>,
## # skill_moves <dbl>, pace <dbl>, shooting <dbl>, passing <dbl>,
## # dribbling <dbl>, defending <dbl>, physic <dbl>, attacking_crossing <dbl>,
## # attacking_finishing <dbl>, attacking_heading_accuracy <dbl>,
## # attacking_short_passing <dbl>, attacking_volleys <dbl>,
## # movement_acceleration <dbl>, movement_sprint_speed <dbl>,
## # movement_agility <dbl>, movement_reactions <dbl>, movement_balance <dbl>,
## # power_shot_power <dbl>, power_jumping <dbl>, power_stamina <dbl>,
## # power_strength <dbl>, power_long_shots <dbl>, mentality_aggression <dbl>,
## # mentality_interceptions <dbl>, mentality_positioning <dbl>,
## # mentality_composure <dbl>, mentality_vision <dbl>,
## # mentality_penalties <dbl>, defending_standing_tackle <dbl>,
## # defending_sliding_tackle <dbl>, skill_curve <dbl>, skill_dribbling <dbl>,
## # skill_fk_accuracy <dbl>, skill_long_passing <dbl>, skill_ball_control <dbl>
```

## DATA ANALYSIS PROCESS:

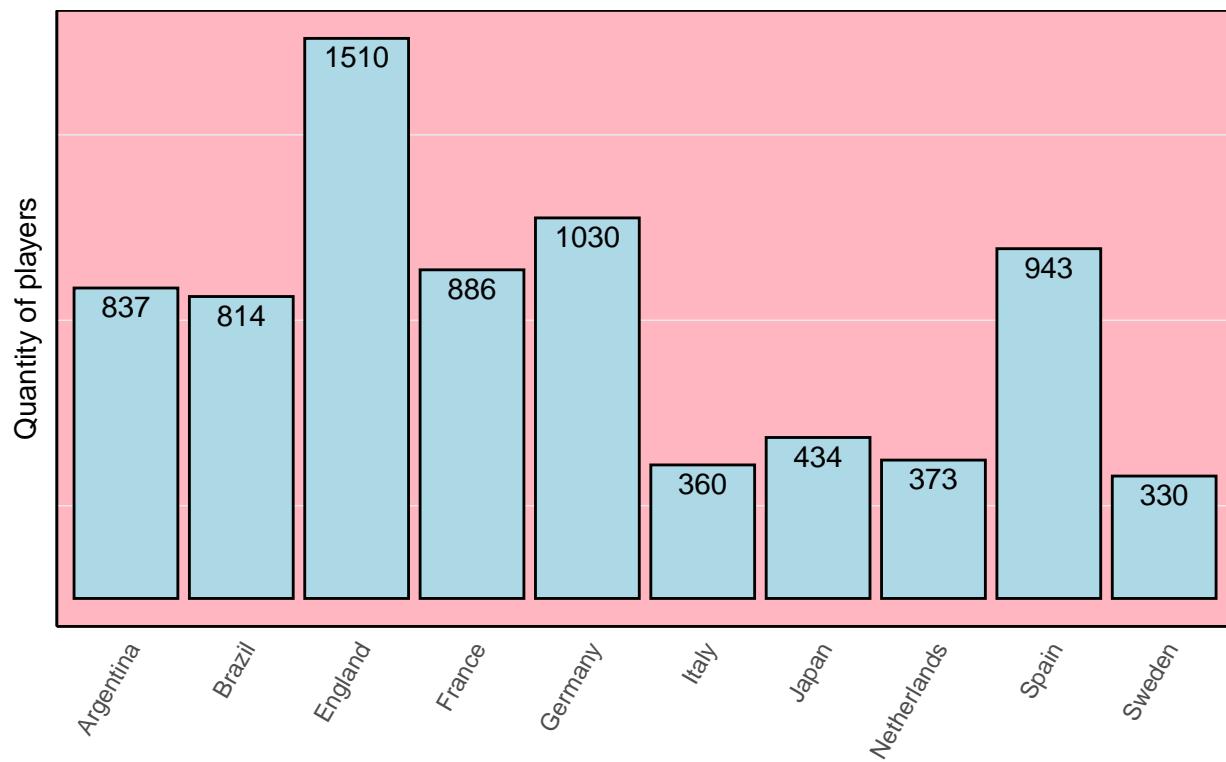
In this process, we applied our management and visualization skills in the previously cleaned data sets that were created in the **DATA CLEANING, FILTERING and SELECTING PROCESS**. In the same, we elaborated bar plots and tables in order to achieve a deeper comprehension of the data. (This process can be considered crucial for the modeling process).

### Top ten nationalities based on quantity:

As you can notice, in the following table and plot for the top ten nationalities based on quantity. The nationality with more players in the data set is **England** with **1510** professional soccer players.

nationality	n
England	1510
Germany	1026
Spain	943
France	886
Argentina	837
Brazil	814
Japan	434
Netherlands	373
Italy	360
Sweden	330

Top 10 nationalities based on quantity

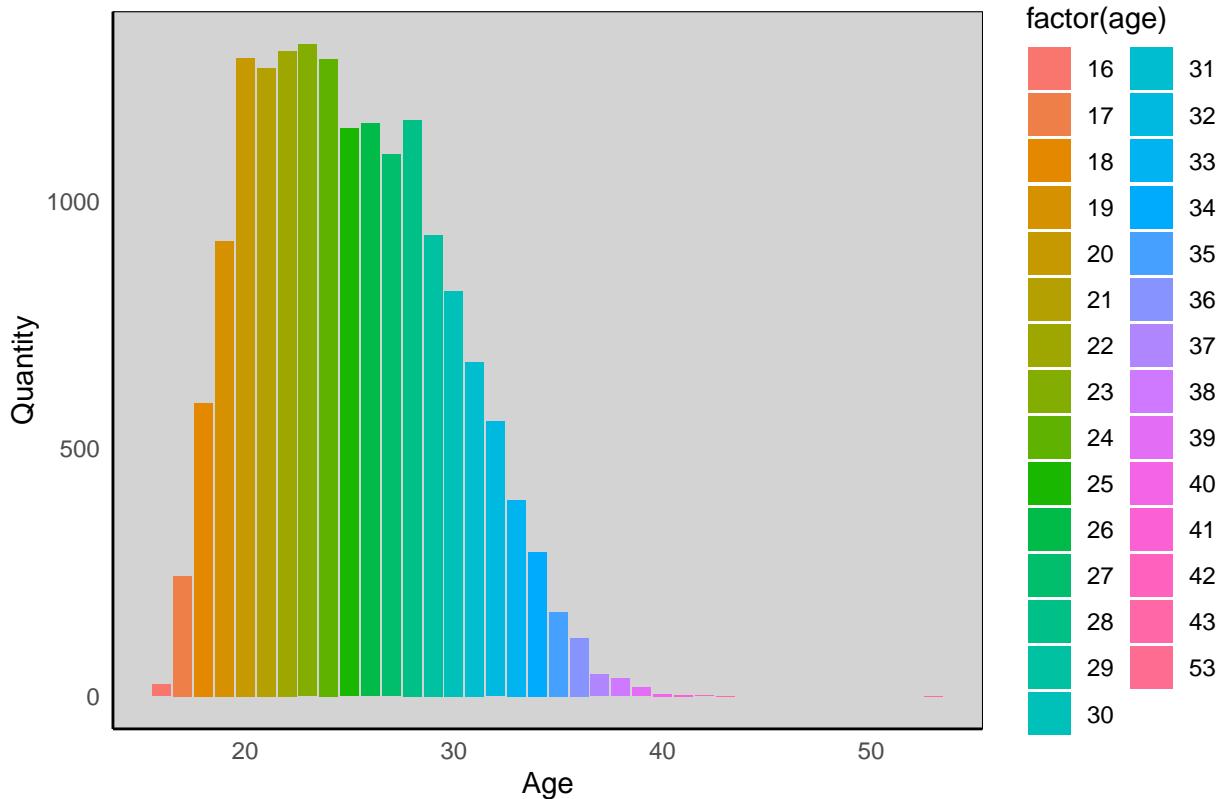


### Age distribution:

In the following table and plot below, we can see that the more recurrent players age in the data set is **23 years** with **1316** appearances and the less recurrent is share between **43 years** and **53 years** with **1** appearances.

age	n
16	24
17	243
18	592
19	919
20	1288
21	1268
22	1302
23	1316
24	1286
25	1146
26	1156
27	1094
28	1163
29	931
30	818
31	675
32	555
33	396
34	291
35	170
36	118
37	44
38	36
39	18
40	5
41	3
42	2
43	1
53	1

## Age distribution



## Top 10 players in terms of potential

In the table below, we can see the top ten players in terms of potential. Also, we can appreciate that the first position in this table is for **K. Mbappé** with a potential overall of **95**.

potential	Name
95	K. Mbappé
93	L. Messi
93	J. Sancho
93	K. Havertz
93	João Félix
93	Vinícius Jr.
92	Cristiano Ronaldo
92	T. Alexander-Arnold
92	M. de Ligt
92	E. Haaland

### Top ten older and youngest players:

If we take a look into the tables below, we can find the top ten older and youngest players in the data set. In the top ten older players table, the first position is for **K. Miura** with an age of **53 years**, meanwhile, for the top ten youngest players table, there isn't any position because of all them count with an age of **16 years**.

age	Name
53	K. Miura
43	H. Sulaimani
42	Hilton
42	S. Nakamura
41	Lee Dong Gook
41	D. Bulman
41	K. Ellison
40	J. Cáceres
40	Nino
40	S. Aquino

age	Name
16	R. Cherki
16	A. Karabec
16	U. Bertelli
16	D. Hoyo-Kowalski
16	L. Gourna-Douath
16	O. Babuscu
16	M. Tanlongo
16	O. Beyaz
16	W. Faghir
16	A. Descotte

### Top ten leagues based on players quantity:

As we can see in table below, most of the players play in the MLS (USA Major League Soccer) league. This league counts with **604 players**.

league_name	n
USA Major League Soccer	604
English League Championship	603
English Premier League	577
Argentina Primera División	560
Italian Serie A	545
Spain Primera Division	541
French Ligue 1	525
English League One	521
English League Two	514
Turkish Süper Lig	504

### Top ten players with best wage:

In the following table and plot, we can appreciate the top ten players with best wage in the data set. The player with the best wage is the FC Barcelona player **L. Messi** (best known as **Lionel Messi**) with **€560,000** in a month of work.

wage_eur	Name
560000	L. Messi
370000	K. De Bruyne
350000	K. Benzema
350000	E. Hazard
310000	Casemiro
310000	T. Kroos
300000	S. Agüero
300000	Sergio Ramos
290000	A. Griezmann
280000	L. Suárez

### Top 10 players with best wage

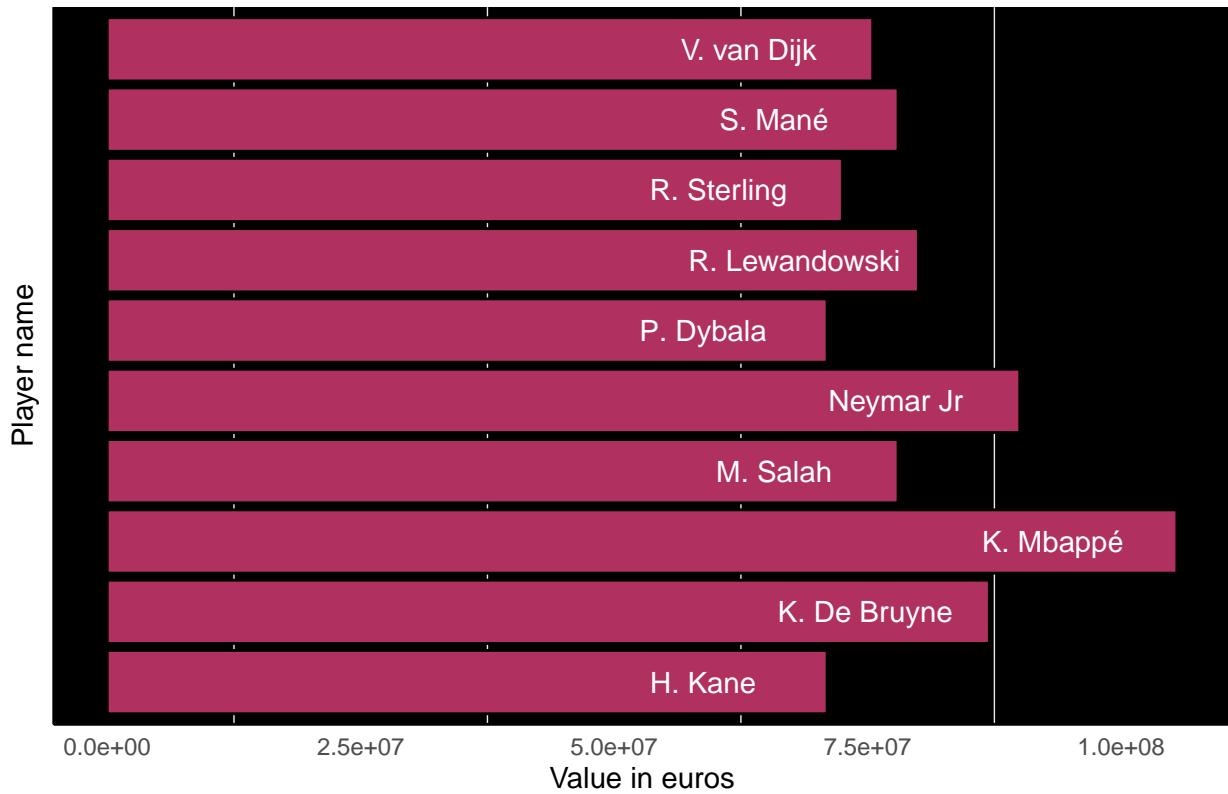


### Top ten more valuable players:

In the following table and plot, we can appreciate the top ten more valuable players in the data set. The player with the best wage is the Paris Saint German player **K. Mbappé** with a high value of **€105,500,000**.

value_eur	Name
105500000	K. Mbappé
90000000	Neymar Jr
87000000	K. De Bruyne
80000000	R. Lewandowski
78000000	S. Mané
78000000	M. Salah
75500000	V. van Dijk
72500000	R. Sterling
71000000	H. Kane
71000000	P. Dybala

Top 10 more valuable players



### Top ten highest and heaviest players:

If we look at the tables below, we can notice that the highest player in the data set is **A. Ba** with **203 cm** and the heaviest player in the data set is **A. Akinfenwa** with **110 kg**.

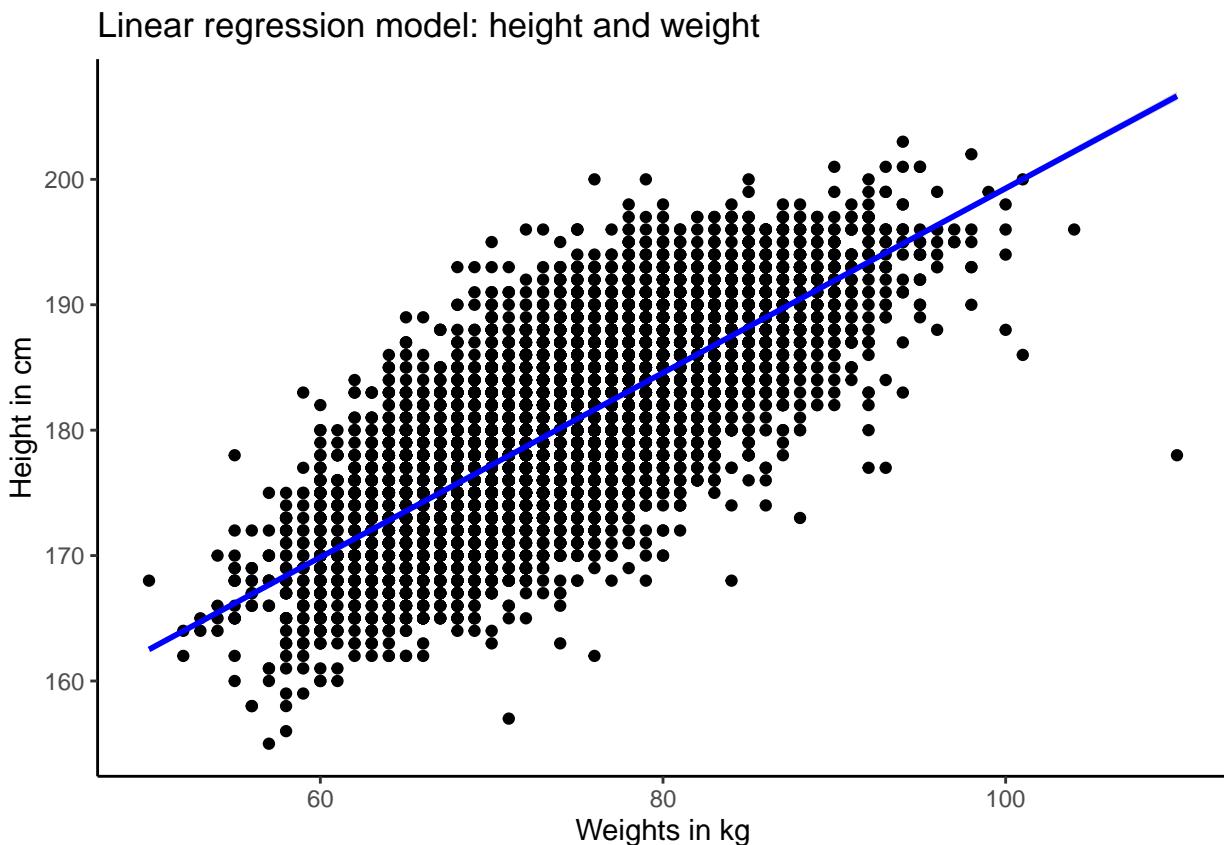
height_cm	Name
203	A. Ba
202	S. Maierhofer
201	P. Onuachu
201	H. Veerman
201	A. Vukotic
201	A. Sjöberg
201	S. Makienok
200	T. Chorý
200	K. Sidibé
200	S. Kalajdžić

weight_kg	Name
110	A. Akinfenwa
104	O. Oularé
101	W. Morgan
101	T. Chorý
100	H. Maguire
100	D. Dike
100	R. Greenidge
100	Ricardo Santos
99	T. Petrášek
98	J. Beauguel

### Linear regression model for players height and weight:

In the result and plot below, we can appreciate a high correlation between the players height and weight of **0.7354**, which means that this two variables are related and depend on each other.

```
##  
## Call:  
## lm(formula = height_cm ~ weight_kg, data = PlayerInfo)  
##  
## Coefficients:  
## (Intercept)    weight_kg  
##     125.7446      0.7354
```

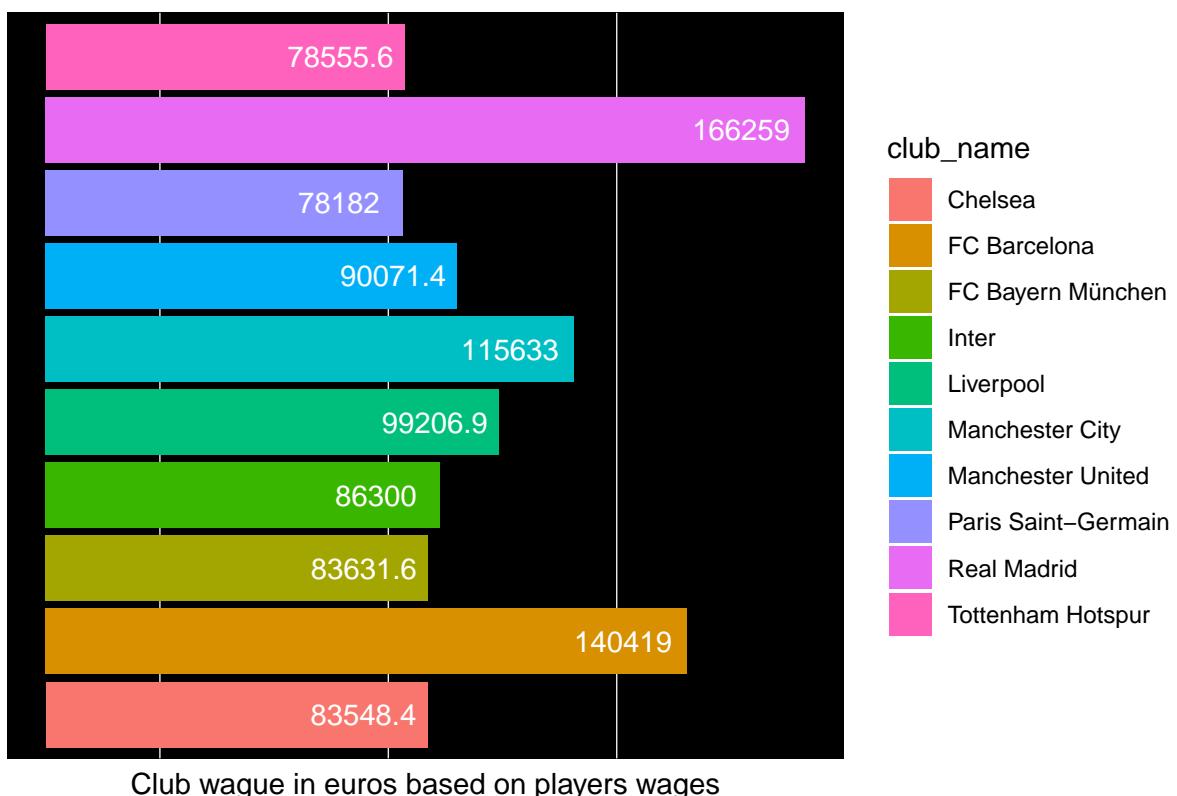


### Top ten clubs with the highest salaries:

As we can see in the top ten clubs with the highest salaries table and plot. The professional soccer club **Real Madrid** have the best paying in comparison to the other clubs in the data set with a mean of **€166,259.26** in month.

club_name	Mean
Chelsea	83548.39
FC Barcelona	140419.35
FC Bayern München	83631.58
Inter	86300.00
Liverpool	99206.90
Manchester City	115633.33
Manchester United	90071.43
Paris Saint-Germain	78182.00
Real Madrid	166259.26
Tottenham Hotspur	78555.56

### Top 10 clubs with the highest salaries

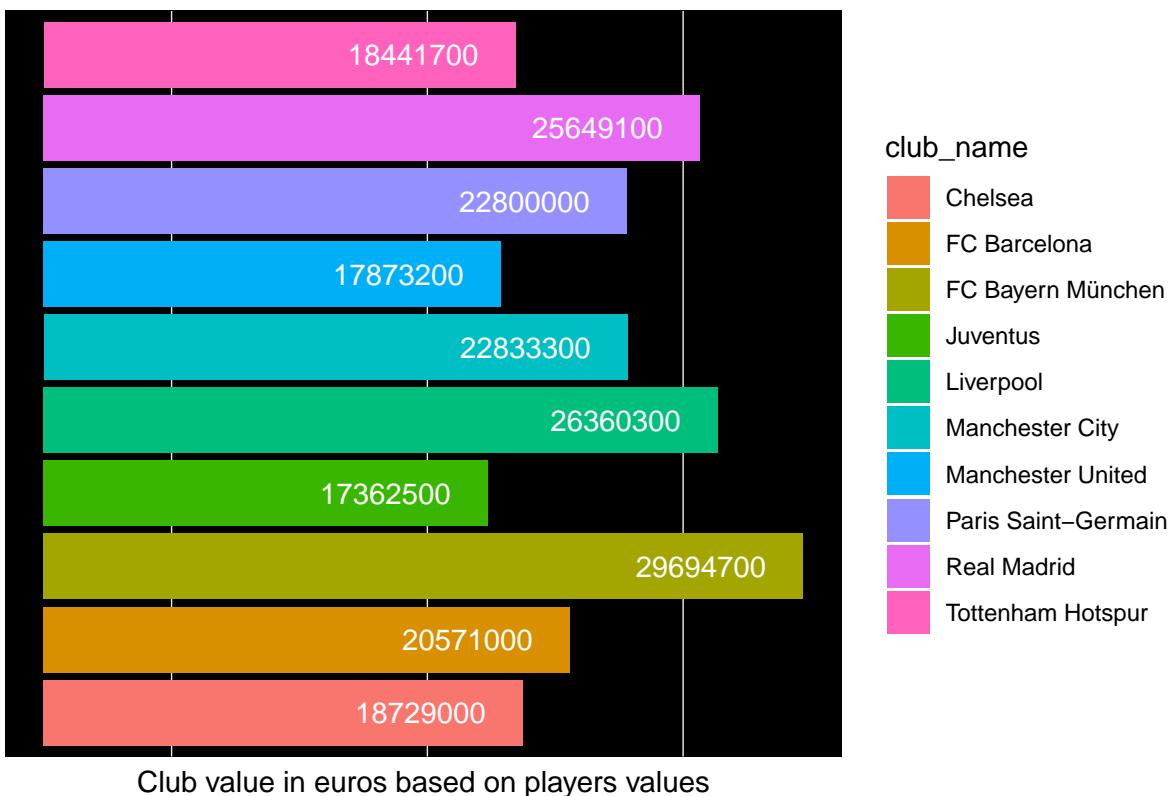


### Top ten more valuable clubs:

As we can see in the top ten more valuable clubs table and plot. The professional soccer club that can be considered as the “more valuable club” is the **FC Bayern München** with a mean value of **€29,694,737**.

club_name	Mean
Chelsea	18729032
FC Barcelona	20570968
FC Bayern München	29694737
Juventus	17362500
Liverpool	26360345
Manchester City	22833333
Manchester United	17873214
Paris Saint-Germain	22800000
Real Madrid	25649074
Tottenham Hotspur	18441667

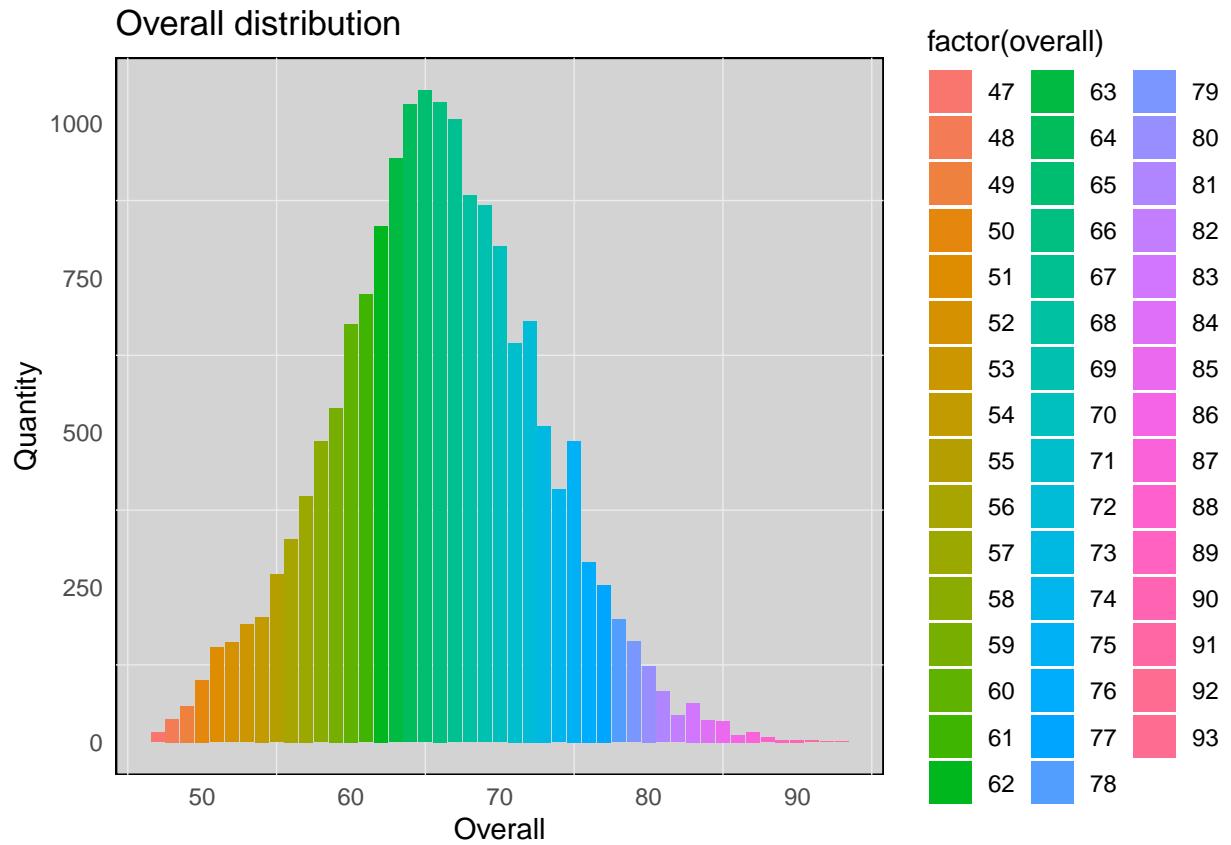
Top 10 more valuable clubs



### Overall distribution:

In the following table and plot below, we can see that the more recurrent players overall in the data set is **65** with **1053** appearances and the less recurrent is **92** and **93** with **1** appearances.

overall	n
47	16
48	38
49	59
50	101
51	153
52	161
53	190
54	202
55	271
56	328
57	398
58	486
59	540
60	676
61	723
62	834
63	943
64	1030
65	1053
66	1034
67	1006
68	883
69	867
70	801
71	645
72	680
73	511
74	408
75	487
76	291
77	254
78	198
79	163
80	123
81	82
82	43
83	63
84	36
85	34
86	12
87	17
88	8
89	4
90	4
91	3
92	1
93	1

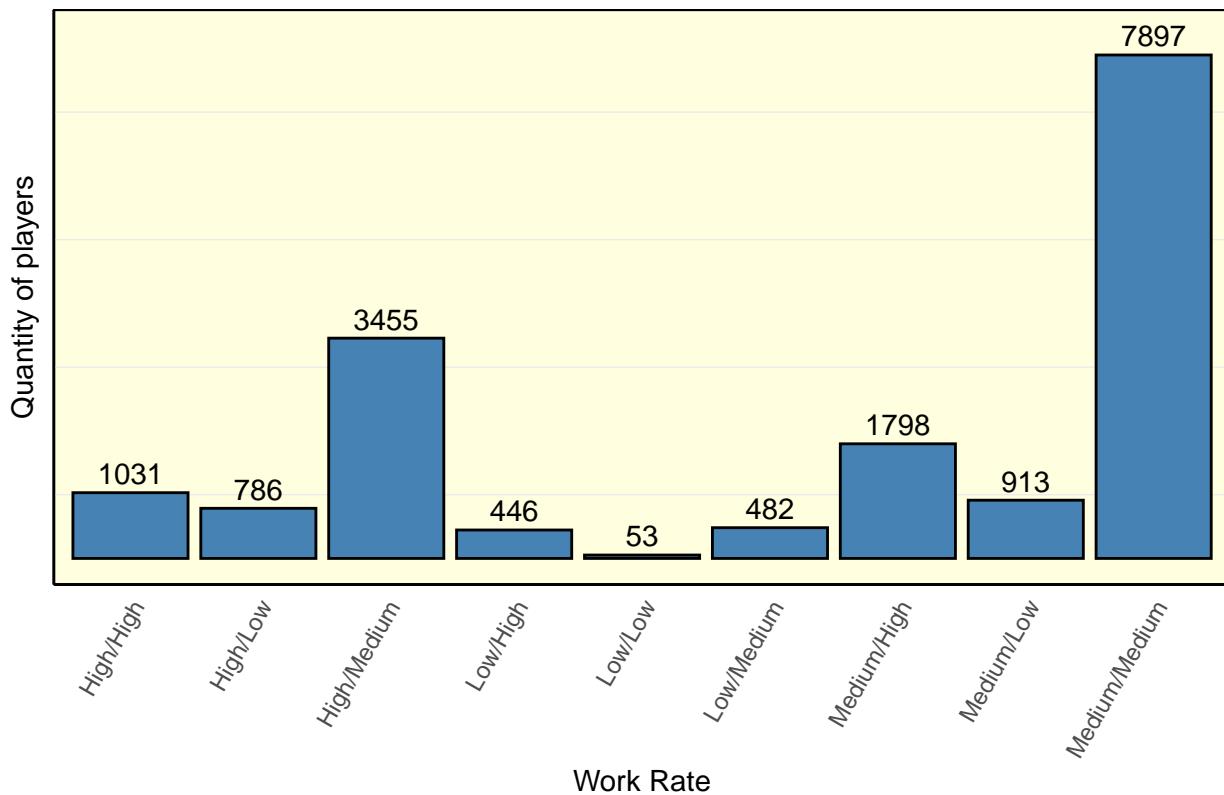


#### Distribution of players work rate:

As we can notice in the table and plot below, most of the soccer players in the data set have a work rate of **Medium/Medium** with **7897** appearances. This work rate can be related to the most recurrent overall which is **65** and with this we can conclude that the more the player work, the more overall the player will have.

work_rate	n
High/High	1031
High/Low	786
High/Medium	3455
Low/High	446
Low/Low	53
Low/Medium	482
Medium/High	1798
Medium/Low	913
Medium/Medium	7897

## Distribution of players work rate



### Top ten players in terms of overall:

In the table below, we can notice that the player with the highest overall is **L. Messi** with an overall of **93**. This explain us the reason why this player receive the highest wage in the data set and have a extremely high market value.

overall	Name
93	L. Messi
92	Cristiano Ronaldo
91	R. Lewandowski
91	Neymar Jr
91	K. De Bruyne
90	K. Mbappé
90	V. van Dijk
90	S. Mané
90	M. Salah
89	S. Agüero

### Quantity of players in every club:

Here we can see an example of this data set with ten clubs. For every clubs should be different the quantity of players (this depends on every league rules).

club_name	n
1. FC Heidenheim 1846	25
1. FC Kaiserslautern	25
1. FC Köln	26
1. FC Magdeburg	23
1. FC Nürnberg	24
1. FC Saarbrücken	24
1. FC Union Berlin	23
1. FSV Mainz 05	29
Aalborg BK	21
Aalesunds FK	23

### Distribution of players positions:

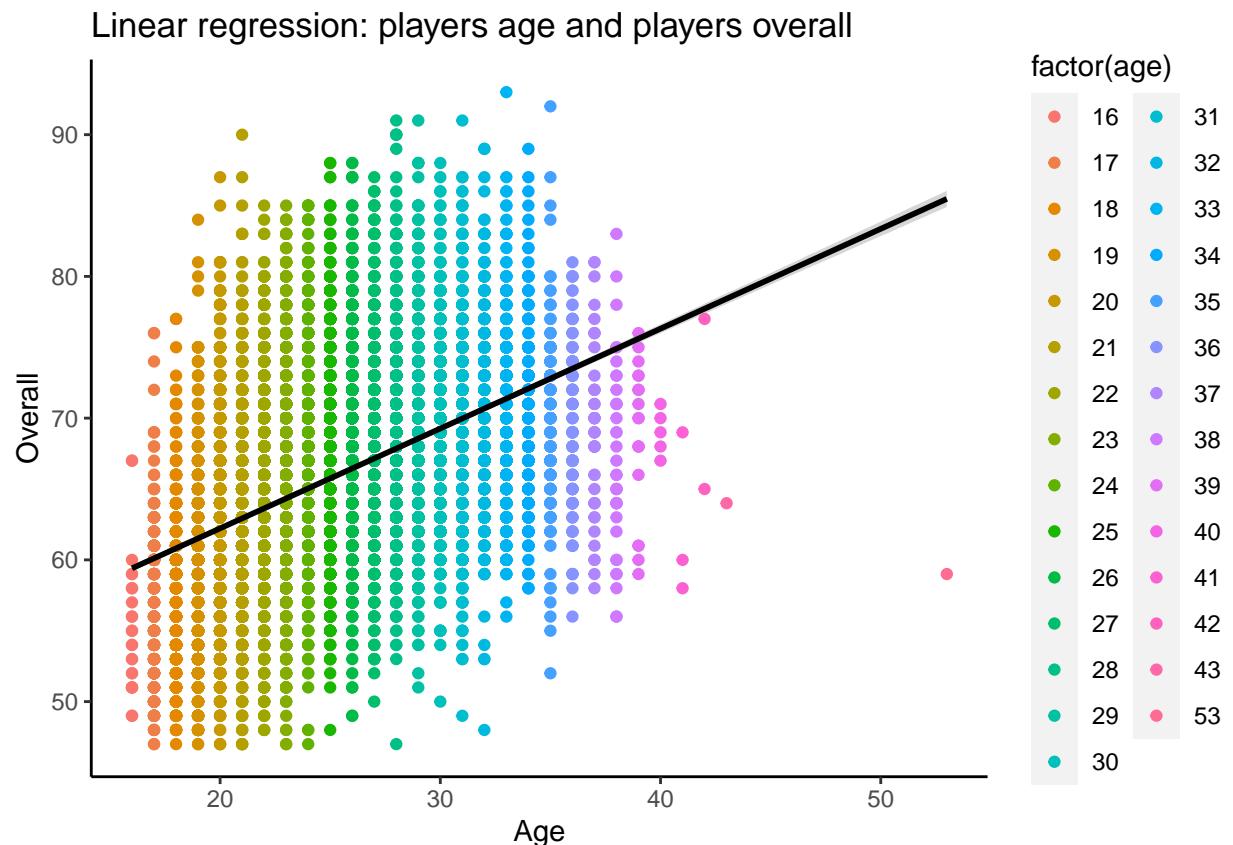
In the table below, we can appreciate the distribution of all the players positions in the data set.

team_position	n
CAM	275
CB	122
CDM	177
CF	13
CM	72
LAM	24
LB	522
LCB	655
LCM	413
LDM	241
LF	12
LM	390
LS	217
LW	160
LWB	75
RAM	25
RB	533
RCB	661
RCM	406
RDM	241
RES	2641
RF	11
RM	407
RS	211
RW	159
RWB	74
ST	423
SUB	6799

### Linear regression model for players age and players overall:

In the result and plot below, we can appreciate a high correlation between the players age and players overall of **0.70**, which means that this two variables are related and depend on each other. After interpreting the result and plot below, we can notice that most of the younger and older players have a lower overall in comparison to the players that can't be considered as old or young players.

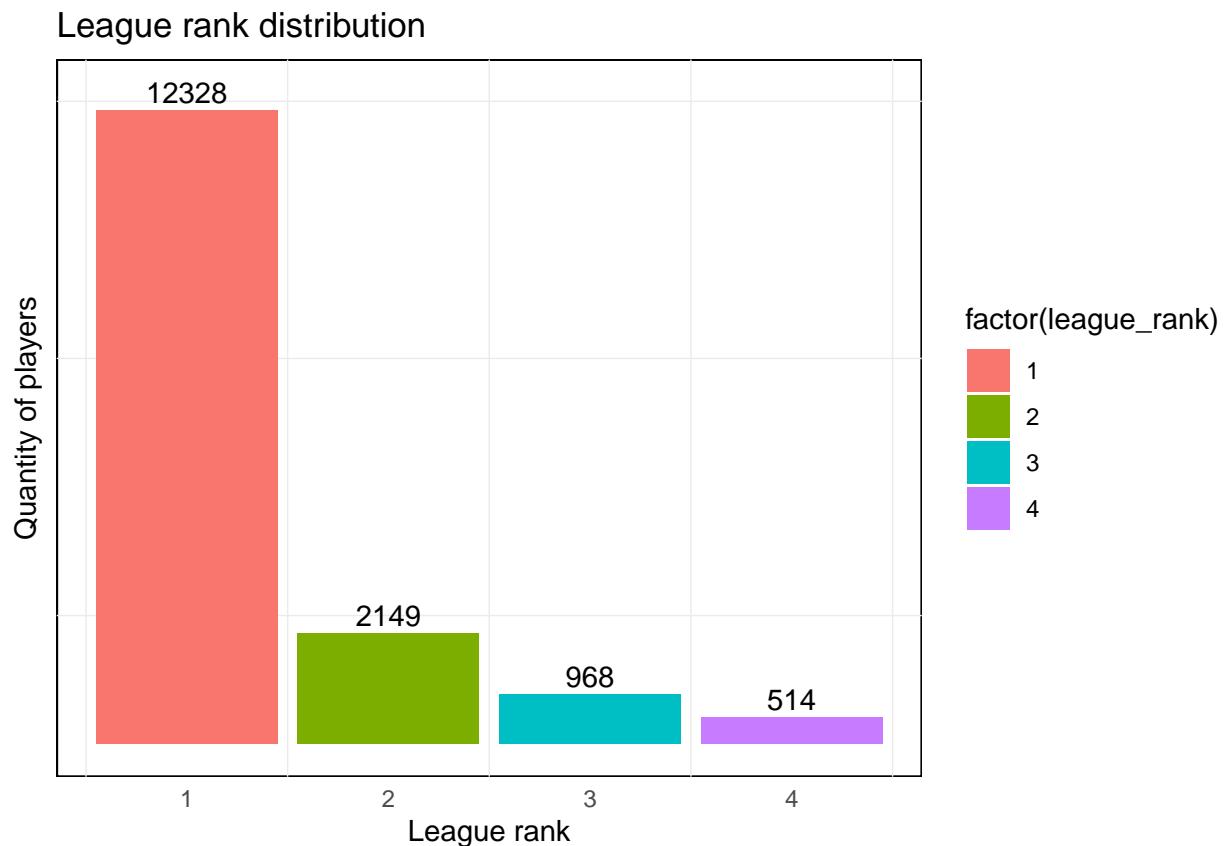
```
##  
## Call:  
## lm(formula = overall ~ age, data = FIFA2021)  
##  
## Coefficients:  
## (Intercept)      age  
##     48.0718      0.6979  
  
##  
## Call:  
## lm(formula = overall ~ age, data = FIFA2021)  
##  
## Coefficients:  
## (Intercept)      age  
##     48.0718      0.6979
```



#### League rank distribution:

In the following plot and table below, we can appreciate that most of the players (**12328 out of 16861**) in the data set play in the first league rank. This league rank can be considered as the better league rank in the data set because it counts with the best salaries and more valuable clubs.

league_rank	n
1	12328
2	2149
3	968
4	514

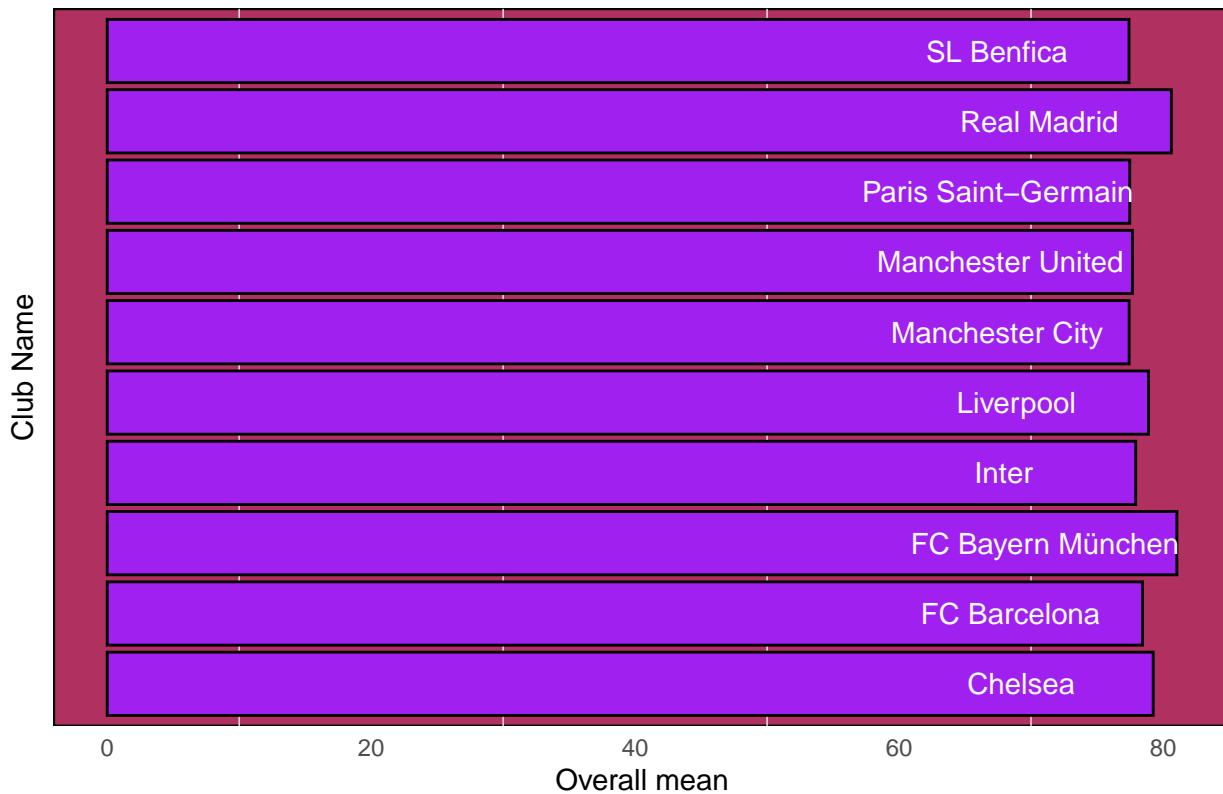


### Top ten clubs based on overall:

In the following plot and table, we can appreciate the top ten clubs in overall terms. In the first position, we can see the professional soccer club **Real Madrid** with a mean overall of **80.62963**.

club_name	Mean
Chelsea	79.25806
FC Barcelona	78.45161
FC Bayern München	81.05263
Inter	77.93333
Liverpool	78.89655
Manchester City	77.43333
Manchester United	77.67857
Paris Saint-Germain	77.48000
Real Madrid	80.62963
SL Benfica	77.42308

### Top 10 clubs based on overall

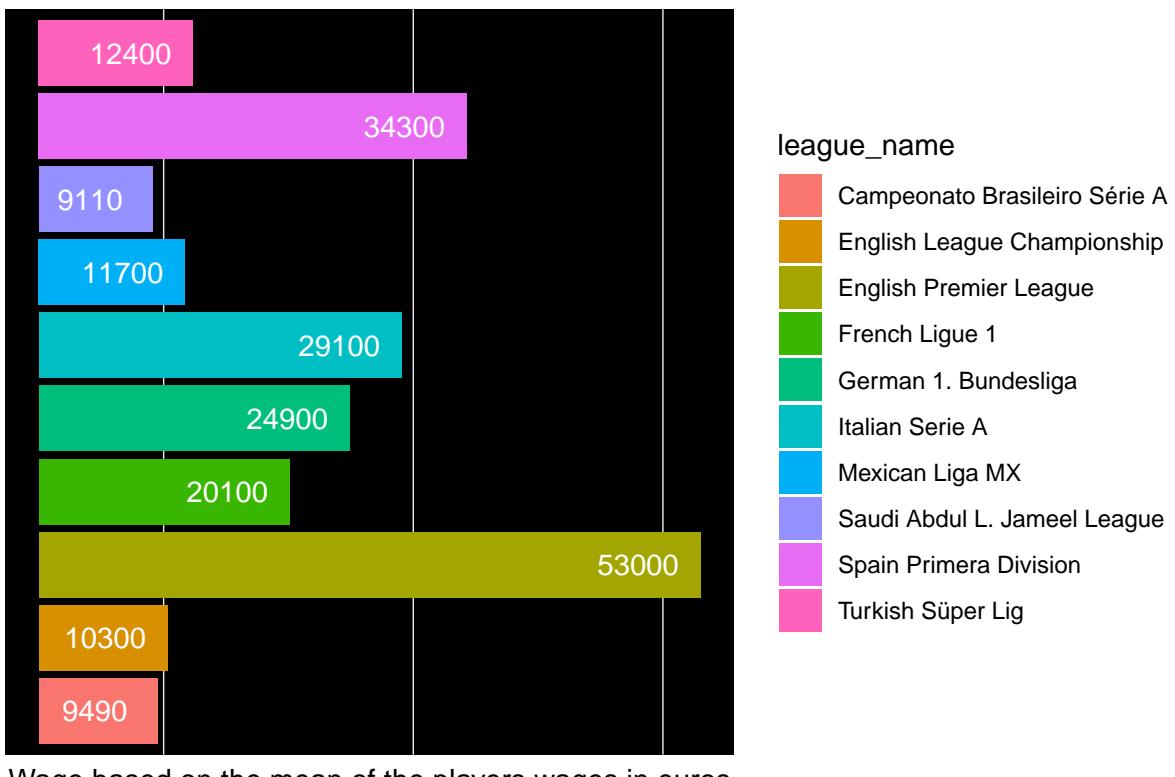


### Top ten leagues with the highest salaries:

If we take a look into the top ten leagues with the highest salaries table and plot, we can notice that the league with the best paying is the **English Premier League** with a monthly pay mean of **€53018.198**.

league_name	Mean
Campeonato Brasileiro Série A	9486.154
English League Championship	10322.554
English Premier League	53018.198
French Ligue 1	20101.143
German 1. Bundesliga	24900.640
Italian Serie A	29074.862
Mexican Liga MX	11716.870
Saudi Abdul L. Jameel League	9105.820
Spain Primera Division	34293.253
Turkish Süper Lig	12351.587

### Top 10 leagues with the highest salaries

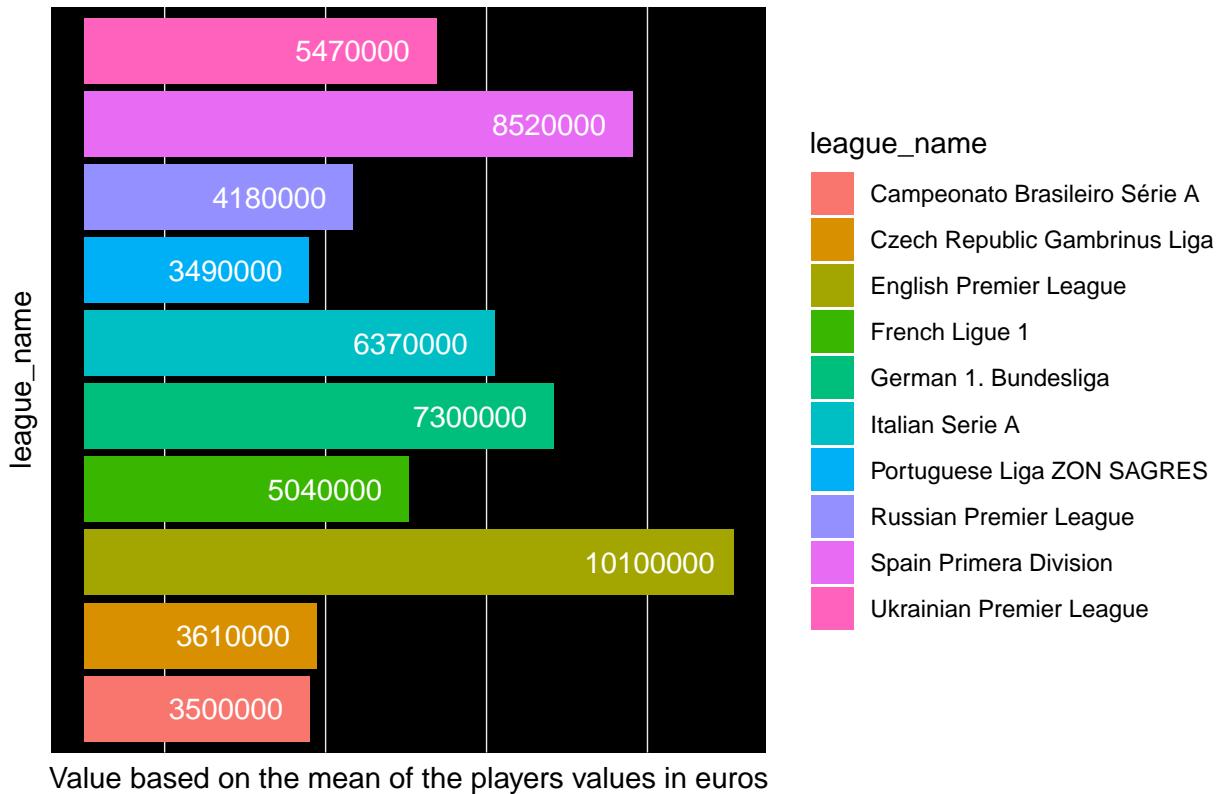


### Top ten more valuable leagues:

In the top ten more valuable leagues table and plot, we can appreciate that the more valuable league is the **English Premier League** with a mean value of **€10,084,749**.

league_name	Mean
Campeonato Brasileiro Série A	3502292
Czech Republic Gambrinus Liga	3612857
English Premier League	10084749
French Ligue 1	5041219
German 1. Bundesliga	7295981
Italian Serie A	6370266
Portuguese Liga ZON SAGRES	3494452
Russian Premier League	4177365
Spain Primera Division	8522357
Ukrainian Premier League	5468627

### Top 10 more valuable leagues



## CREATE FIFA SET (TRAINING SET) AND VALIDATION SET (TEST SET):

In this process, we created a training and validation set in order to start making our modeling process for the FIFA 2022 overalls predictions.

```

## Warning in set.seed(123, sample.kind = "Rounding"): non-uniform 'Rounding'
## sampler used

## [1] 0

## # A tibble: 10,115 x 43
##   sofifa_id short_name    potential work_rate international_repu~ overall    age
##       <dbl> <chr>          <dbl> <chr>                  <dbl>    <dbl> <dbl>
## 1     158023 L. Messi        93 Medium/Low            5      93    33
## 2     20801 Cristiano Ro~    92 High/Low             5      92    35
## 3     188545 R. Lewandows~   91 High/Medi~           4      91    31
## 4     190871 Neymar Jr      91 High/Medi~           5      91    28
## 5     203376 V. van Dijk    91 Medium/Me~          3      90    28
## 6     209331 M. Salah       90 High/Medi~           3      90    28
## 7     153079 S. Agüero      89 High/Medi~           4      89    32
## 8     155862 Sergio Ramos   89 High/Medi~           4      89    34
## 9     165153 K. Benzema     89 Medium/Low           4      89    32
## 10    200145 Casemiro       89 High/High            3      89    28
## # ... with 10,105 more rows, and 36 more variables: weak_foot <dbl>,
## #   skill_moves <dbl>, pace <dbl>, shooting <dbl>, passing <dbl>,
## #   dribbling <dbl>, defending <dbl>, physic <dbl>, attacking_crossing <dbl>,
## #   attacking_finishing <dbl>, attacking_heading_accuracy <dbl>,
## #   attacking_short_passing <dbl>, attacking_volleys <dbl>,
## #   movement_acceleration <dbl>, movement_sprint_speed <dbl>,
## #   movement_agility <dbl>, movement_reactions <dbl>, movement_balance <dbl>,
## #   power_shot_power <dbl>, power_jumping <dbl>, power_stamina <dbl>,
## #   power_strength <dbl>, power_long_shots <dbl>, mentality_aggression <dbl>,
## #   mentality_interceptions <dbl>, mentality_positioning <dbl>,
## #   mentality_composure <dbl>, mentality_vision <dbl>,
## #   mentality_penalties <dbl>, defending_standing_tackle <dbl>,
## #   defending_sliding_tackle <dbl>, skill_curve <dbl>, skill_dribbling <dbl>,
## #   skill_fk_accuracy <dbl>, skill_long_passing <dbl>, skill_ball_control <dbl>

## [1] 0

## # A tibble: 6,746 x 43
##   sofifa_id short_name    potential work_rate international_repu~ overall    age
##       <dbl> <chr>          <dbl> <chr>                  <dbl>    <dbl> <dbl>
## 1     192985 K. De Bruy~    91 High/High            4      91    29
## 2     231747 K. Mbappé     95 High/Low             3      90    21
## 3     208722 S. Mané       90 High/Medium          3      90    28
## 4     211110 P. Dybala     89 Medium/Med~          3      88    26
## 5     182521 T. Kroos       88 Medium/Med~          4      88    30
## 6     201024 K. Kouliba~   88 Medium/High           3      88    29
## 7     216267 A. Roberts~   89 High/High            3      87    26
## 8     192387 C. Immobile   87 High/Medium          3      87    30
## 9     194765 A. Griezma~   87 Medium/Med~          4      87    29
## 10    200104 H. Son        87 High/High            3      87    27

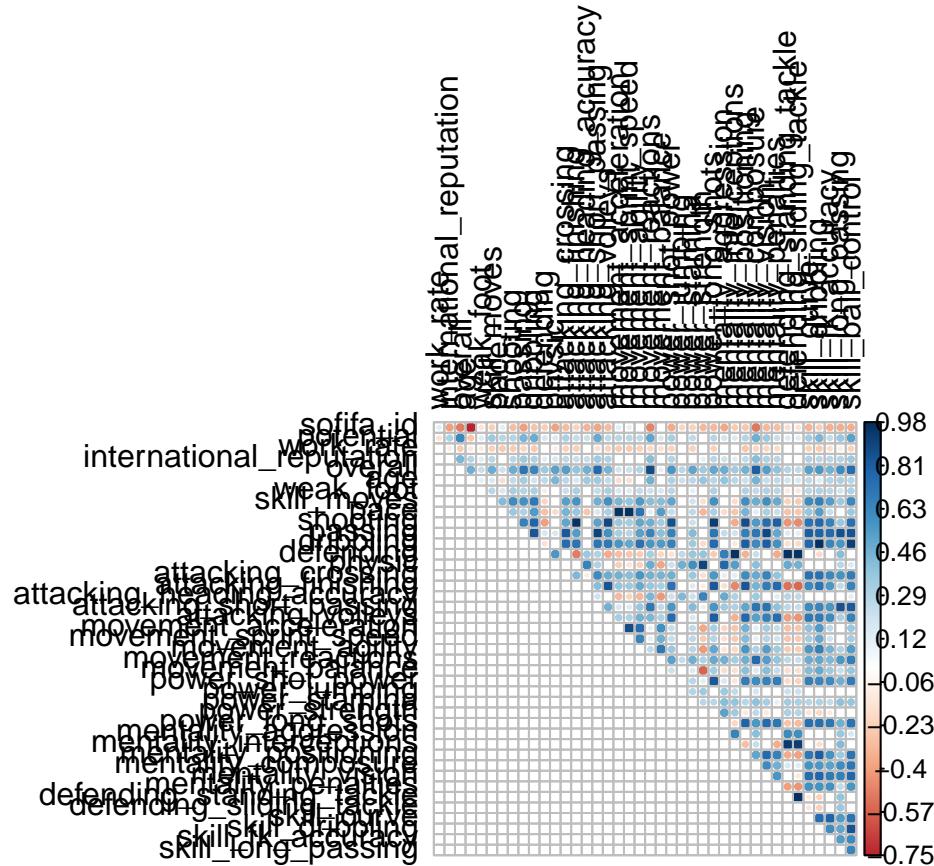
```

```

## # ... with 6,736 more rows, and 36 more variables: weak_foot <dbl>,
## #   skill_moves <dbl>, pace <dbl>, shooting <dbl>, passing <dbl>,
## #   dribbling <dbl>, defending <dbl>, physic <dbl>, attacking_crossing <dbl>,
## #   attacking_finishing <dbl>, attacking_heading_accuracy <dbl>,
## #   attacking_short_passing <dbl>, attacking_volleys <dbl>,
## #   movement_acceleration <dbl>, movement_sprint_speed <dbl>,
## #   movement_agility <dbl>, movement_reactions <dbl>, movement_balance <dbl>,
## #   power_shot_power <dbl>, power_jumping <dbl>, power_stamina <dbl>,
## #   power_strength <dbl>, power_long_shots <dbl>, mentality_aggression <dbl>,
## #   mentality_interceptions <dbl>, mentality_positioning <dbl>,
## #   mentality_composure <dbl>, mentality_vision <dbl>,
## #   mentality_penalties <dbl>, defending_standing_tackle <dbl>,
## #   defending_sliding_tackle <dbl>, skill_curve <dbl>, skill_dribbling <dbl>,
## #   skill_fk_accuracy <dbl>, skill_long_passing <dbl>, skill_ball_control <dbl>

```

Before starting the **modeling process**, we have to analyze our variables. The most common and accurate form to evaluate our variables was creating a function that make correlations for all of them. After having the results in a plot, we had to choose the variables with the best results in order to build more precise models.



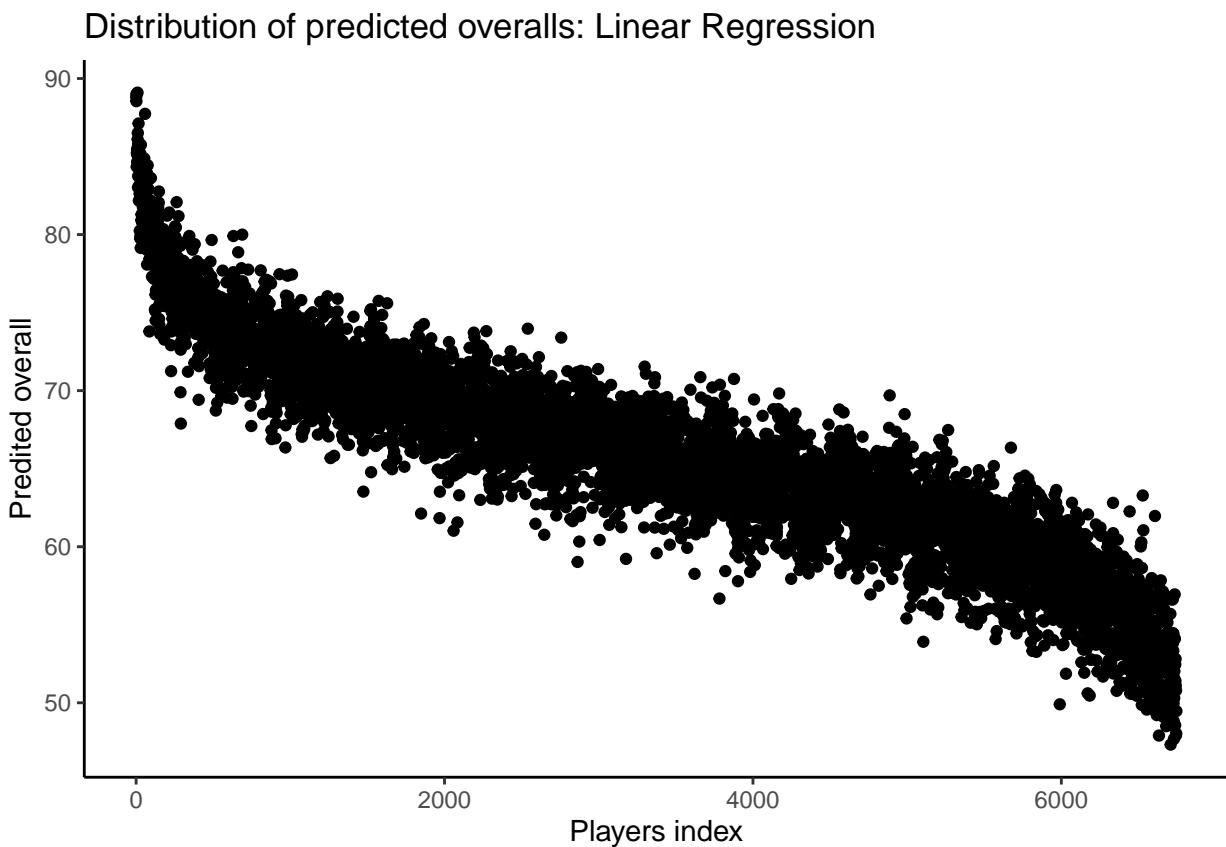
## MODELING PROCESS:

In this process, we trained and built different types of models in order to create predictions for the FIFA 2022 players overall. The purpose of this process was to achieve the most accurate predictions using the RMSE as our evaluation metric and demonstrate data visualization, analytics, training and precision skills. All models were retrieved from the caret package.

For the **Linear Regression** model, we did predictions based on the method “lm” and created plots in order to visualize the results. This model can be considered one of the most common for predictive models.

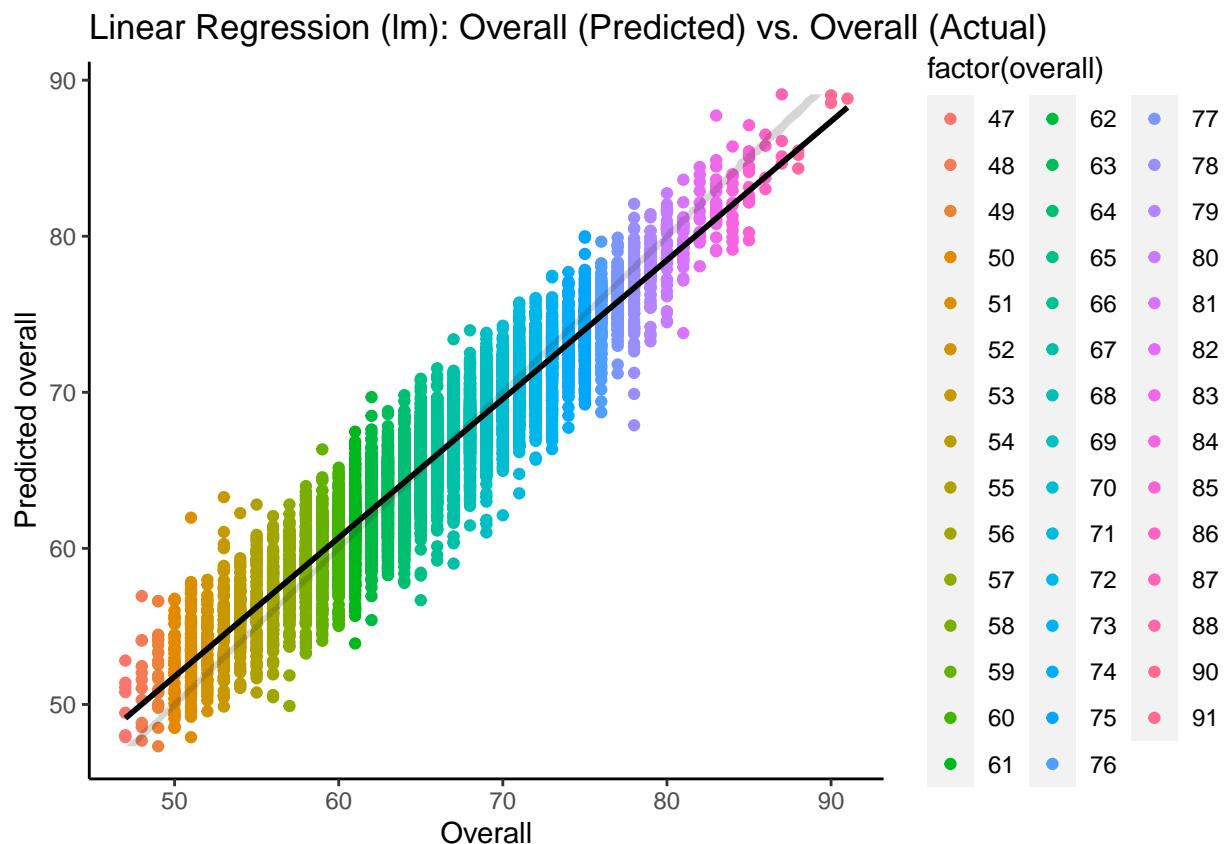
**Formula used:**  $Y = a + bX + e$

After we made the predictions for the **Linear Regression** model, we elaborated a plot in order to visualize the distribution of predicted overalls for all the players in the training set.



In the following plot, we can see a linear regression graph for our **Linear Regression** model. This plot shows us the relationship between the predicted overall and the actual overall. In this plot we can see a high correlation/relationship with a value of **0.8894** in this two variables. However, this graph and value can be improved depending on the prediction made for each model.

```
##  
## Call:  
## lm(formula = fit ~ overall, data = Outcome_linear_regression)  
##  
## Coefficients:  
## (Intercept)      overall  
##       7.3219        0.8894
```



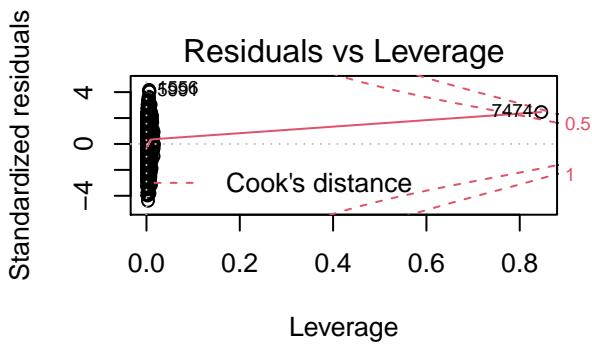
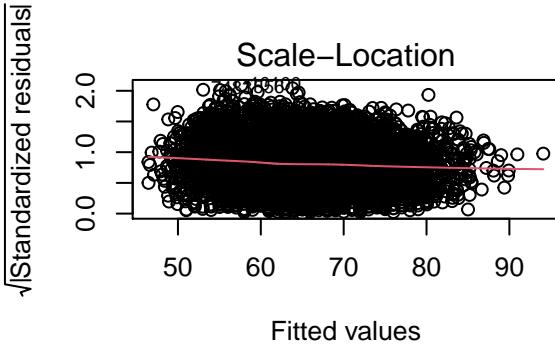
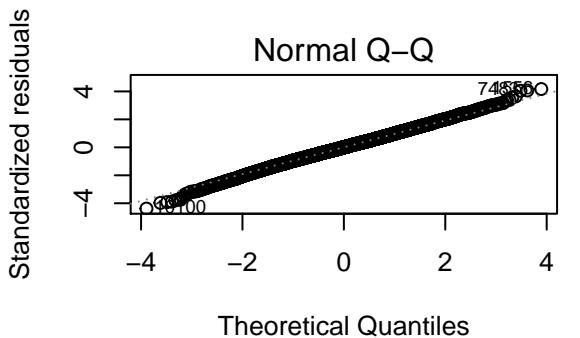
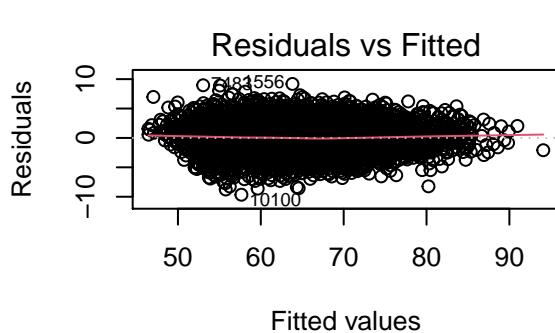
RESULT (RMSE) FOR THE LINEAR REGRESSION MODEL:

$$\overline{\overline{x}} = 2.240636$$

## EXTRA GRAPHS:

In the following plot, we can see the main four graphs for the **Linear Regression** model.

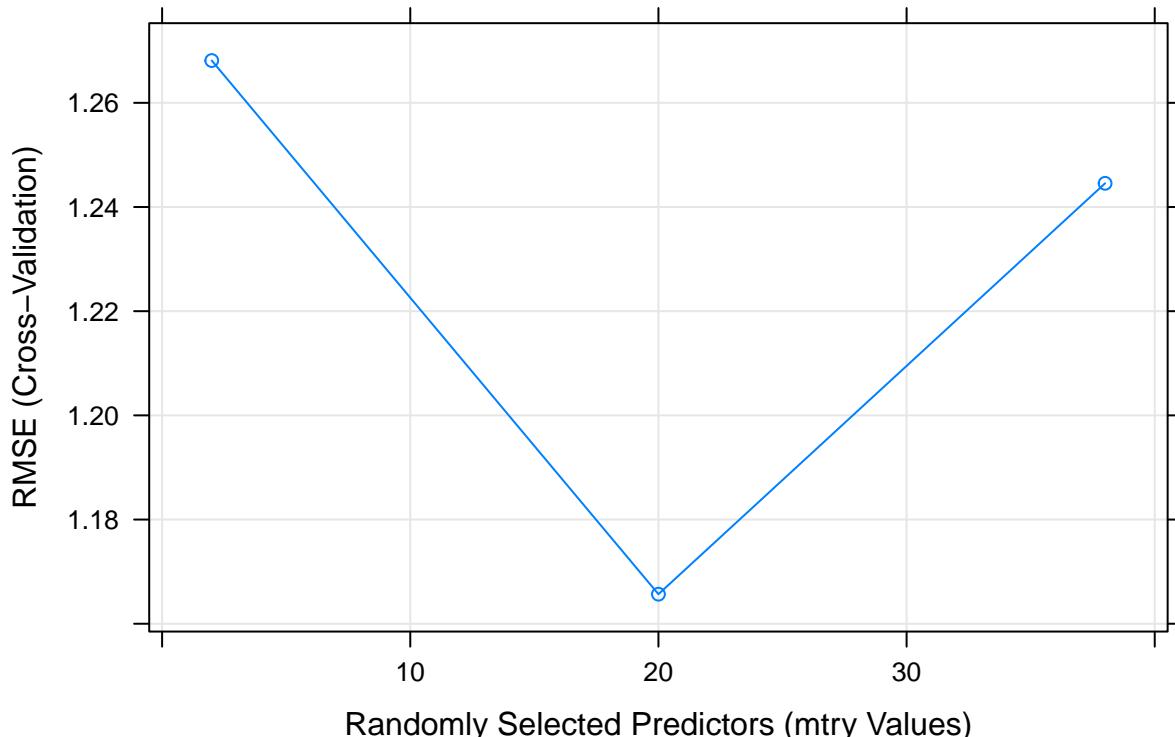
Graphs created: Residuals vs. Fitted, Normal Q-Q, Scale-Location and Residuals vs. Leverage.



For the **Random Forest** model, we have done predictions based on the method “rf” and created plots in order to visualize the results. This model can be considered one of the most common and accurate for predictive models and is defined as an ensemble learning method for classification.

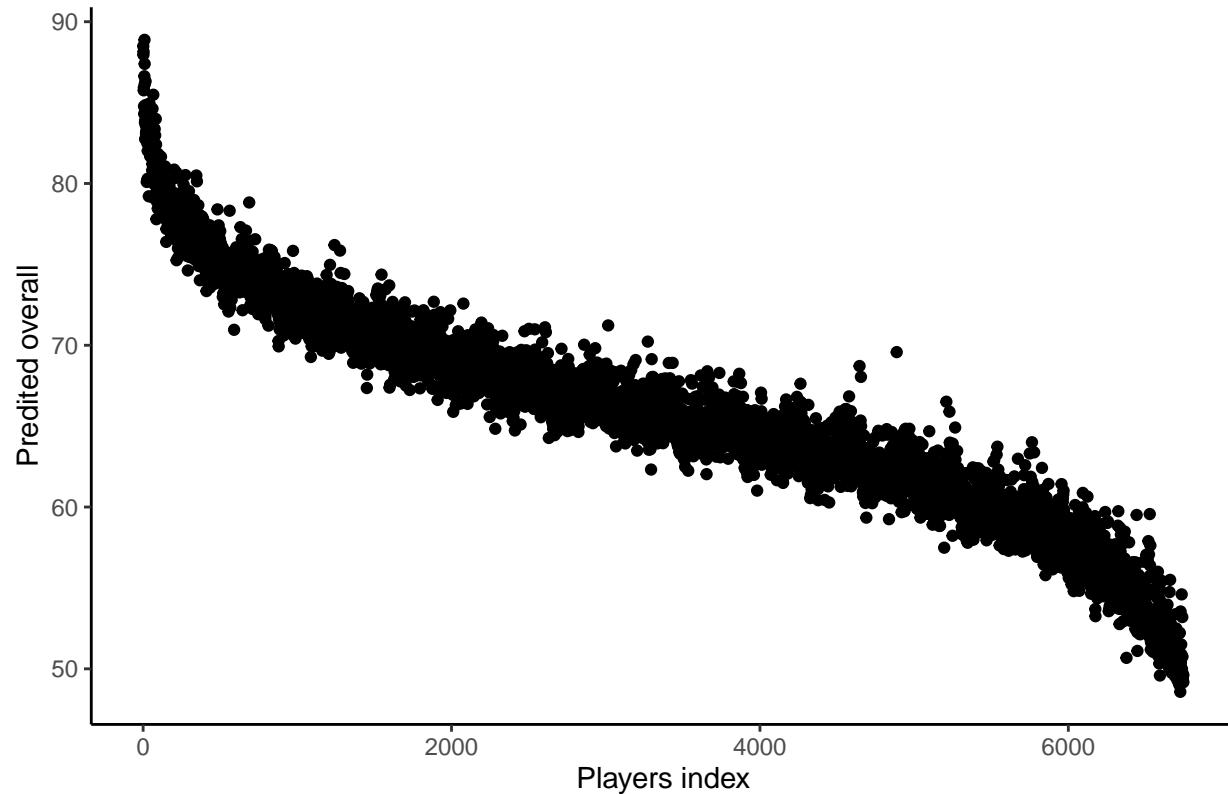
```
## Random Forest
##
## 10115 samples
##     38 predictor
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 9103, 9102, 9104, 9104, 9104, 9105, ...
## Resampling results across tuning parameters:
##
##     mtry   RMSE    Rsquared   MAE
##     2      1.268120 0.9716293 0.9616607
##     20     1.165671 0.9727149 0.8846786
##     38     1.244548 0.9683474 0.9367475
##
## RMSE was used to select the optimal model using the smallest value.
## The final value used for the model was mtry = 20.
```

## Random Forest Results



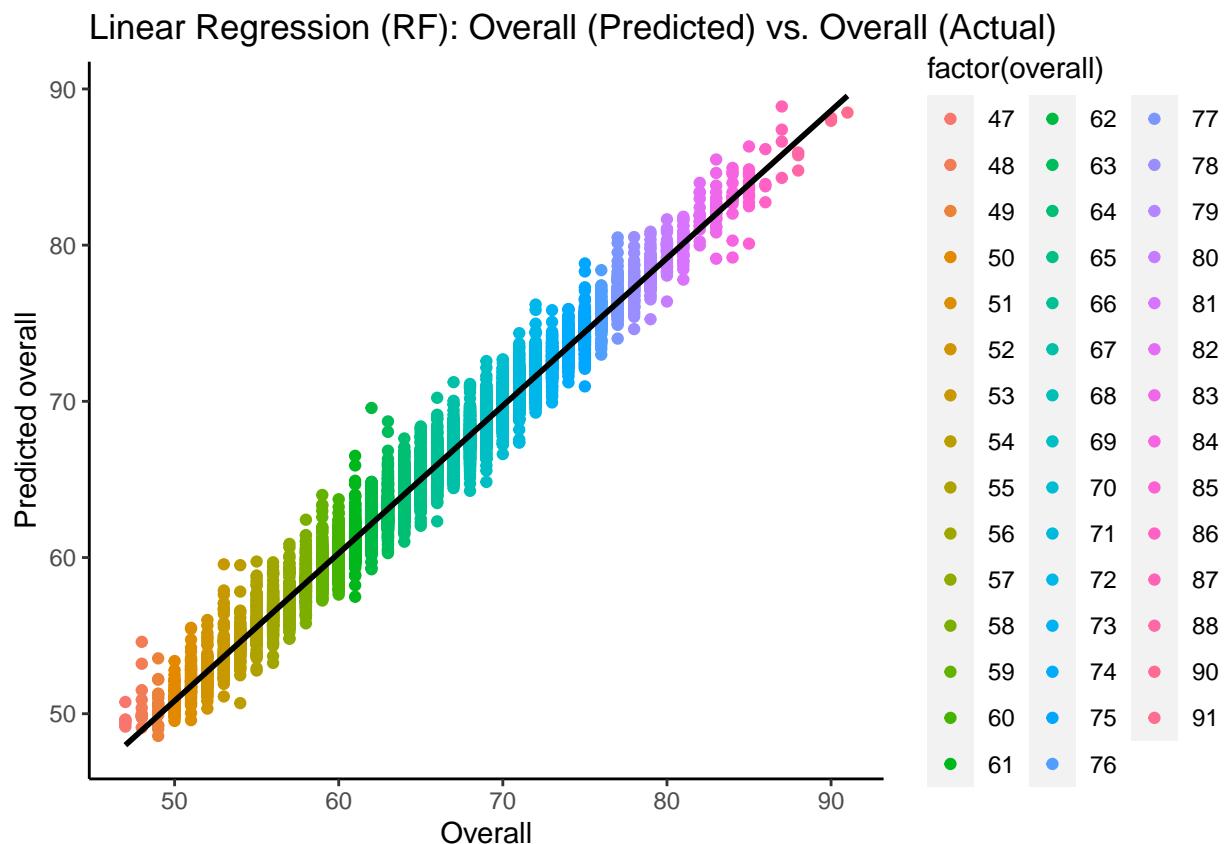
After we made the predictions for the **Random Forest** model, we elaborated a plot in order to visualize the distribution of predicted overalls for all the players in the training set.

Distribution of predicted overalls: Random Forest



In the following plot, we can see a linear regression graph for our **Random Forest** model. This plot show us the relationship between the predicted overall and the actual overall. In this plot we can see an almost perfect correlation/relationship with a value of **0.9447** in this two variables. The reason why this model got a better correlation in comparison to the linear regression model was because the predictions for this model were more precise.

```
##  
## Call:  
## lm(formula = Prediction_RandomForest ~ overall, data = Outcome_RandomForest)  
##  
## Coefficients:  
## (Intercept)      overall  
##       3.5871      0.9447
```



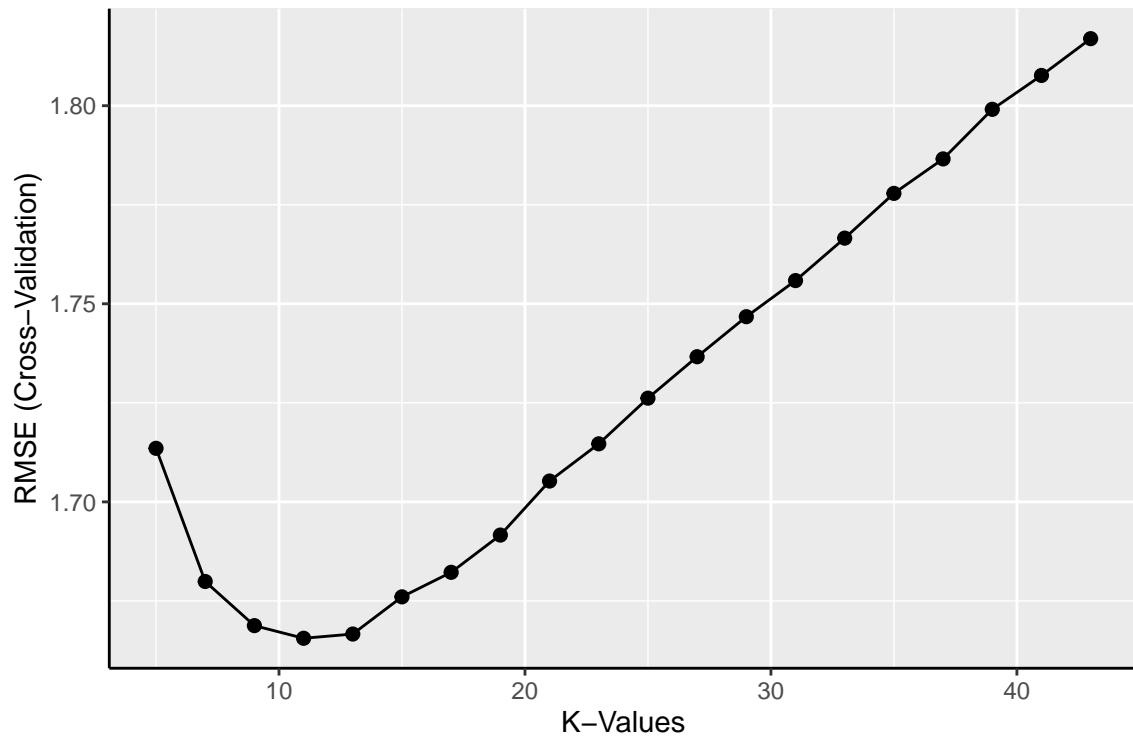
RESULT (RMSE) FOR THE RANDOM FOREST MODEL:

$$\overline{x} \\ \underline{1.14096}$$

For the (**Knn**) **K-nearest neighbors** model, we have done predictions based on the method “knn” and created plots in order to visualize the results. This model can be considered one of the most common for predictive models but is not more accurate than the random forest model. This model can be defined as a supervised machine learning algorithm that can be used to solve both classification and regression tasks.

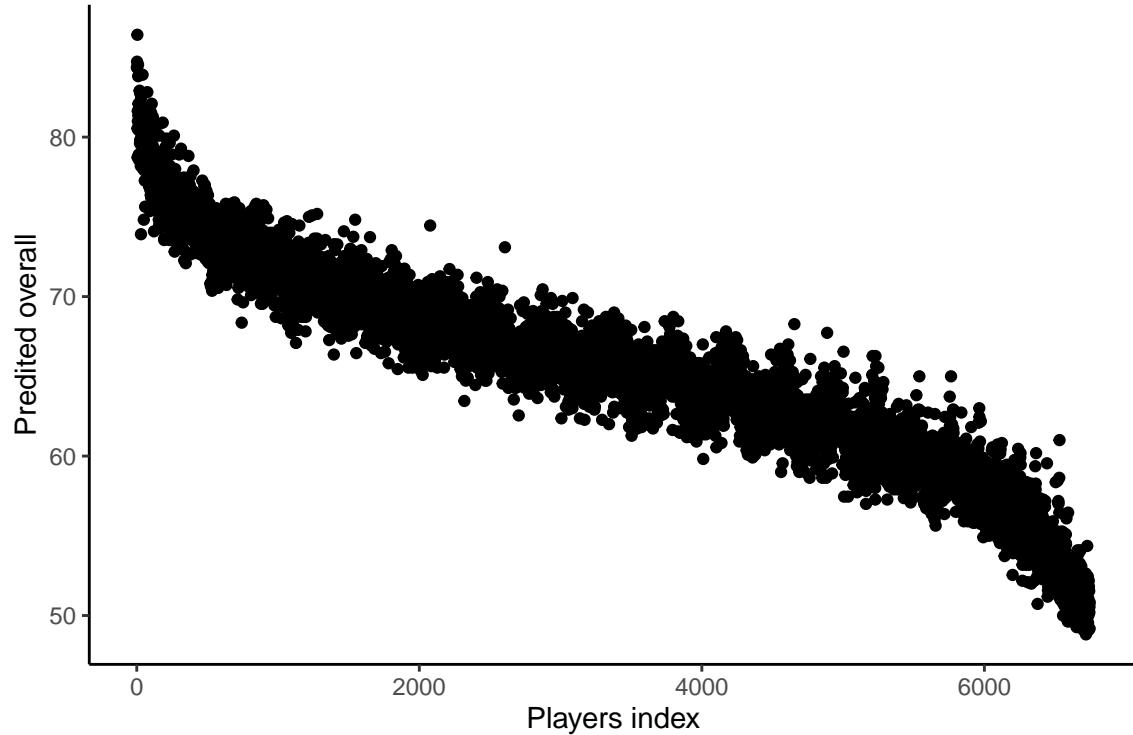
```
## k-Nearest Neighbors
##
## 10115 samples
##     38 predictor
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 9105, 9103, 9104, 9103, 9103, 9105, ...
## Resampling results across tuning parameters:
##
##     k    RMSE    Rsquared   MAE
##      5  1.713521  0.9441665  1.334563
##      7  1.679907  0.9475465  1.306158
##      9  1.668770  0.9496132  1.297946
##     11  1.665579  0.9507730  1.289880
##     13  1.666648  0.9514611  1.289381
##     15  1.676034  0.9514207  1.292446
##     17  1.682223  0.9515792  1.297111
##     19  1.691626  0.9514586  1.304457
##     21  1.705263  0.9511615  1.313350
##     23  1.714656  0.9510088  1.320147
##     25  1.726180  0.9506663  1.328733
##     27  1.736641  0.9504521  1.336610
##     29  1.746748  0.9501983  1.343150
##     31  1.755858  0.9499960  1.350069
##     33  1.766573  0.9496953  1.358895
##     35  1.777869  0.9493380  1.367318
##     37  1.786572  0.9491557  1.372746
##     39  1.799085  0.9487060  1.381453
##     41  1.807636  0.9484921  1.387764
##     43  1.816925  0.9482497  1.394178
##
## RMSE was used to select the optimal model using the smallest value.
## The final value used for the model was k = 11.
```

Knn Model:



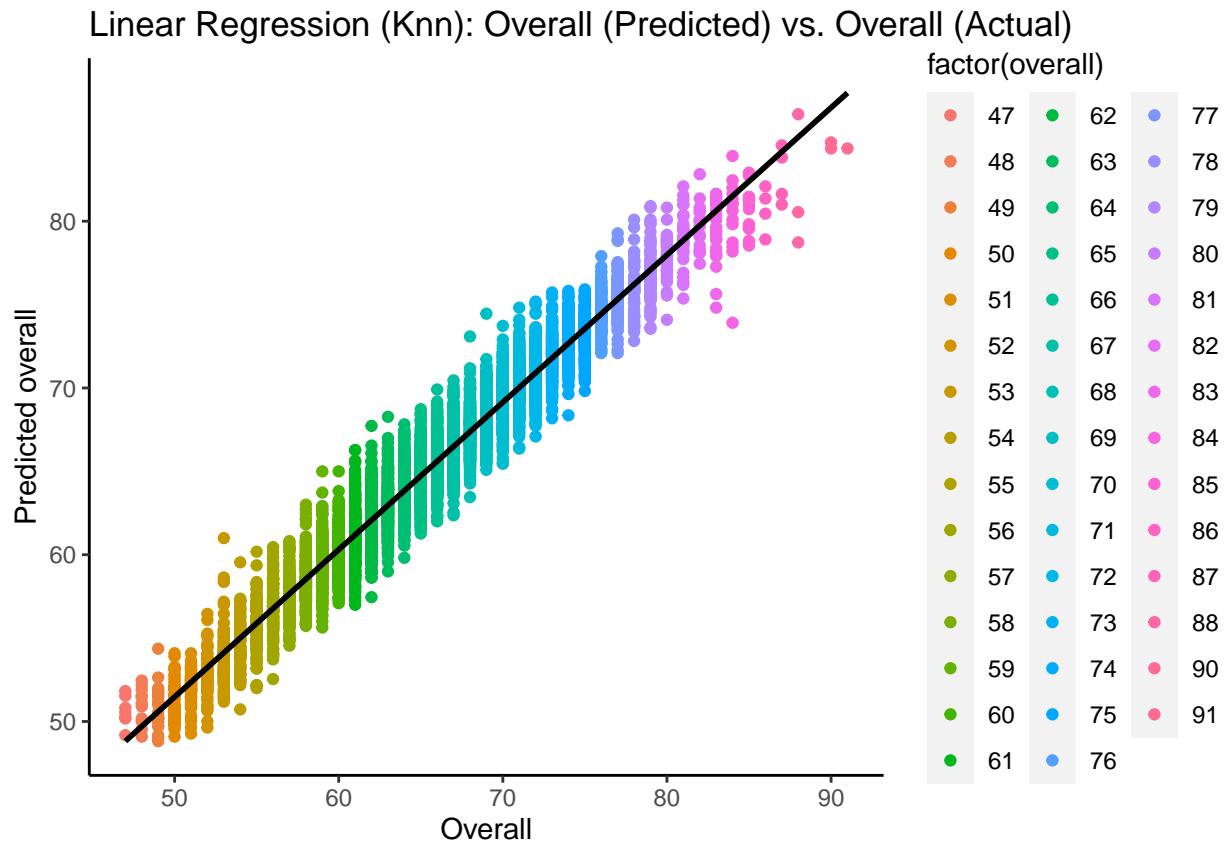
After we made the predictions for the **(Knn) K-nearest neighbors** model, we elaborated a plot in order to visualize the distribution of predicted overalls for all the players in the training set.

Distribution of predicted overalls: Knn



In the following plot, we can see a linear regression graph for our (**Knn**) **K-nearest neighbors** model. This plot show us the relationship between the predicted overall and the actual overall. In this plot we can see a high correlation/relationship with a value of **0.8835** in this two variables.

```
##  
## Call:  
## lm(formula = Prediction_Knn ~ overall, data = Outcome_Knn)  
##  
## Coefficients:  
## (Intercept)      overall  
##       7.2991        0.8835
```

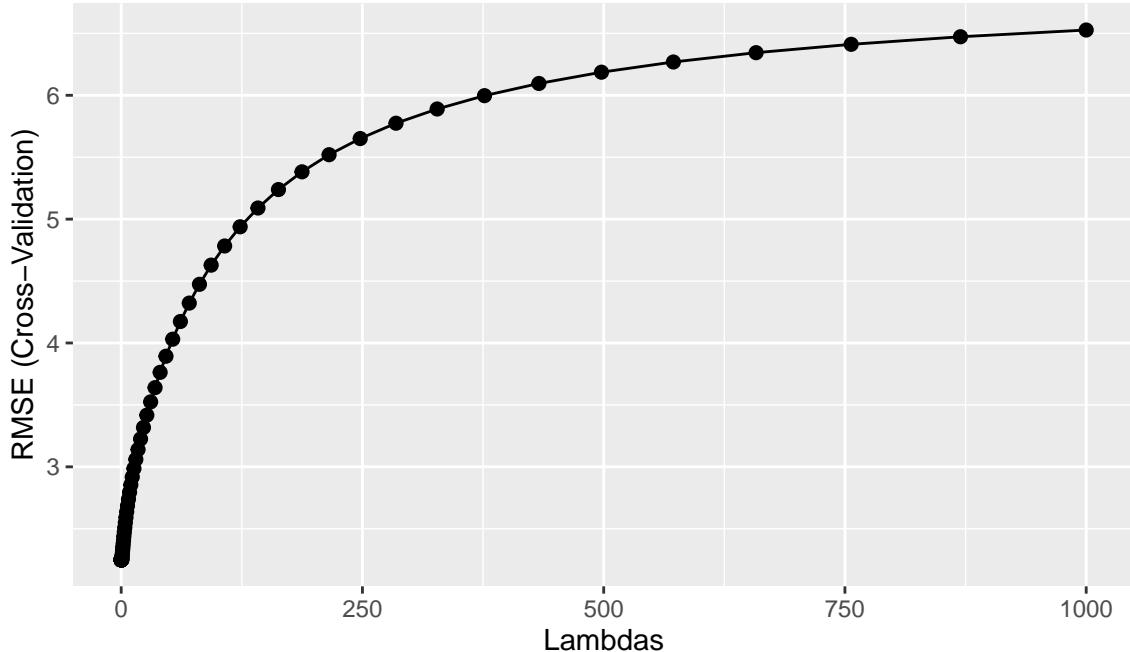


RESULT (RMSE) FOR THE (KNN) K-NEAREST NEIGHBOR MODEL:

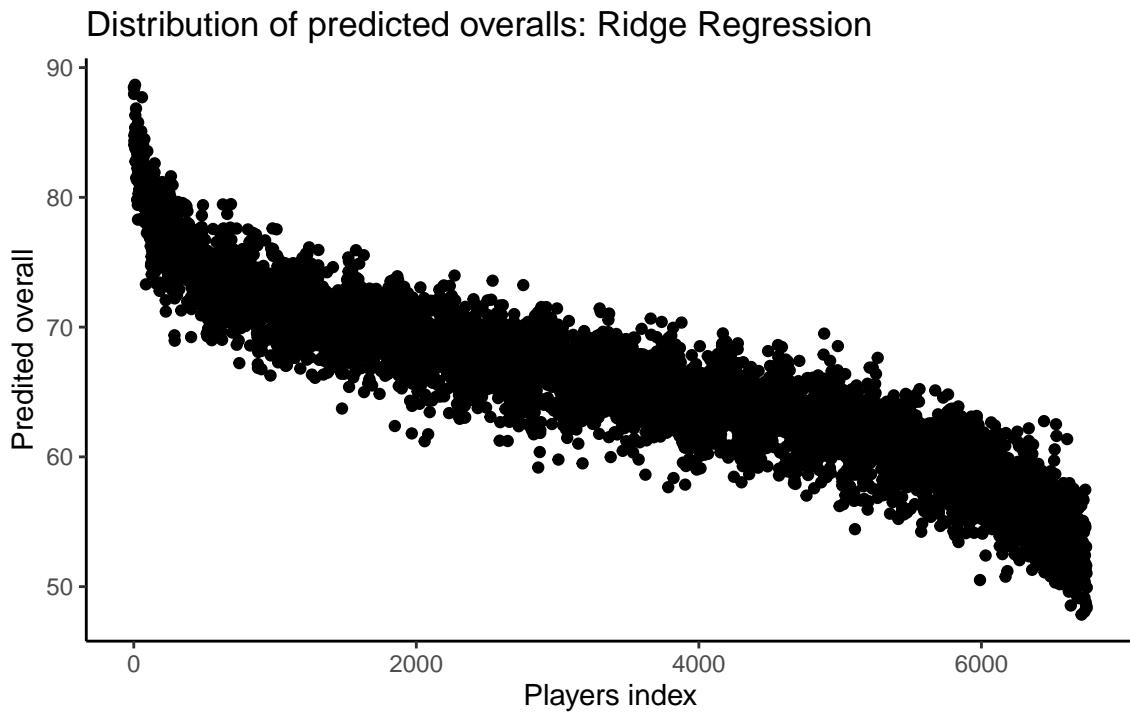
$$\overline{x} \\ 1.652653$$

For the **Ridge Regression** model, we have done predictions based on the method “glm” and created plots in order to visualize the results. This is a model tuning method that is used to analyze any data that suffers from multicollinearity. Also, this model can be considered common for predicting models.

### Ridge regression Model:

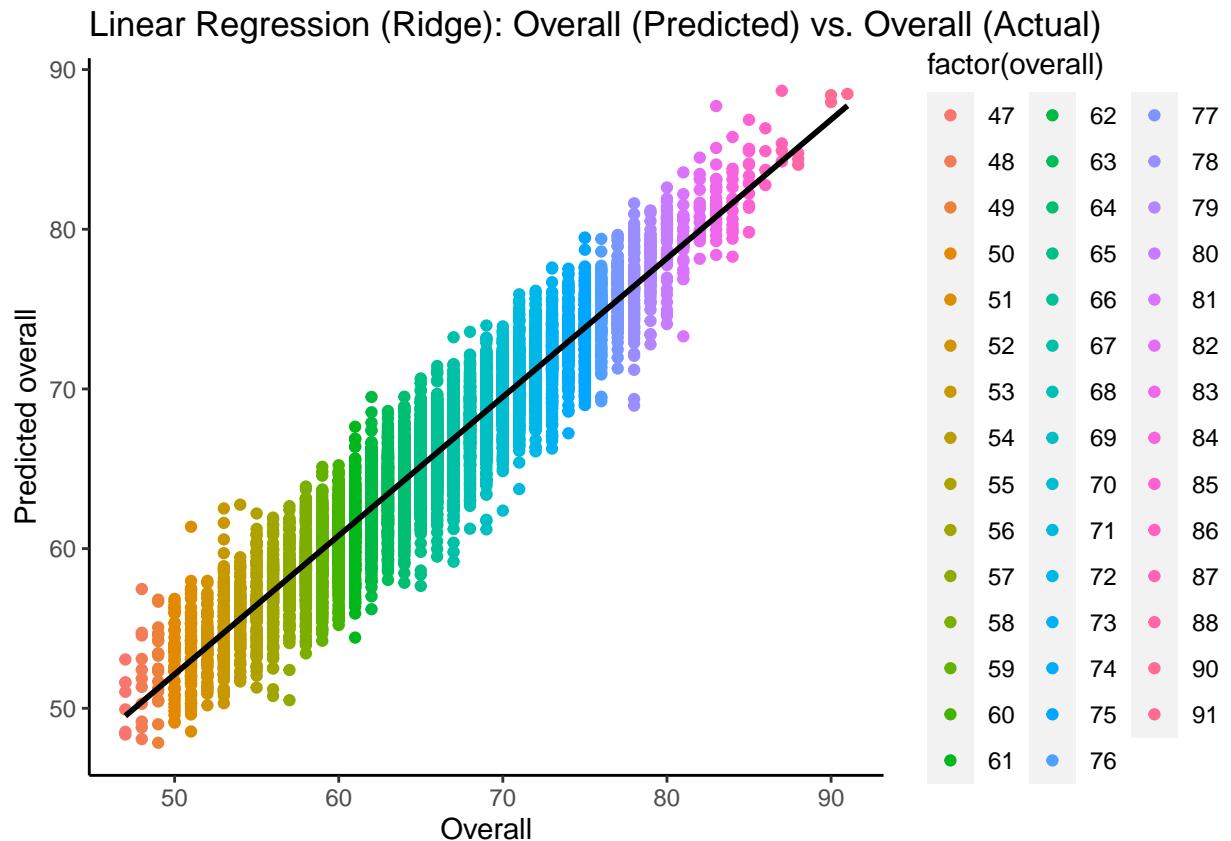


After we made the predictions for the **Ridge Regression** model, we elaborated a plot in order to visualize the distribution of predicted overalls for all the players in the training set.



In the following plot, we can see a linear regression graph for our **Ridge Regression** model. This plot shows us the relationship between the predicted overall and the actual overall. In this plot we can see a high correlation/relationship with a value of **0.8685** in this two variables.

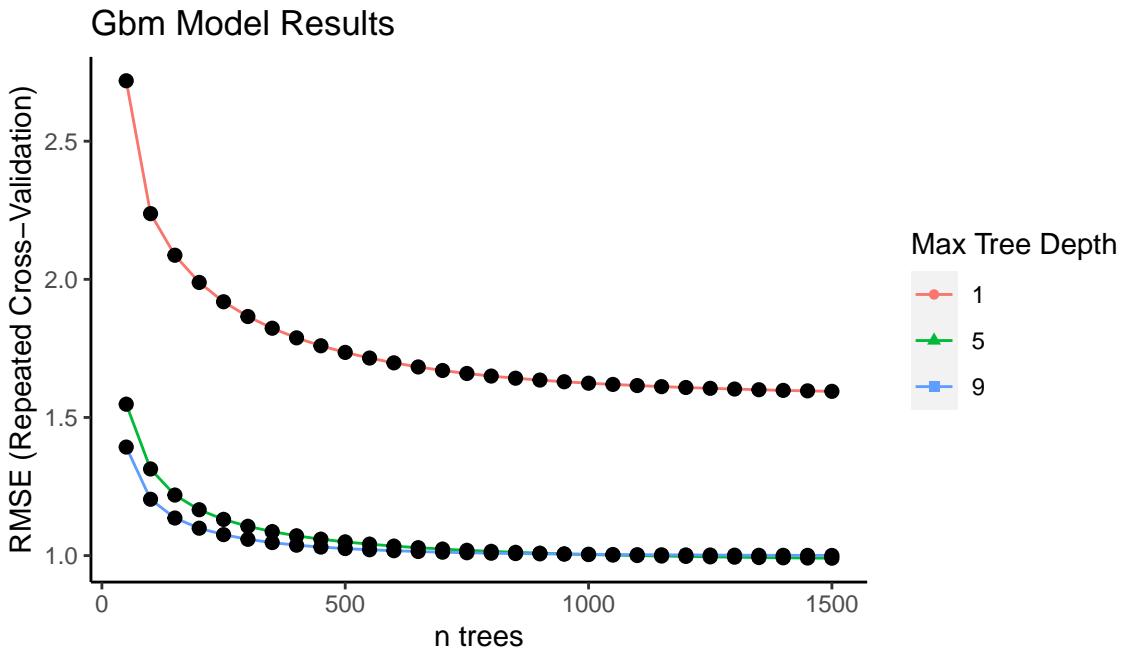
```
##  
## Call:  
## lm(formula = Prediction_Ridge ~ overall, data = Outcome_Ridge)  
##  
## Coefficients:  
## (Intercept)      overall  
##     8.7028        0.8685
```



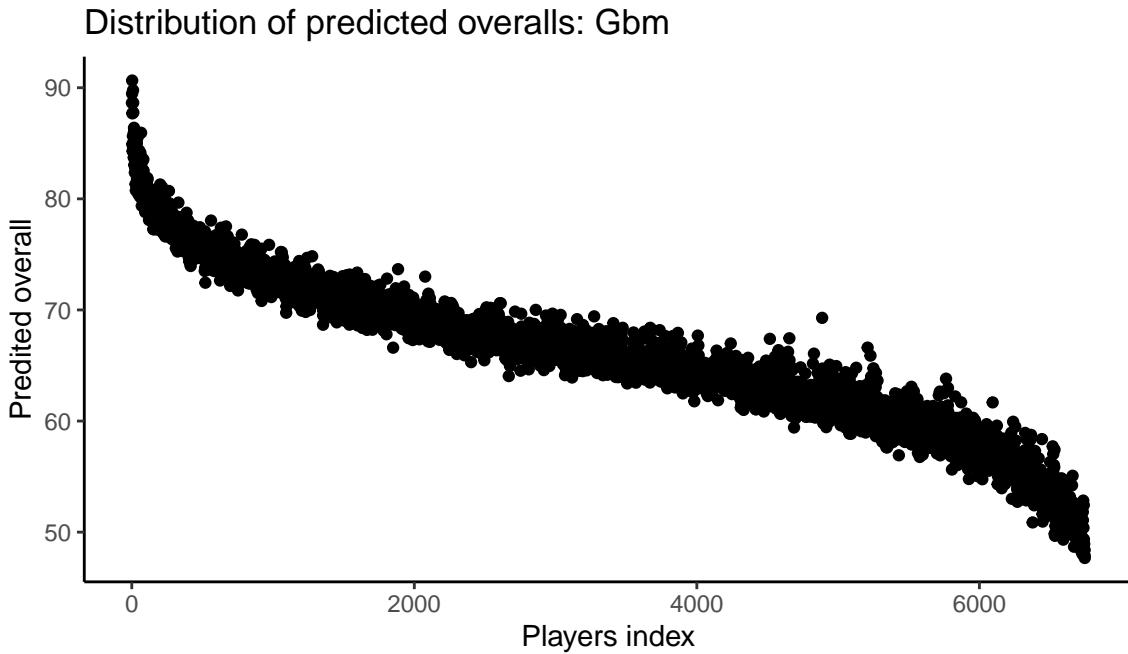
RESULT (RMSE) FOR THE RIDGE REGRESSION MODEL:

$$\overline{x} \\ 2.254696$$

For the **GBM (Gradient Boosting Machine)** model, we have done predictions based on the method “gbm” and created plots in order to visualize the results. This model combines the predictions from multiple decision trees to generate the final predictions. Also, this model can be considered as one of the hardest for predicting models because much of the parameters must be done correctly. However, can be a great option if it’s built in the right form.

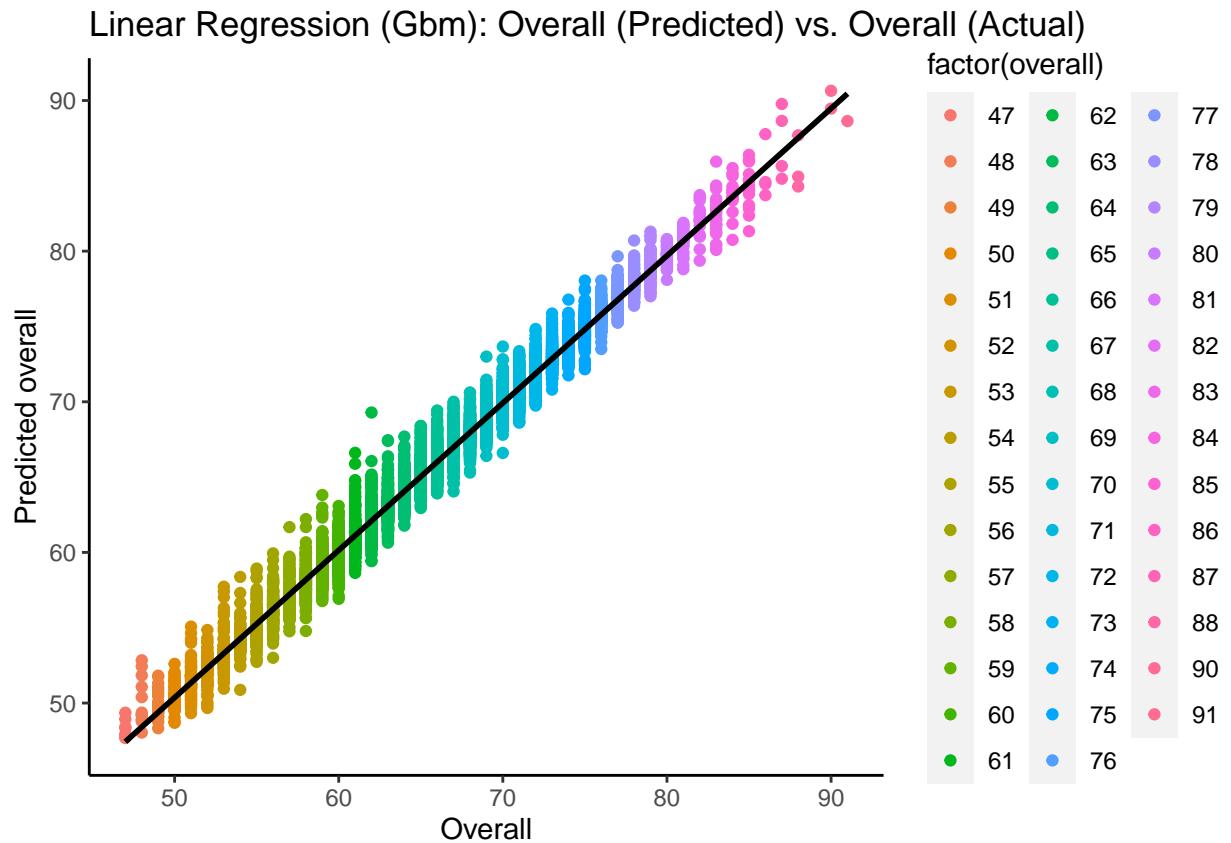


After we made the predictions for the **(GBM) Gradient Boosting Machine** model, we elaborated a plot in order to visualize the distribution of predicted overalls for all the players in the training set.



In the following plot, we can see a linear regression graph for our **(GBM) Gradient Boosting Machine** model. This plot shows us the relationship between the predicted overall and the actual overall. In this plot we can see an almost perfect correlation/relationship with a value of **0.9778** in this two variables.

```
##  
## Call:  
## lm(formula = Prediction_gbm ~ overall, data = Outcome_gbm)  
##  
## Coefficients:  
## (Intercept)      overall  
##       1.4742        0.9778
```



RESULT (RMSE) FOR THE (GBM) GRADIENT BOOSTING MACHINE MODEL:

$$\overline{\overline{x}} \\ \overline{0.9756678}$$

### Results and interpretation for the models:

- The “Linear regression” model evaluated in the validation set gave us a RMSE of **2.24**. The performance of this model can be considered poor and not very precise.
- The “Random Forest” model evaluated in the validation set gave us a RMSE of **1.14**. The performance of this model can be considered very good and a lot more precise in comparison to the linear regression model.
- The “K-nearest neighbors” model evaluated in the validation set gave us a RMSE of **1.65**. The performance of this model can be considered good and more precise than the linear regression model. However, the random forest model realized a better performance than this model.
- The “Ridge regression” model evaluated in the validation set gave us a RMSE of **2.25**. The performance of this model can be considered poor and not very precise, actually this model has the worst RMSE of all the models.
- The “Gradient Boosting Machine” model evaluated in the validation set gave us a RMSE of **0.98**. The performance of this model can be considered great and very precise. In fact, this model has surpass by a lot the other models RMSE results.

Model	RMSE
Linear Regression	2.2406356
Random Forest	1.1409596
Knn	1.6526533
Ridge Regression	2.2546961
Gbm	0.9756678

### Conclusions:

-Based on the results obtained, the (**GBM**) **Gradient Boosting Machine model**, which was more tuned and worked, achieved a lower RMSE than the others models. In consequence, we can conclude that, the more tuned the model is, the more exact results will be.

-After comparing the linear regression graphs for all the models, it was possible to observe that the more accurate the prediction is, the more **correlation/relationship** will exist for the players overall.

-After analyzing the results, it is possible to conclude that the **“GBM (Gradient Boosting Machine) model** was the one that gave us the lowest RMSE, which allowed us to determine that this model was the most accurate and precise in terms of predictions.

-Finally, in the process of creating predictions for the FIFA 2022 players overall, a limitation was found, this limitation was a column named **“defending marking”** that had NA VALUES, this column was very important because it could have contribute a lot in the process of building and evaluating the predictions. However, predictions were good despite the fact that this limitation make our models be less accurate. For future work purposes, I recommend to fill the columns with NA values and omit irrelevant information like, for example: players traits, body type, real face, etc.

## Bibliography:

- Donges, N. (2020, September 3). A complete guide to the random forest algorithm. Built In. <https://builtin.com/data-science/random-forest-algorithm>
- Glen. S (2017, July 29) Ridge Regression; Simple Definition. Retrieved from: <https://www.statisticshowto.com/ridge-regression/>
- Harrison, O. (2019, July 14). Machine Learning Basics with the K-Nearest Neighbors Algorithm. Medium. <https://towardsdatascience.com/machine-learning-basics-with-the-k-nearest-neighbors-algorithm-6a6e71d01761#:~:text=KNN%20works%20by%20finding%20the,in%20the%20case%20of%20regression.>
- Irizarry, R. A. (2021, 20 february). Introduction to Data Science. <https://rafalab.github.io/dsbook/>. <https://rafalab.github.io/dsbook/>
- Murphy, R. (2019, September 12). FIFA player ratings explained: How are the card number & stats decided? Goal.Com. <https://www.goal.com/en-ae/news/fifa-player-ratings-explained-how-are-the-card-number-stats/1hszd2fgr7wgf1n2b2yjdpwynu>
- Panchotia, R. (2020, August 5). Predictive Modelling Using Linear Regression - The Startup. Medium. [https://medium.com/swlh/predictive-modelling-using-linear-regression-e0e399dc4745#:~:text=Linear%20regression%20is%20one%20of,given%20predictor%20variable\(s\).](https://medium.com/swlh/predictive-modelling-using-linear-regression-e0e399dc4745#:~:text=Linear%20regression%20is%20one%20of,given%20predictor%20variable(s).)
- ZevRoss (2018, October 2). Predictive modeling and machine learning in R with the caret package. Technical Tidbits From Spatial Analysis & Data Science. <http://zevross.com/blog/2017/09/19/predictive-modeling-and-machine-learning-in-r-with-the-caret-package/>