



UNIVERSITY OF TORONTO
SCHOOL OF CONTINUING STUDIES

Exploratory Analysis European League Football

Kaustubh Saha
Steven Tong
Raka Dhar
Joydip Dutta



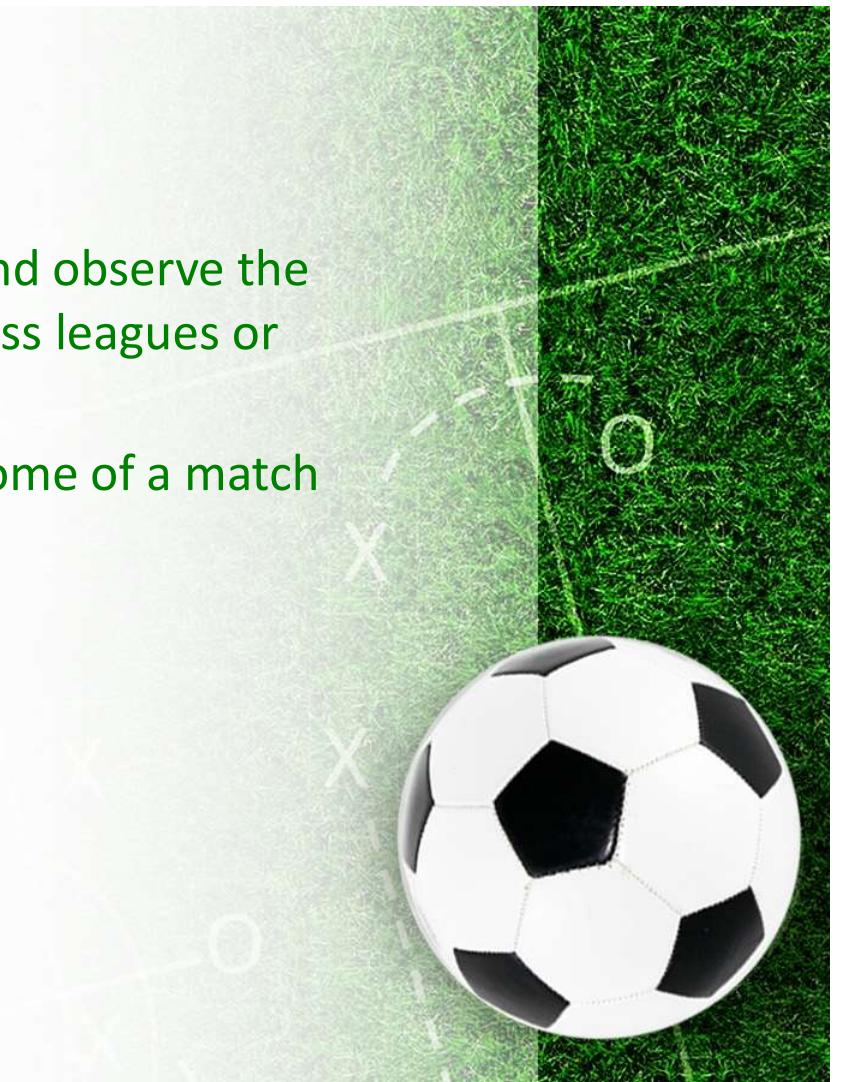
Speakers

- Data source, structure and cleanup- Raka
- Feature analysis- Joydip
- Factors and impact- Steven
- Event Data analysis- Kaustubh



Objective

- Explore the various features in the data and observe the trend in movement of those features across leagues or across years
- Identify the features influencing the outcome of a match
- Build a predictive model and test it



Data Source – League match

EPL : <https://datahub.io/sports-data/english-premier-league>

LaLiga : <https://datahub.io/sports-data/spanish-la-liga>

Serie A : <https://datahub.io/sports-data/italian-serie-a>

BundesLiga : <https://datahub.io/sports-data/german-bundesliga>

Ligue One : <https://datahub.io/sports-data/french-ligue-1>

Data Source – EU events

Kaggle (<https://www.kaggle.com/secareanualin/football-events>)

Structure of data- Match & Event

Date
HomeTeam
AwayTeam
FullTime_HomeTeam_Goals
FullTime_AwayTeam_Goals
FullTime_Result
HalfTime_HomeTeam_Goals
HalfTime_AwayTeam_Goals
HalfTime_Result
HomeTeam_Shots
AwayTeam_Shots
HomeTeam_ShotsOnTarget
AwayTeam_ShotsOnTarget
HomeTeam_FoulsCommitted
AwayTeam_FoulsCommitted
HomeTeam_Corners
AwayTeam_Corners
HomeTeam_YellowCards
AwayTeam_YellowCards
HomeTeam_RedCards
AwayTeam_RedCards

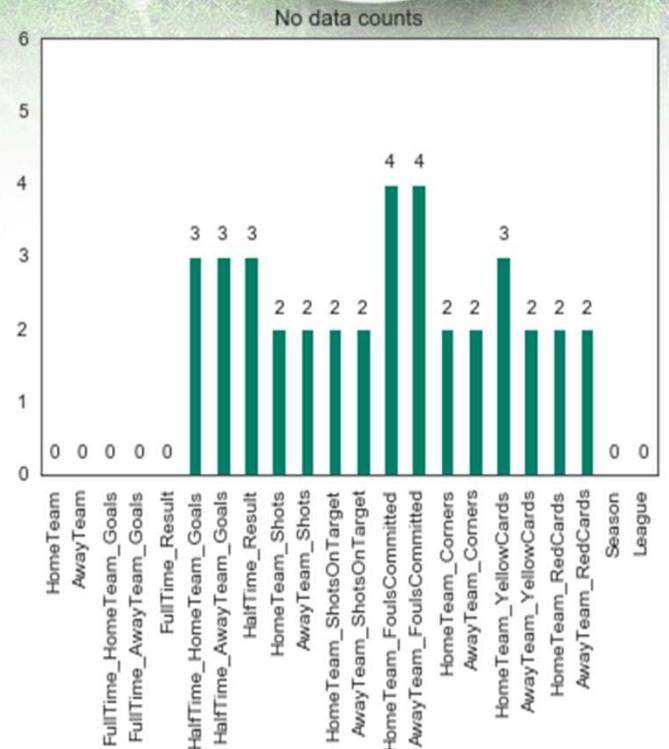
id_odsp
link_odsp
adv_stats
date
league
season
country
ht
at
ftgh
ftag
odd_h
odd_d
odd_a
odd_over
odd_under
odd_bts
odd_bts_n

Data Load and Prep

- Renamed columns to make it intuitive.
- Changed index to date of match played.
- Built a function to load data.
- To begin with, we loaded one of the files and see if the data layout meets our expectations and followed the same approach for all the 5 of them and all of the data has similar structure.
- Concatenated all league match data for each of the leagues including 2018-19 season which is still underway into a single data frame.

Data Cleansing

- There were no duplicate records!
- Missing records populated as follows-
 - ✓ Half time result was in line with full time result
 - ✓ Half time goals is half of full time goals.
 - ✓ Total number of shots on target is total number of goals.
 - ✓ Fouls, corners, red/yellow cards were taken as average of other matches for the same team.



Engineered Data

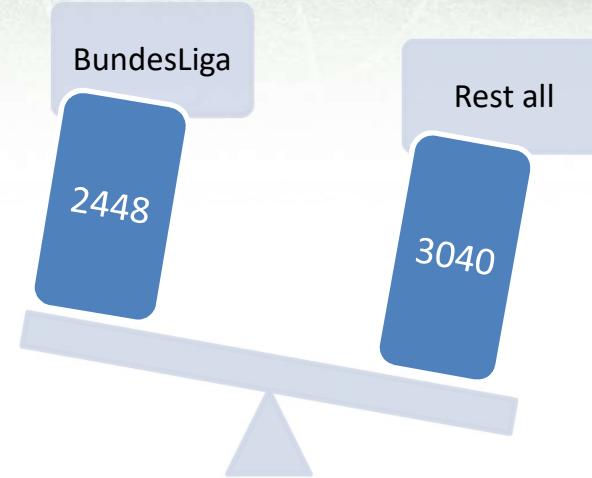
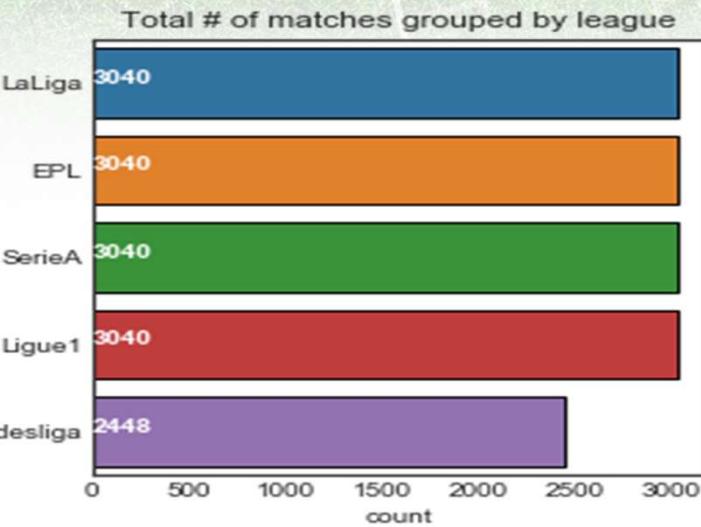
- **Winner** : if FullTime_Result = 'H', then its Home Team. if FullTime_Result = 'A', then its Away Team. Else its null (implying draw)
- **Loser** : if FullTime_Result = 'H', then its Home Team. if FullTime_Result = 'A', then its Away Team. Else its null (implying draw)
- **Total cards (Home Team)** = Home Team yellow cards + Home Team red cards
- **Total cards (Away Team)** = AT yellow cards + AT red cards
- **Total cards** = Total cards (HT) + Total cards (AT)
- **% of shots on target (Home Team)** = HT shots on target / HT total shots
- **% of shots on target (Away Team)** = AT shots on target / AT total shots
- **Home team goal saves** = AT shots on target – Full time AT goals
- **Away team goal saves** = HT shots on target – Full time HT goals
- **Total Goals** = Full time HT goals + Full time AT goals

Final data structure

HomeTeam
AwayTeam
FullTime_HomeTeam_Goals
FullTime_AwayTeam_Goals
FullTime_Result
HalfTime_HomeTeam_Goals
HalfTime_AwayTeam_Goals
HalfTime_Result
HomeTeam_Shots
AwayTeam_Shots
HomeTeam_ShotsOnTarget
AwayTeam_ShotsOnTarget
HomeTeam_FoulsCommitted
AwayTeam_FoulsCommitted
HomeTeam_Corners
AwayTeam_Corners

HomeTeam_YellowCards
AwayTeam_YellowCards
HomeTeam_RedCards
AwayTeam_RedCards
Season
League
Winner
Loser
TotalGoals
HomeTeam_TotalCards
AwayTeam_TotalCards
TotalCards
HomeTeam_ShotsOnTarget_Percent
AwayTeam_ShotsOnTarget_Percent
HomeTeam_GoalSaves
AwayTeam_GoalSaves

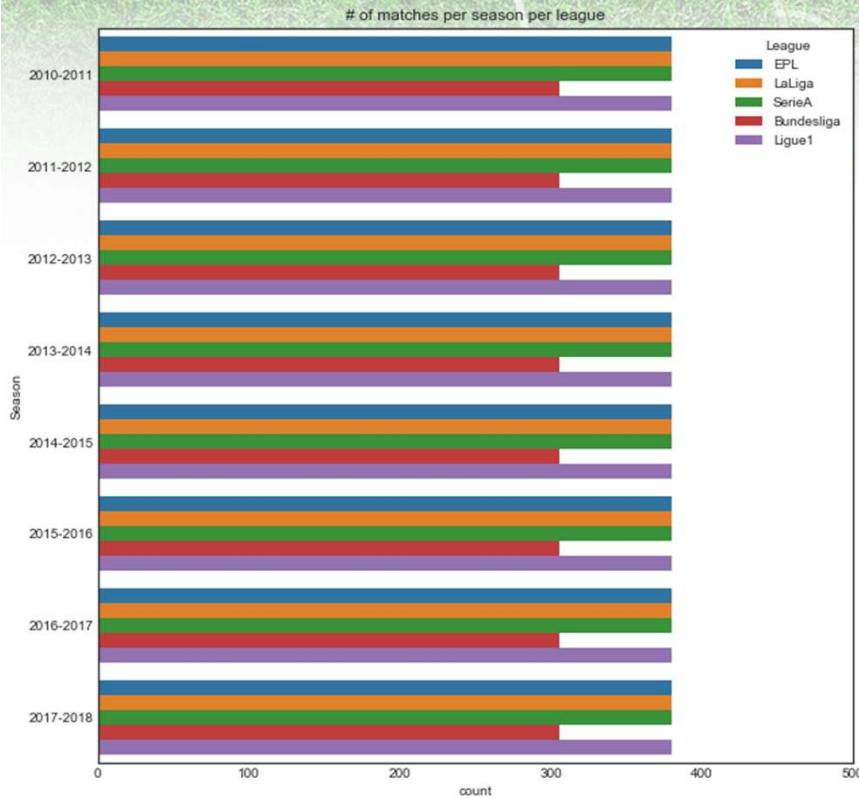
Total number of matches by League



Lets have a look at the number of matches per season per league. If its not uniform, then we might need to adjust the aggregate data accordingly before comparing (We'll ignore the current season as not every league starts at the same time)

Clearly Bundesliga has lesser number of matches as compared to other leagues.

Team and number of matches

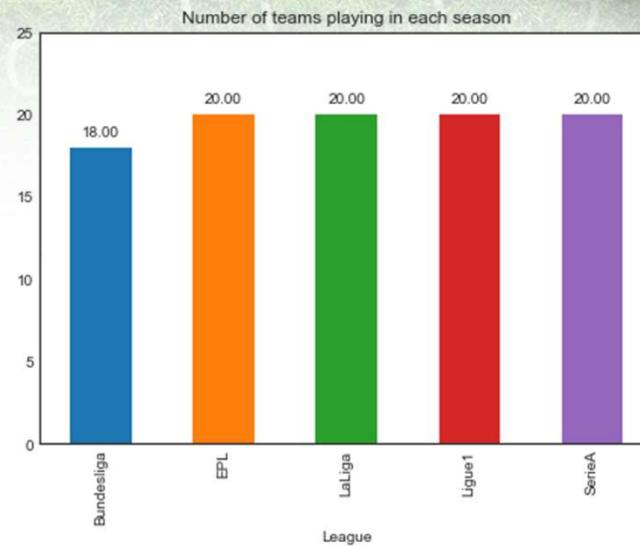
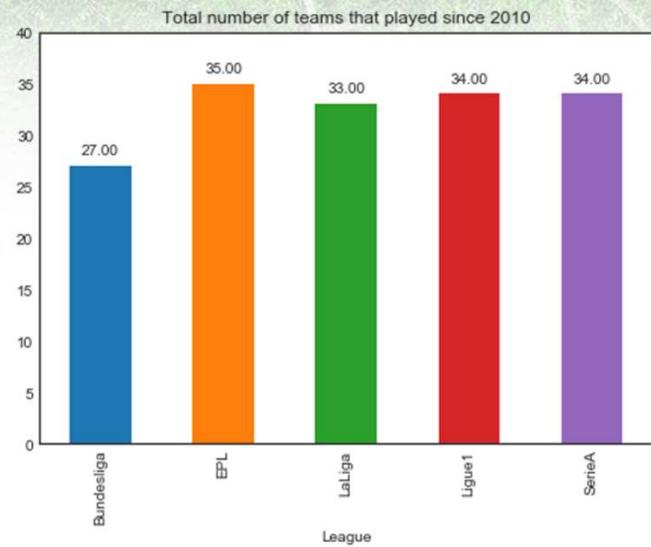


- ✓ For a particular league, the number of matches per season has remained consistent over the years.
- ✓ However when we compare other league aggregate data with Bundesliga, the data might need some adjustment.

BundesLiga data adjustment

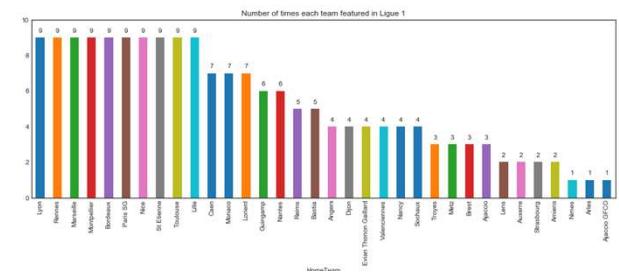
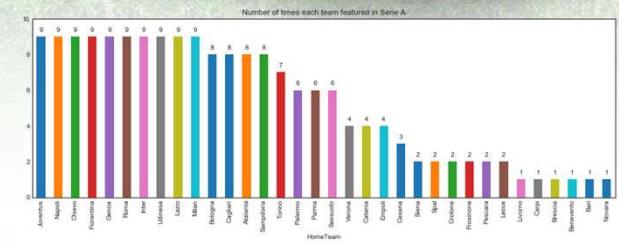
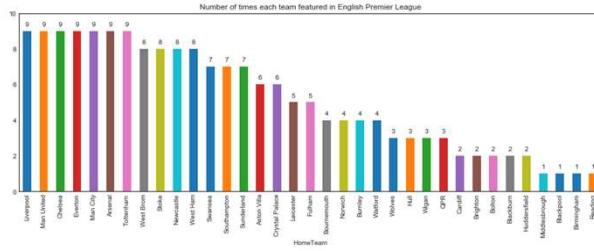
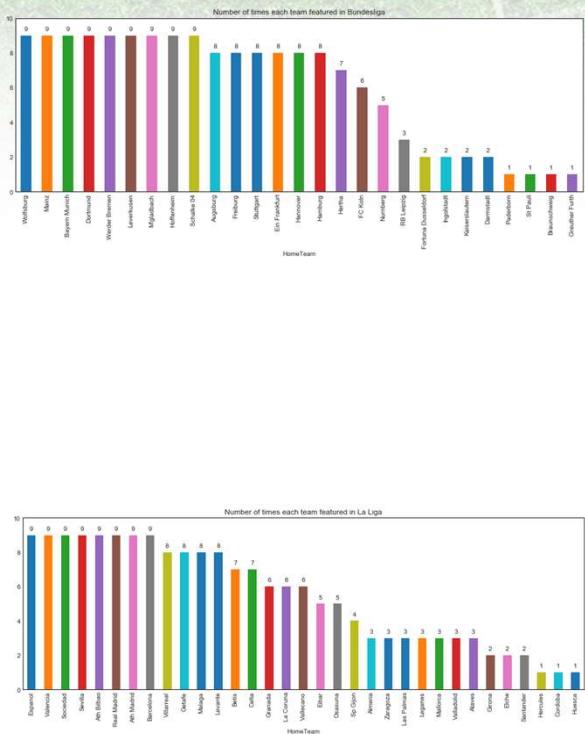
- ✓ Bundesliga has 18 teams and 306 games (34 games for each team) per season whereas other leagues have 20 teams and 380 games a season (38 games for each team)
- ✓ Bundesliga has been scaled up accordingly in order to do a fair comparison wherever aggregate was compared across leagues
 - BUNDESLIGA_TEAM_SCALEUP = 38/34
 - BUNDESLIGA_SEASON_SCALEUP = 380/306

Feature analysis – Team and number of matches

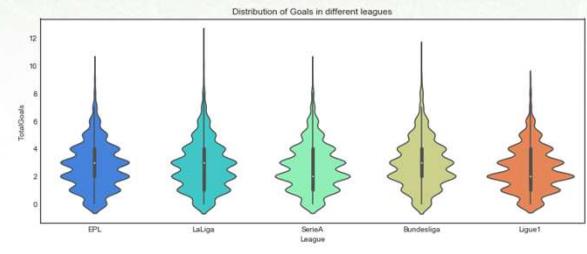
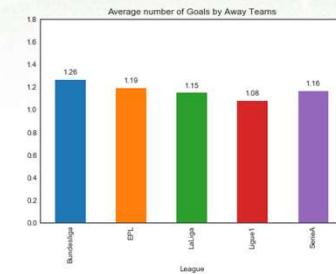
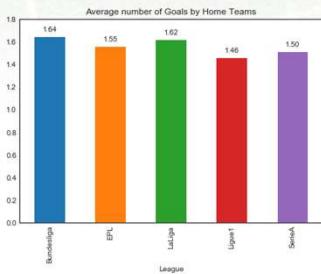
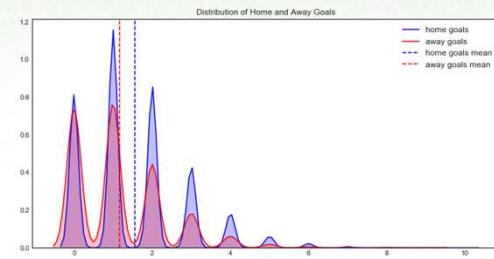


The difference between the two numbers indicates that there were occasions where teams were relegated to lower tier league or top teams from lower tier leagues were promoted

Feature analysis – Relegation



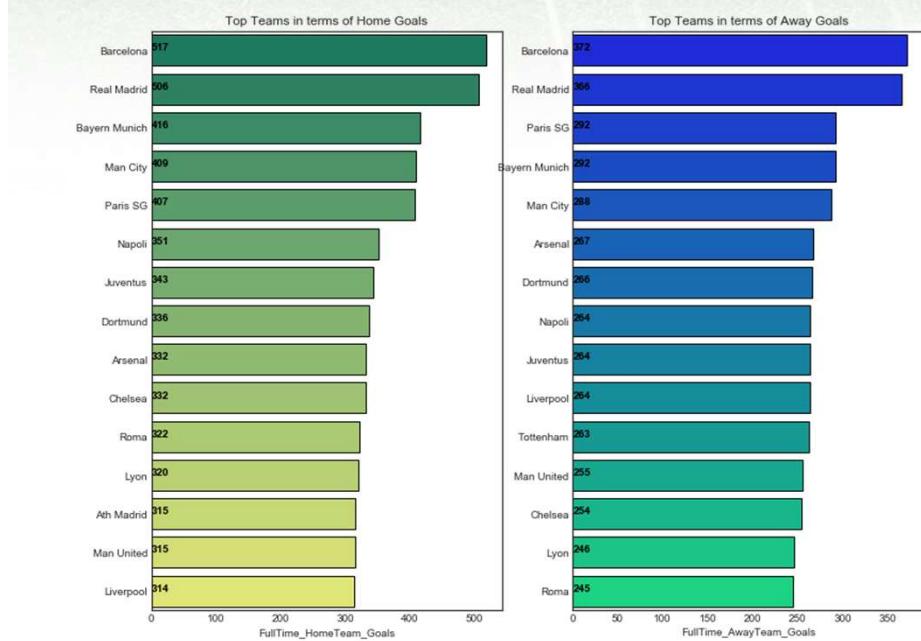
Feature analysis – Goals



- ✓ The mean value for home goals is higher than that of away goals. So there might be some home advantage.

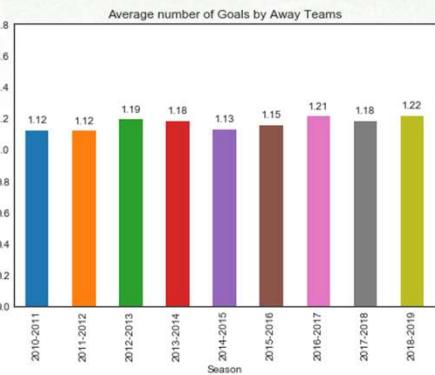
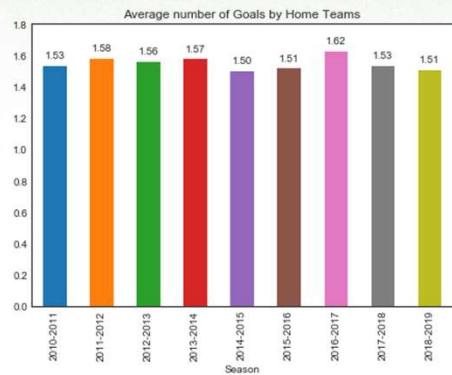
- ✓ Total number of goals vary between 1 and 4. It's quite rare to see more than 6 goals in a match in any league

Feature analysis – Top teams and Goals



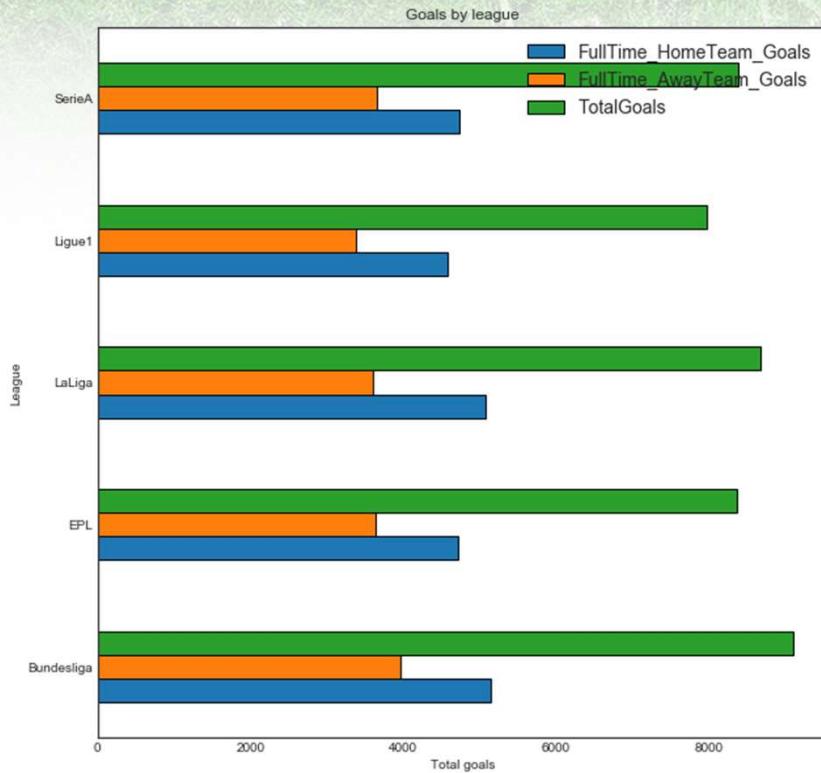
- ✓ These two lists have a lot of teams in common. This indicates that good teams usually do well at home as well as away (However they score better at home than at away).

Feature analysis – Over the year performance



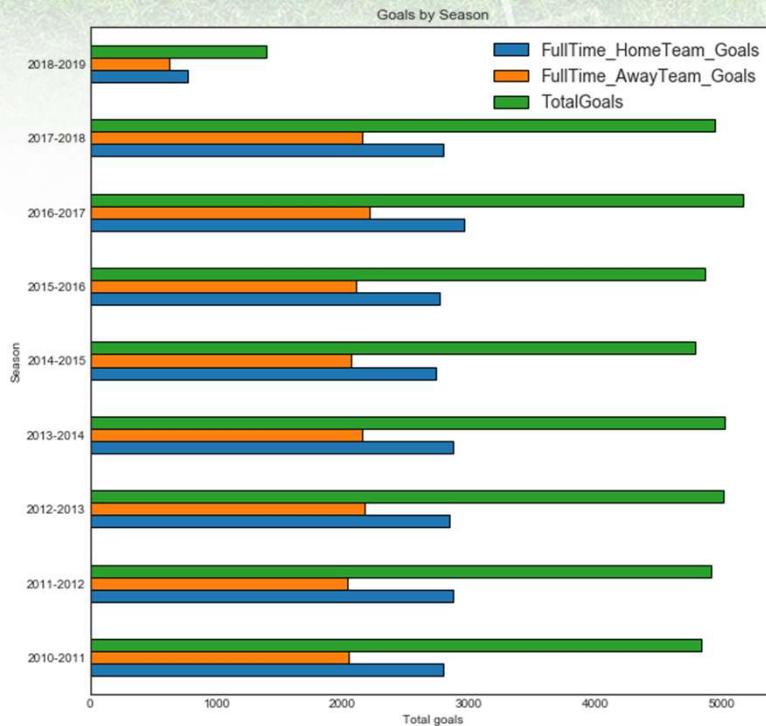
- ✓ Over 10 years, away team have marginally improved from 1.12 goals to 1.22 goals per match.

Feature analysis – Total goals for all seasons



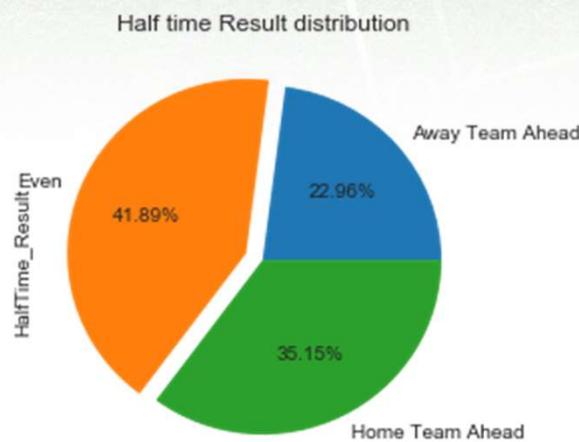
- ✓ Ligue One produces less goals as compared to other leagues. Post adjustment, Bundesliga produces more goals than any other leagues.

Feature analysis – Season wise goals



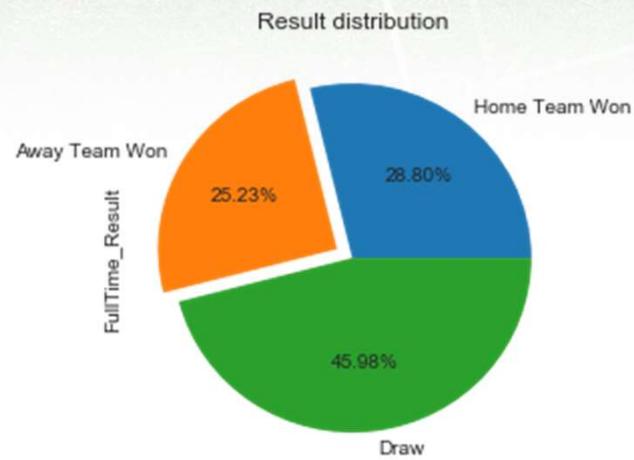
- ✓ From 2010-11 to 2013-14, number of goals increased each season but stopped growing after 2013-14 season (except for a spike in 2016-17)

Feature analysis – Half time result



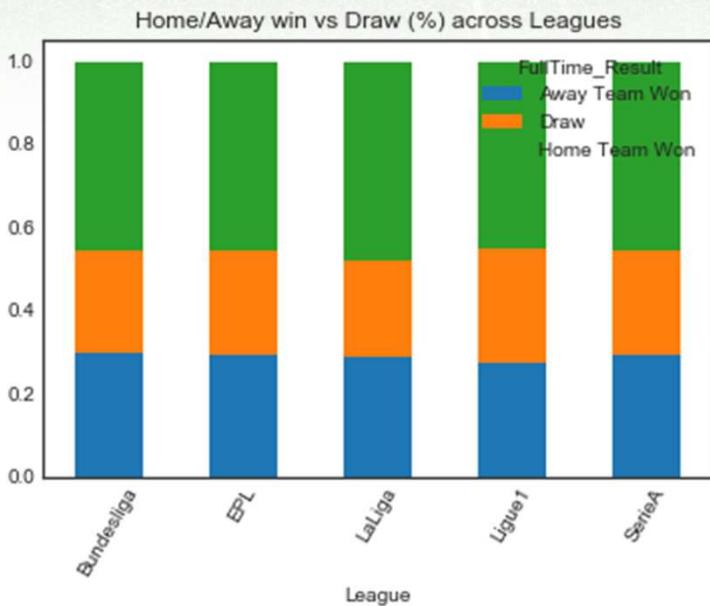
- ✓ Home team has advantage right from the beginning.

Feature analysis – Full time result



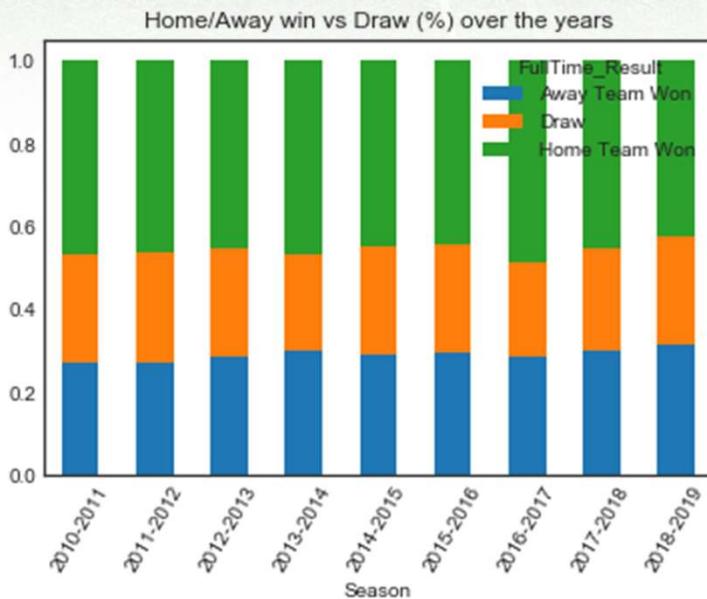
- ✓ Home team advantage fell marginally but still overall result favors home team

Feature analysis – Result variation across leagues



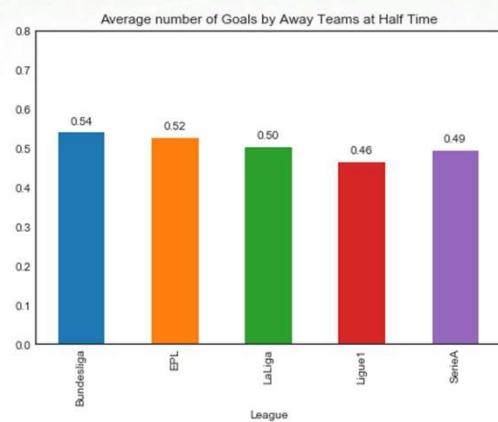
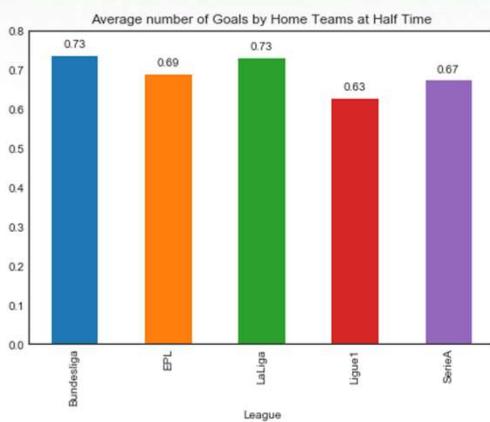
- ✓ Home team wins more than 40% of matches in every league.
- ✓ Home advantage seems to be more prevalent in La Liga.
- ✓ Ligue1 produces more draws than other leagues.

Feature analysis – Result variation across seasons



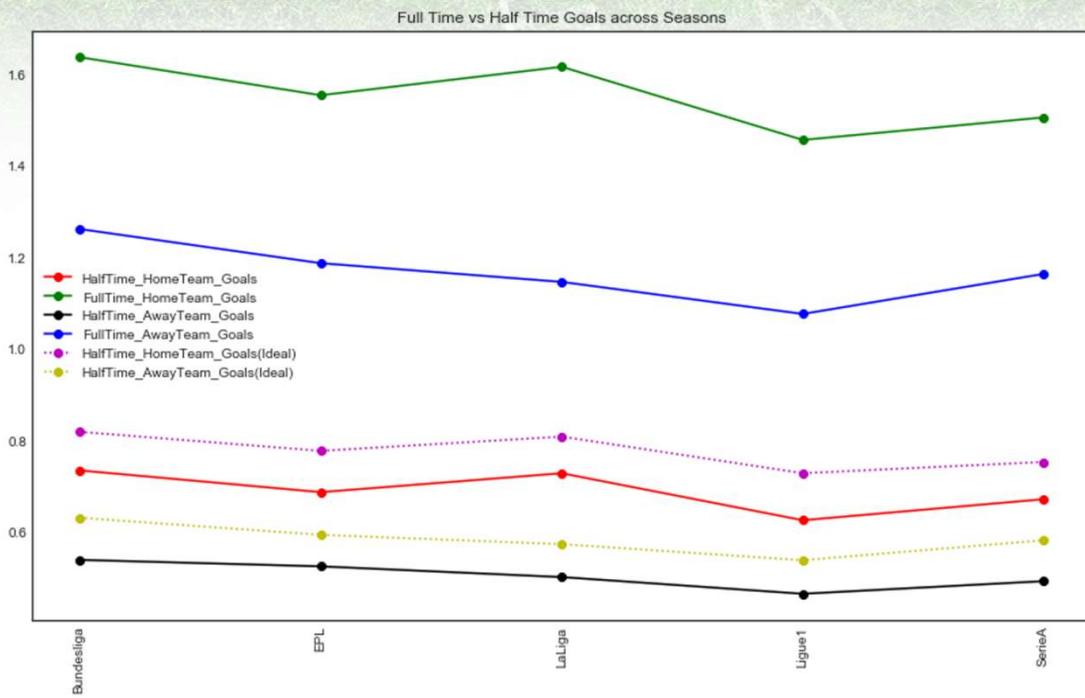
- ✓ No significant change.
- ✓ Home advantage seen across all seasons.
- ✓ Away teams' performance improved marginally.

Feature analysis – Half time goals



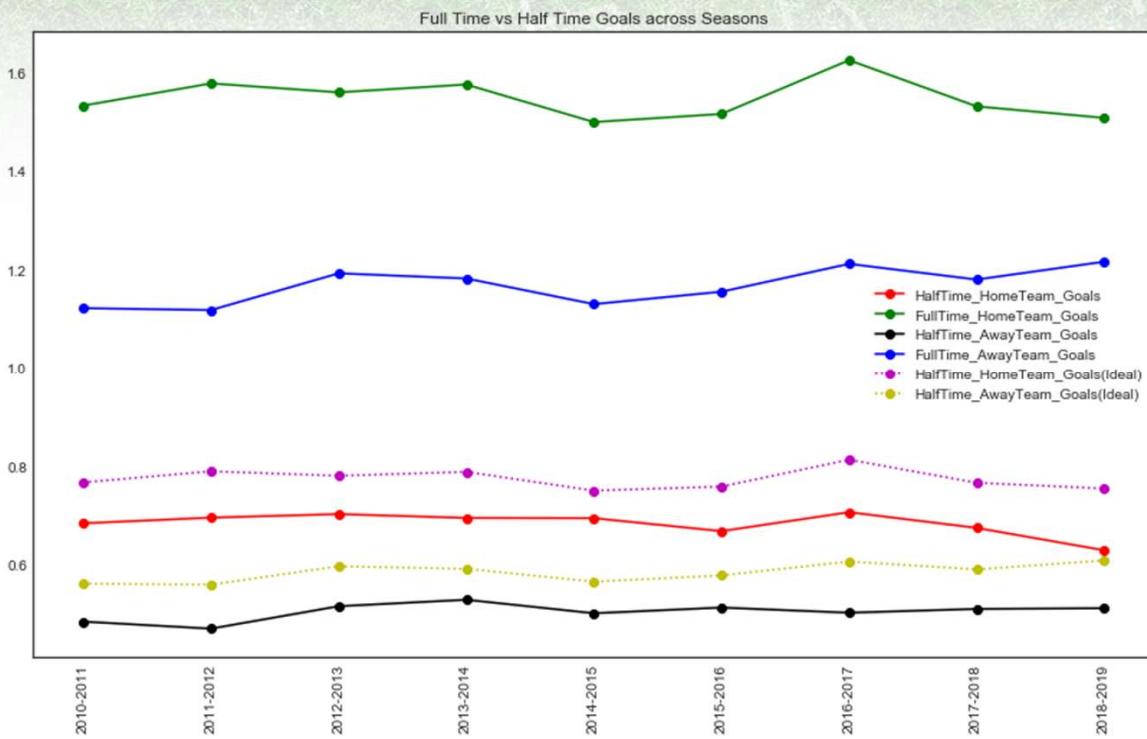
- ✓ Home team dominates from the beginning of the match
- ✓ Home advantage is clear here

Feature analysis – Full time vs Half time goals by League



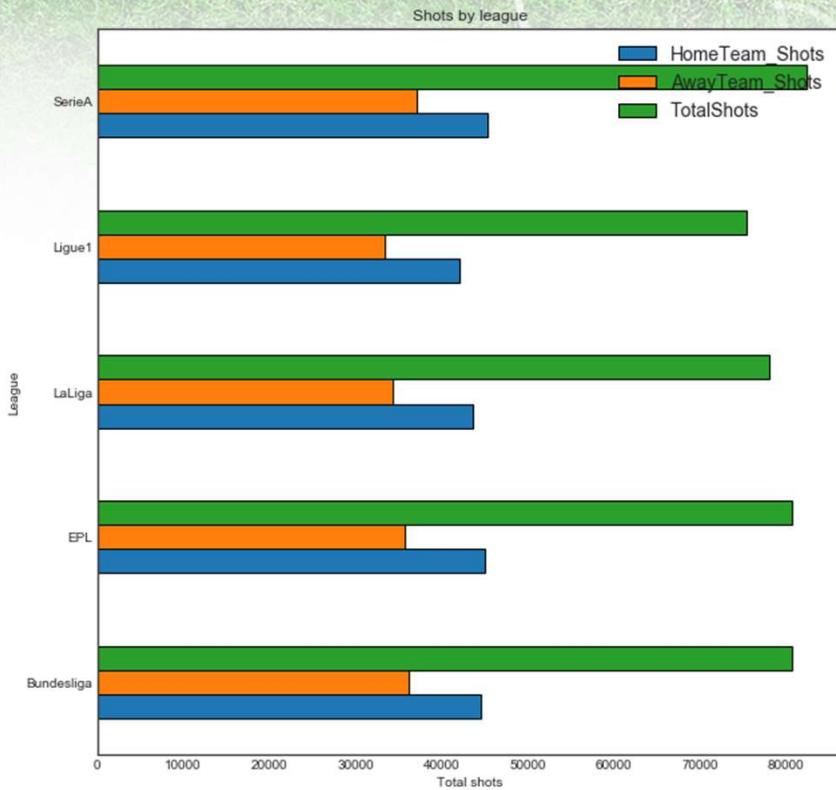
- ✓ For both Home teams and Away teams, half time goals is significantly less than full time goals.
- ✓ Both Home and Away teams aggressive during the second half of the match

Feature analysis – Full time vs Half time goals by Season



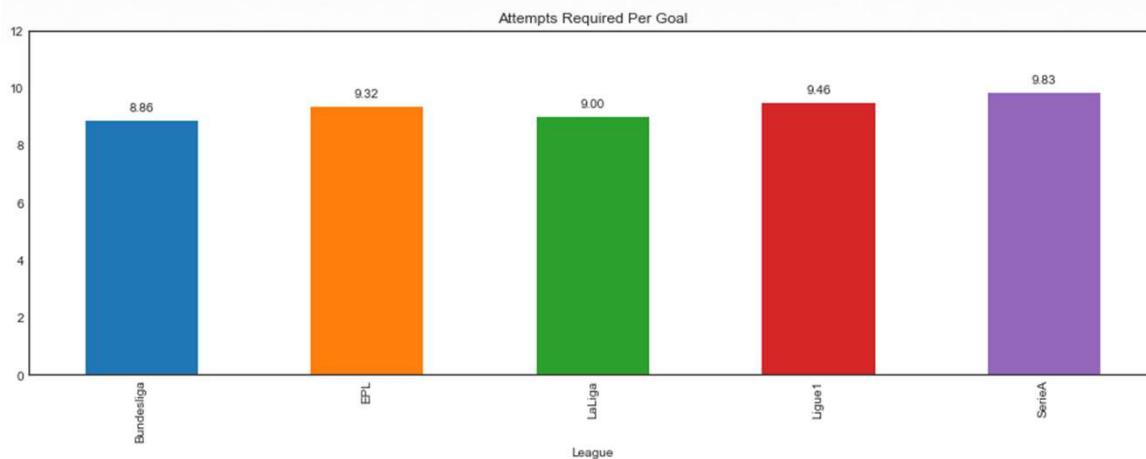
- ✓ Over the years, both home and away teams have remained defensive during the first half of a match and became aggressive during the second half.

Feature analysis – Shots across all seasons



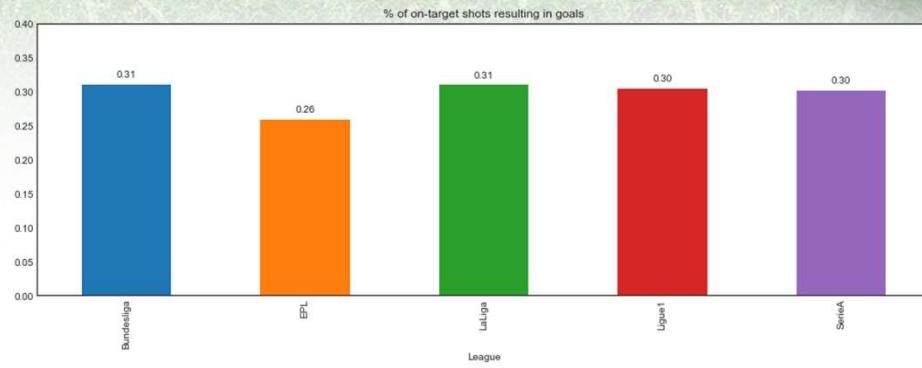
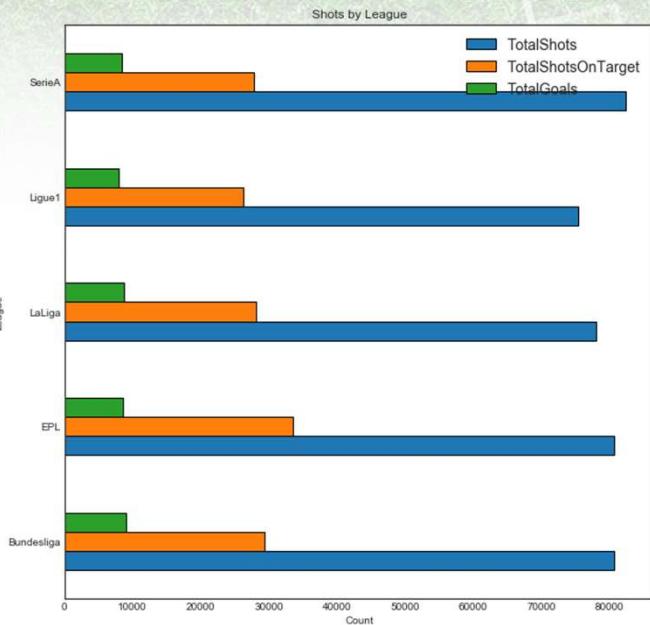
- ✓ Serie A, Bundesliga (adjusted) and EPL has max number of shots at goal per season.
- ✓ Interestingly, La Liga leads in number of goals per season but not at number of shots whereas it is other way round for Serie A
- ✓ So, we can conclude that La Liga forwards are more accurate in general as compared to Serie A forwards

Feature analysis – Attempts per goal



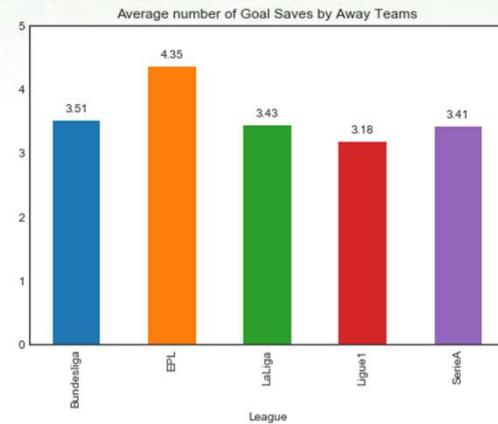
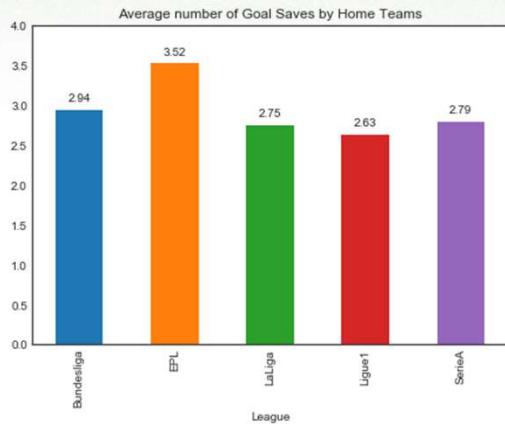
- ✓ Forwards in Bundesliga and LaLiga are more accurate as compared to others. Forwards in SerieA are the least effective

Feature analysis – Shots on target and yield



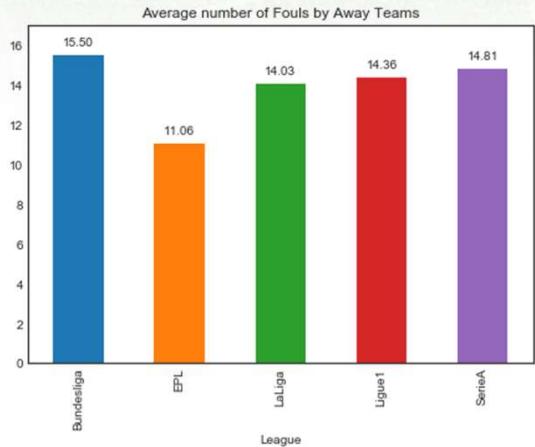
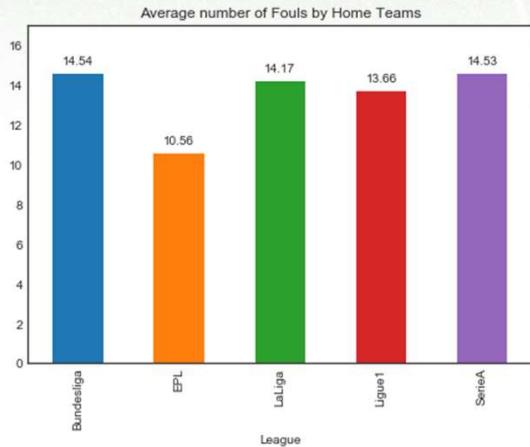
- ✓ Even though EPL has more number of shots on target, its yield is the least.

Feature analysis – Goals saved



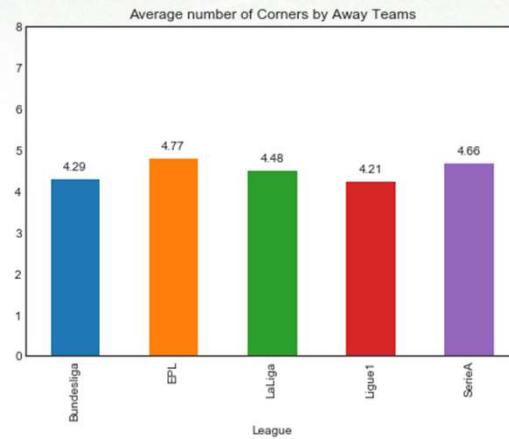
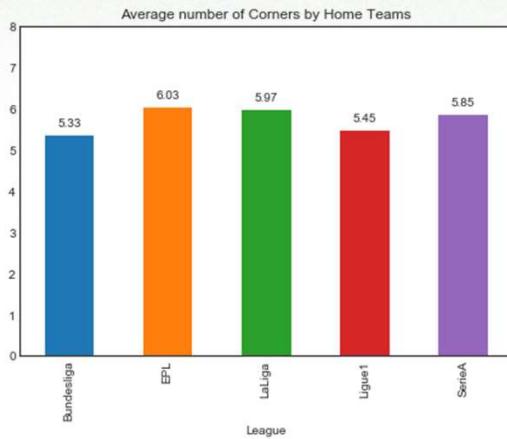
- ✓ Interestingly Away teams save more goals than home team in almost every league. But that might also be because home teams usually are more effective with number of shots on target. EPL probably keeps goalkeepers busier than other leagues

Feature analysis – Fouls



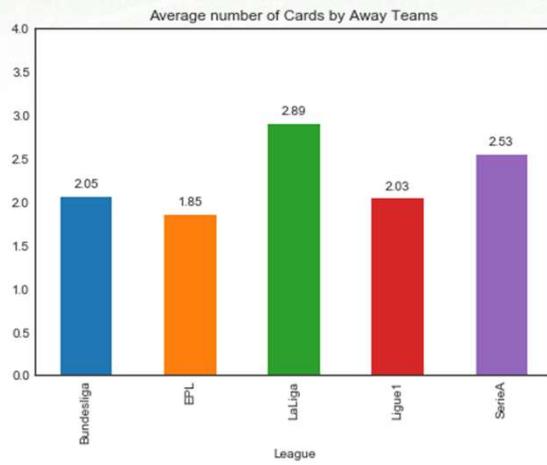
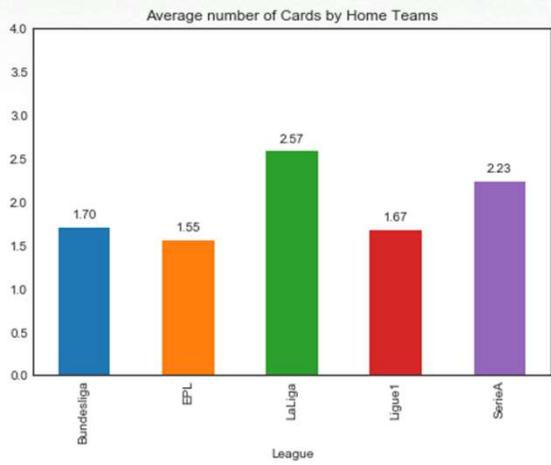
- ✓ EPL is apparently the most disciplined league.
- ✓ Bundesliga has highest number of fouls but that could be because it was scaled up.

Feature analysis – Corners



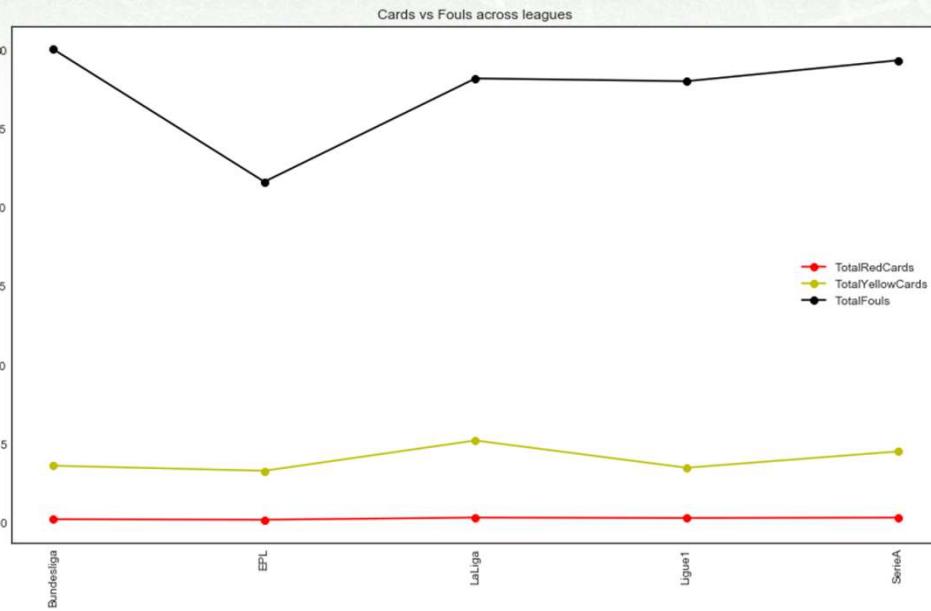
- ✓ Home team is awarded with most corners. This aligns with number of shots on target.
- ✓ This also indicates home team is more aggressive.

Feature analysis – Red/Yellow cards



✓ Away team faces more red/yellow cards showing that they become aggressive to catch up.

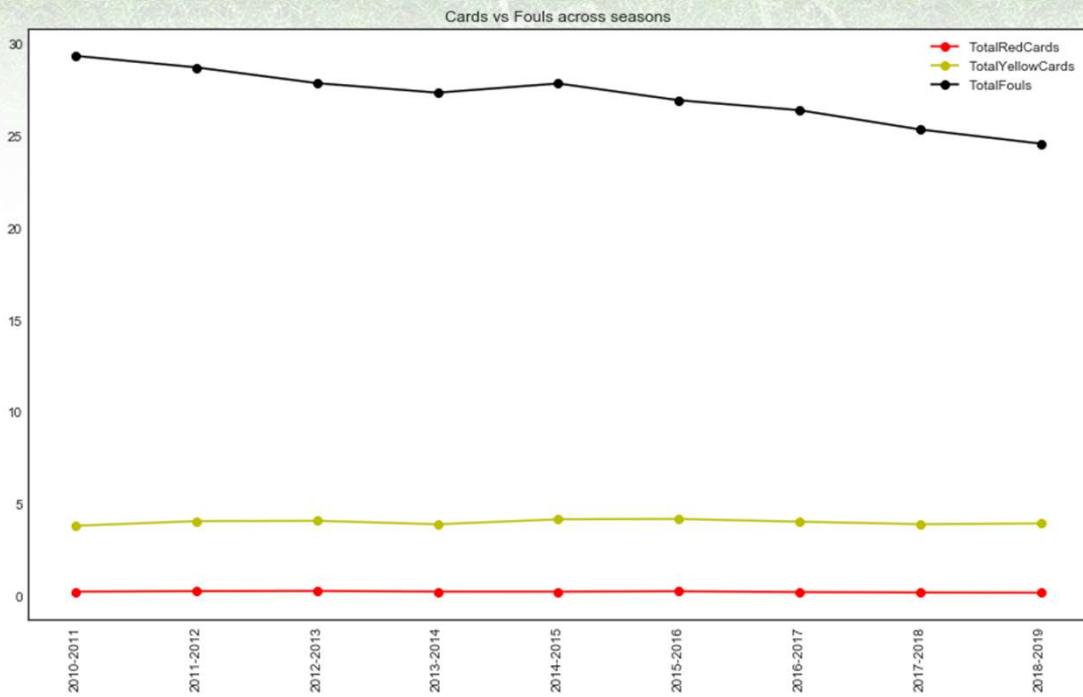
Feature analysis – Correlation between cards and fouls



- ✓ Except in EPL, there seems to be a linear correlation between number of cards and number of fouls



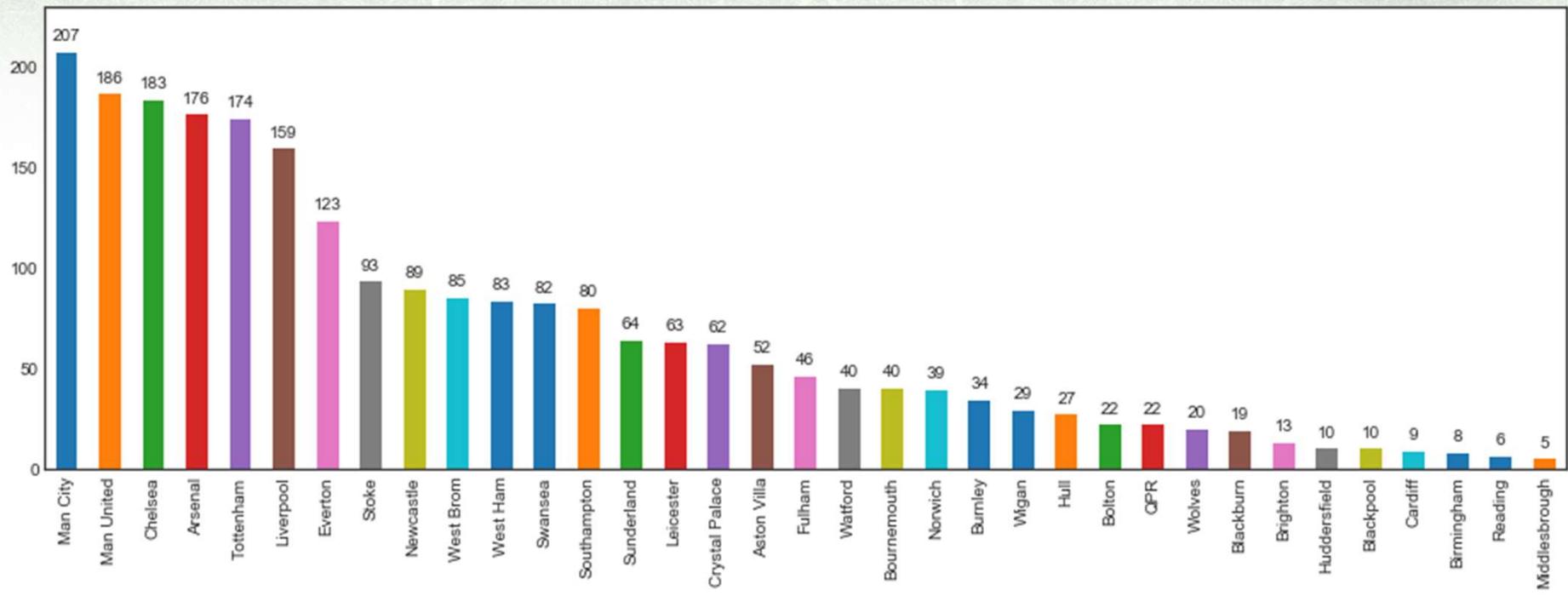
Feature analysis – Fouls and cards trend



- ✓ Average number of fouls per game has come down significantly however the average number of cards haven't changed much. This indicates that referees have become more strict over the years

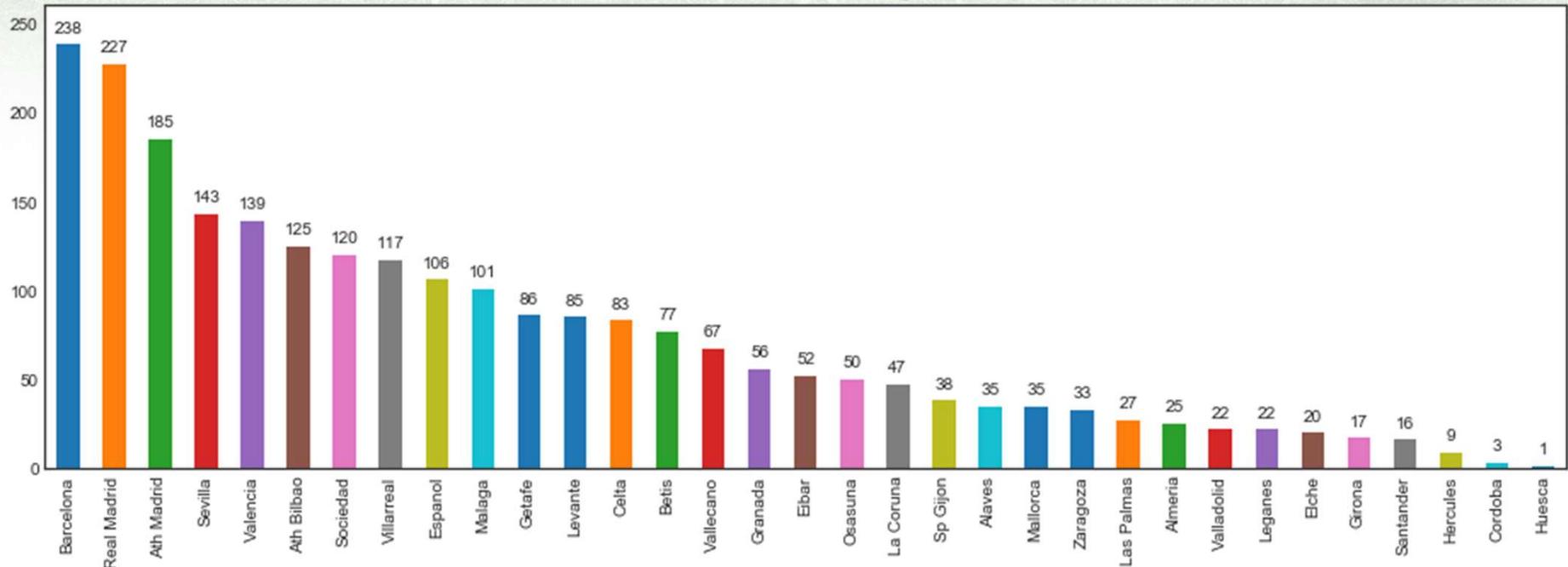
Feature analysis – Winning teams- EPL

Number of wins in EPL



Feature analysis – Winning teams- La Liga

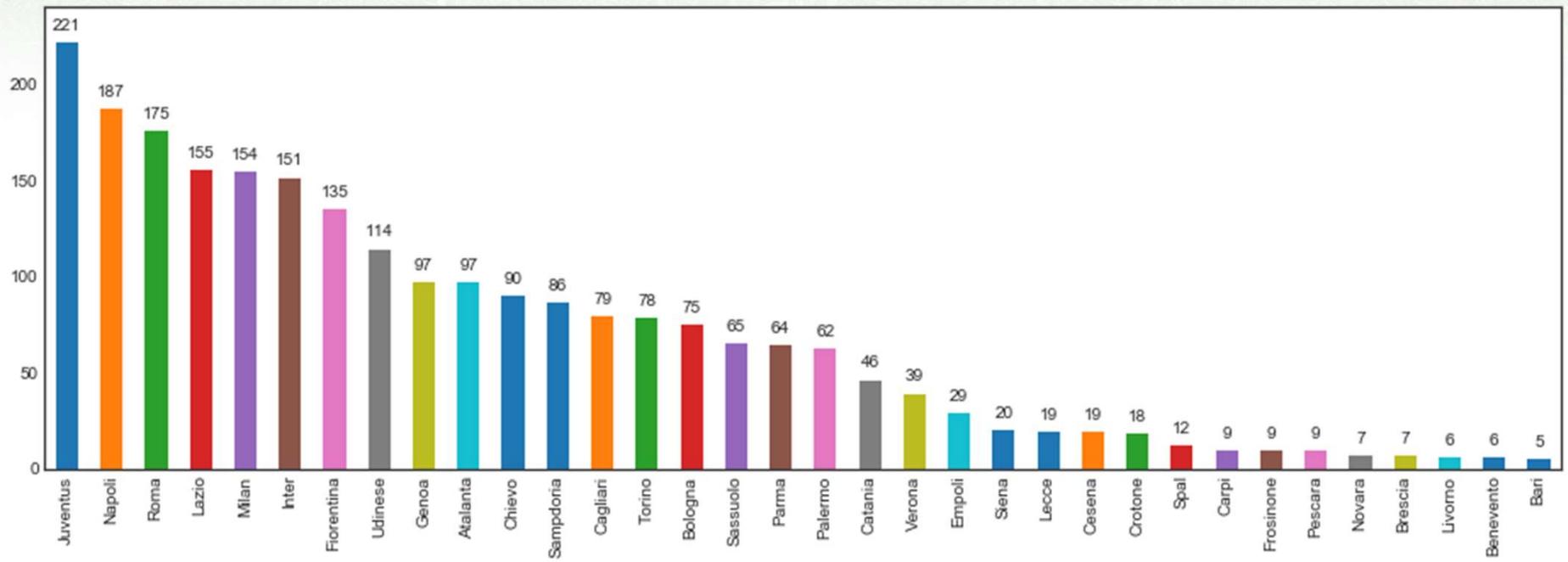
Number of wins in La Liga



Feature analysis – Winning teams- Serie A

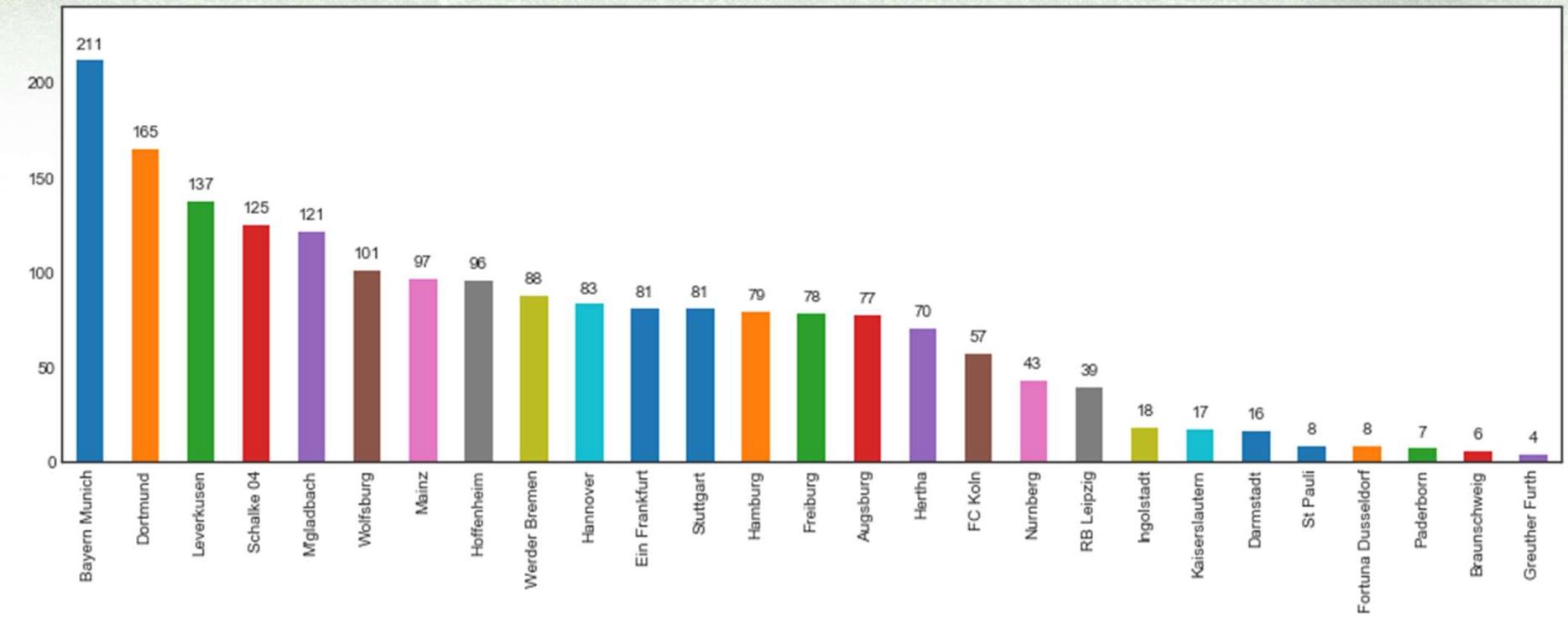


Number of wins in Serie A



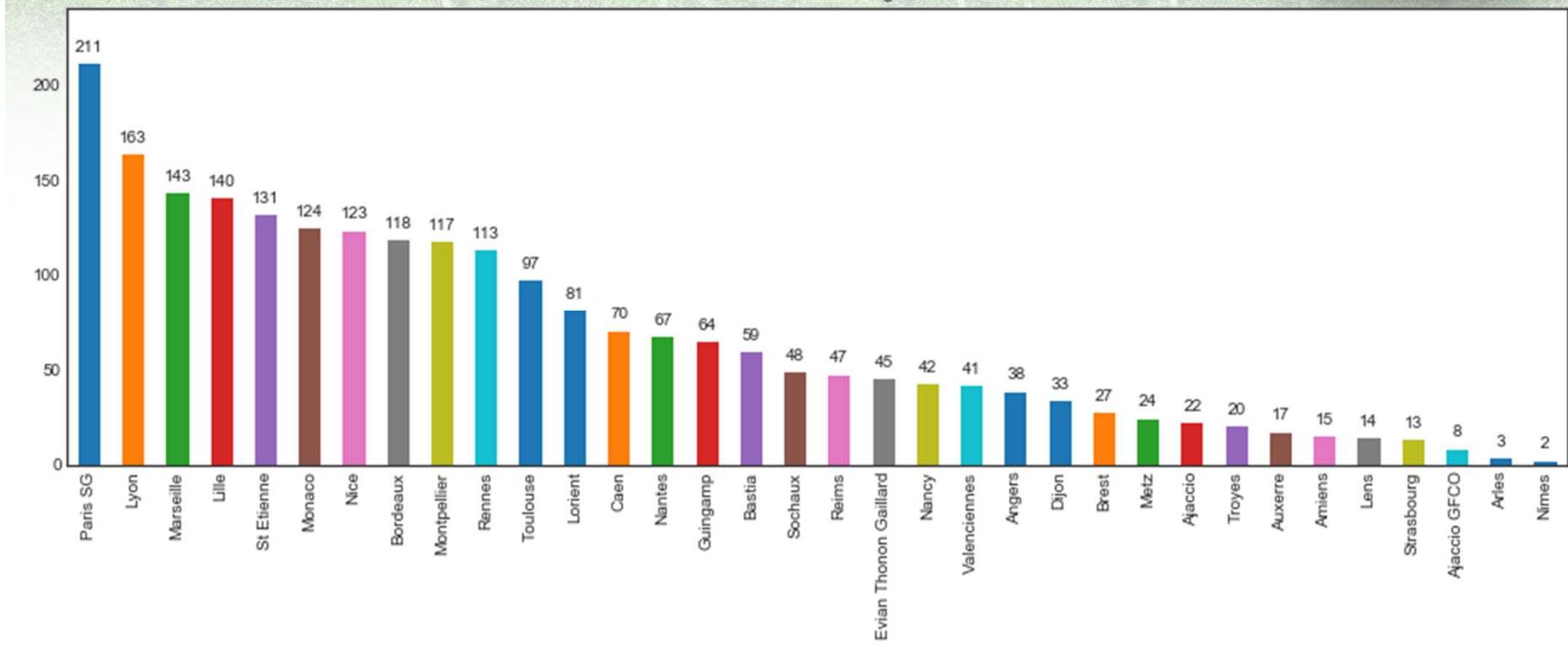
Feature analysis – Winning teams- Bundesliga

Number of wins in Bundesliga

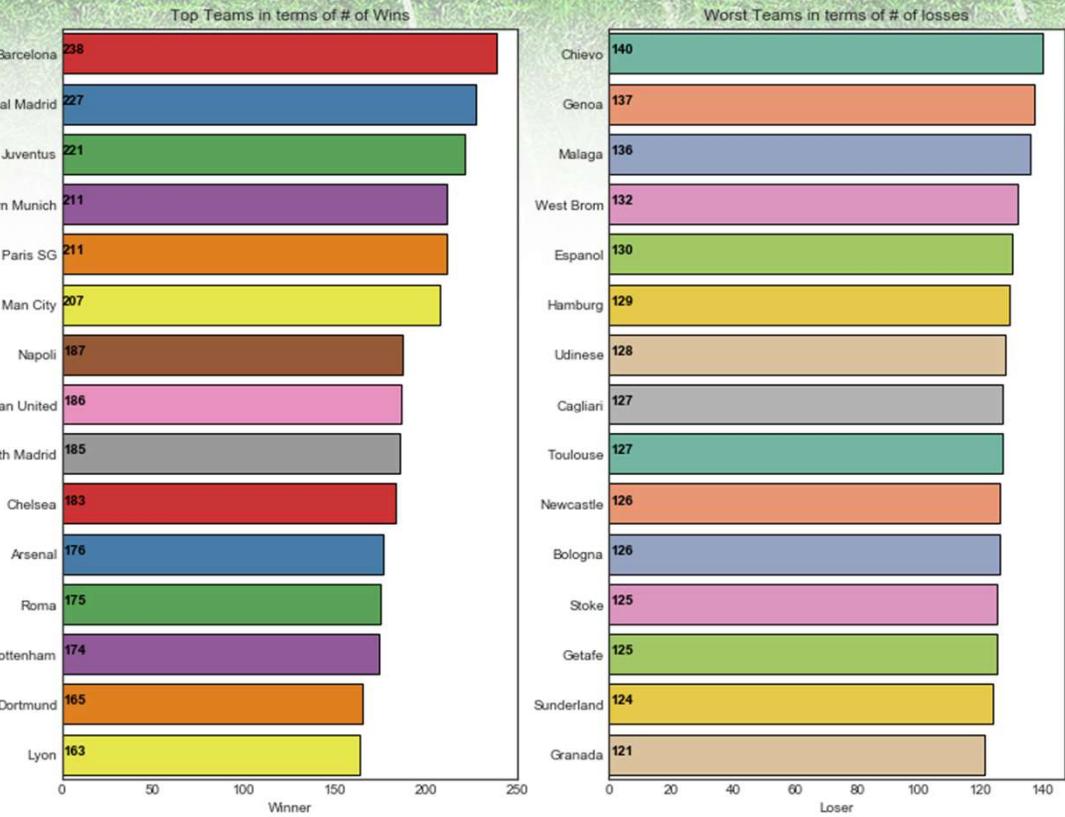


Feature analysis – Winning teams- Ligue1

Number of wins in Ligue 1



Feature analysis – Top and bottom 15



- ✓ Two La Liga giants- Barcelona and Real Madrid have won the maximum number of matches in last 9 seasons

Feature analysis – Conclusion



- EPL is the most volatile league in terms of league qualification. Only 35 % of teams have been consistently present in last 10 years
- There's a significant home advantage in every league. Away performance is improving in general but not at a very rapid pace
- EPL and LaLiga produces more goals compared to other leagues
- Both Home and Away teams produce more goals during second half
- Forwards in La Liga and Bundesliga are the most effective. Forwards in Serie A are the least effective
- EPL clearly is the most disciplined and least violent league with least number of fouls and cards. La Liga leads in number of red and yellow cards
- Average number of fouls have fallen across all leagues over last 10 years. However average number of cards haven't
- EPL forwards land maximum number of shots on target. However, EPL goalkeepers save maximum number of goals as well

Match Outcome Feature Analysis

- What features affect the outcome of a game?
 - Shots?
 - Fouls? Red Cards? Corners?
 - Does those feature change over seasons?
 - This analysis focused on the feature on winners
- What features make a team to won the Champion?
 - Home Goals?
 - Away Goals?
 - Cards and Fouls?

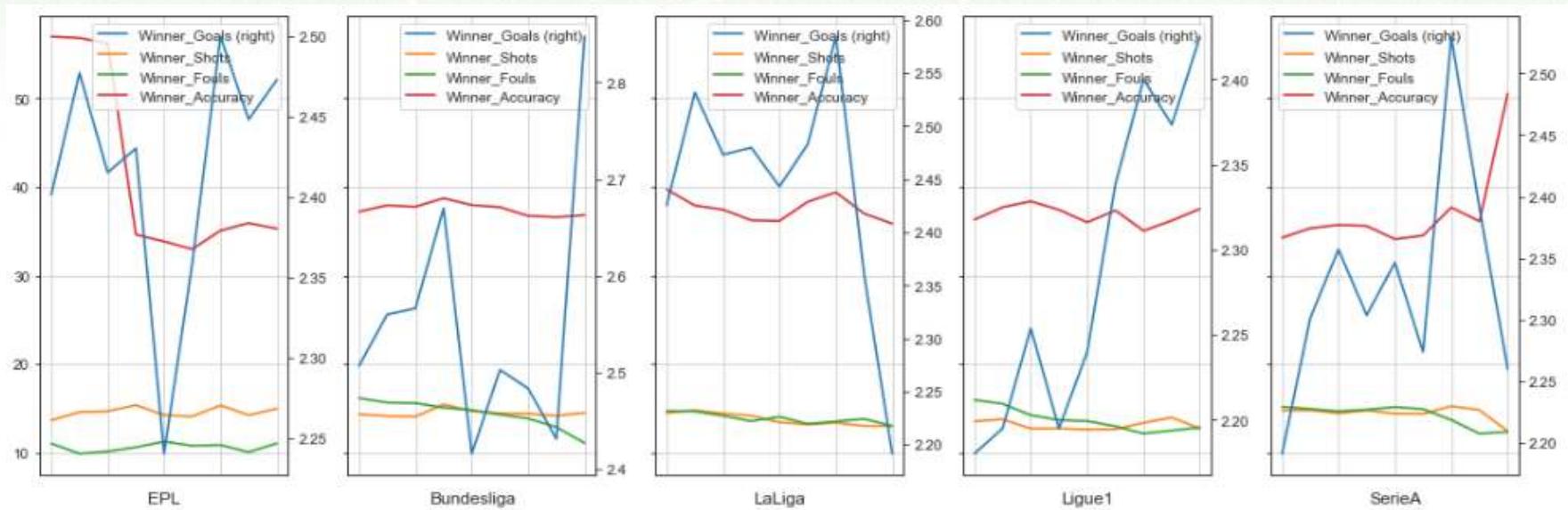
Match Outcome Feature Analysis

- ✓ The average winner goal is similar for different league
- ✓ The EPL winner has the highest accuracy, and least fouls.

League	Winner_Goals	Winner_Shots	Winner_Fouls	Winner_Accuracy
Bundesliga	2.582399	15.264884	14.191544	41.749881
EPL	2.454735	15.518106	10.390669	46.191782
LaLiga	2.554817	14.653821	13.990698	42.320281
Ligue1	2.319014	13.857746	13.553345	41.146722
SerieA	2.395528	15.464011	14.132075	40.289794

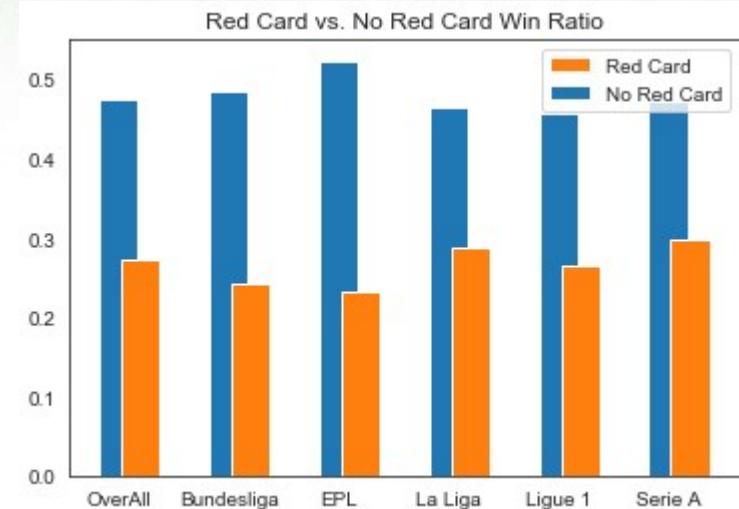


Match Outcome Feature Analysis



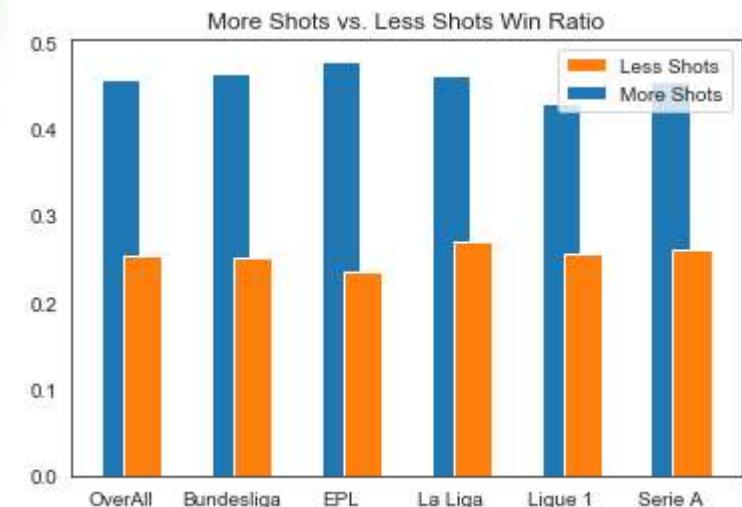
Match Outcome Feature Analysis

- ✓ Generally, the team without red card will have much higher chance to win compare to the team with a red card
- ✓ This is true for all 5 leagues, especially EPL. This may indicated EPL has higher competition.



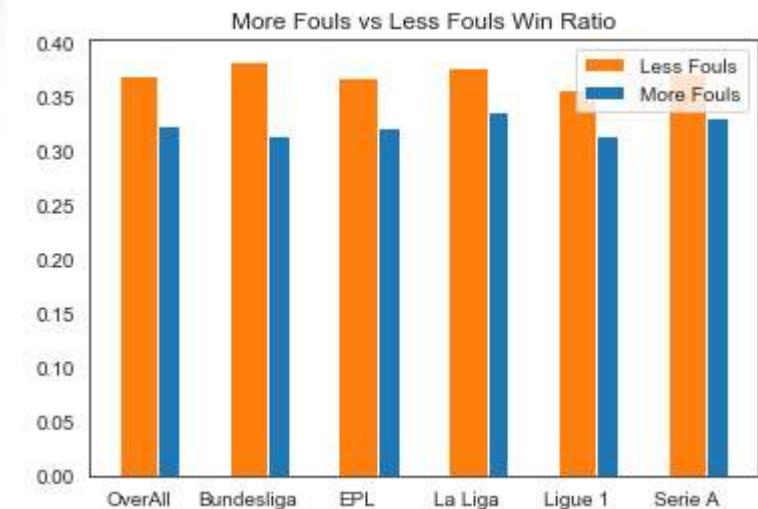
Match Outcome Feature Analysis

- ✓ Generally, the team with more shots will have higher chance to win the game
- ✓ This is true for all 5 leagues, especially EPL.



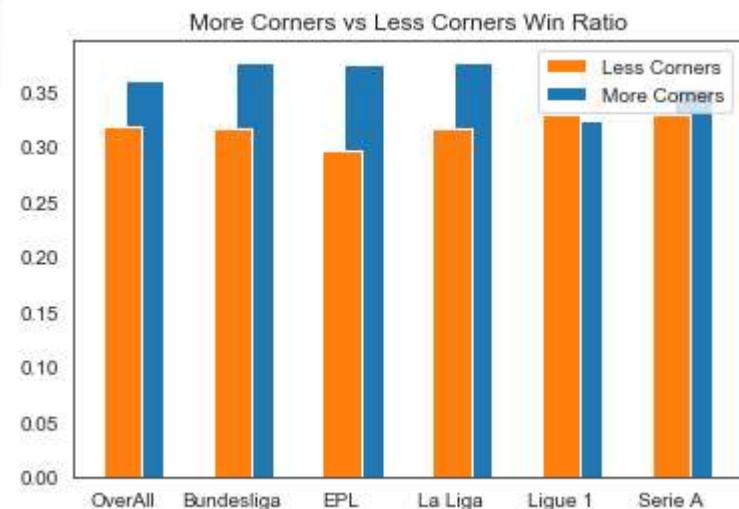
Match Outcome Feature Analysis

- ✓ The team with less fouls will have a high chance to win the game. However, the difference is small.



Match Outcome Feature Analysis

- ✓ The team with more corners will have a high chance to win the game. However, the difference is small.
- ✓ This is not true for Ligue 1.



Match Outcome Feature Analysis Summary

✓ Team will have higher chance to win the game if they have following features:

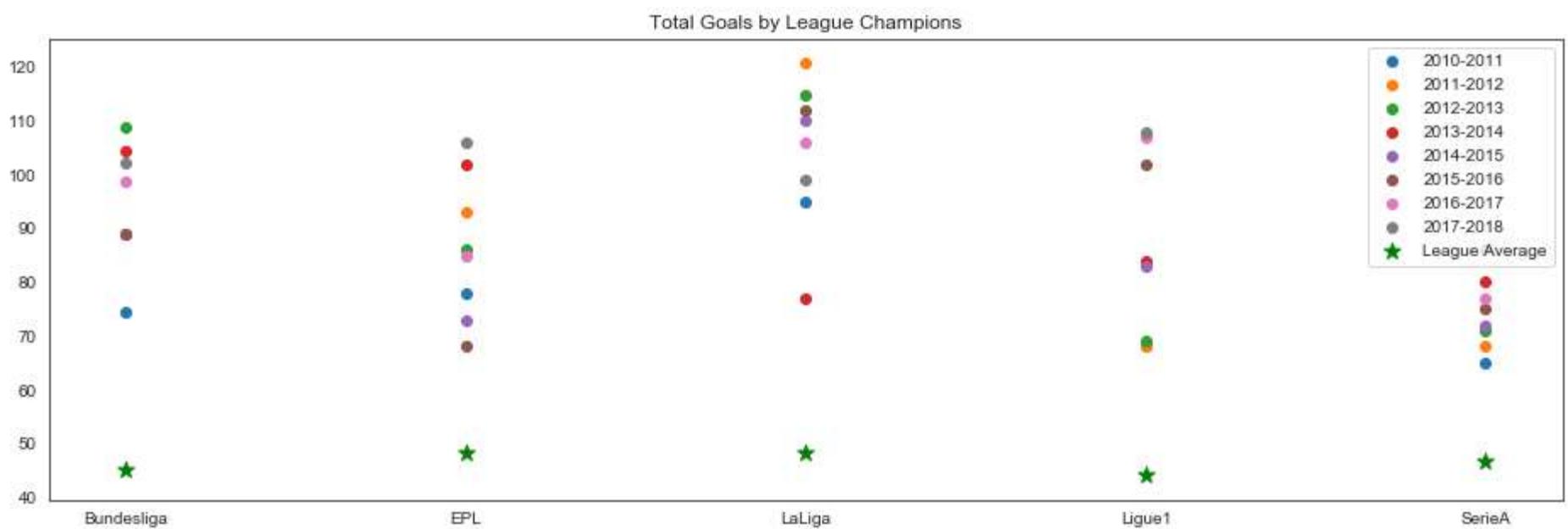
- ✓ more shots
- ✓ no red card
- ✓ less fouls
- ✓ more corners

Champion Feature Analysis

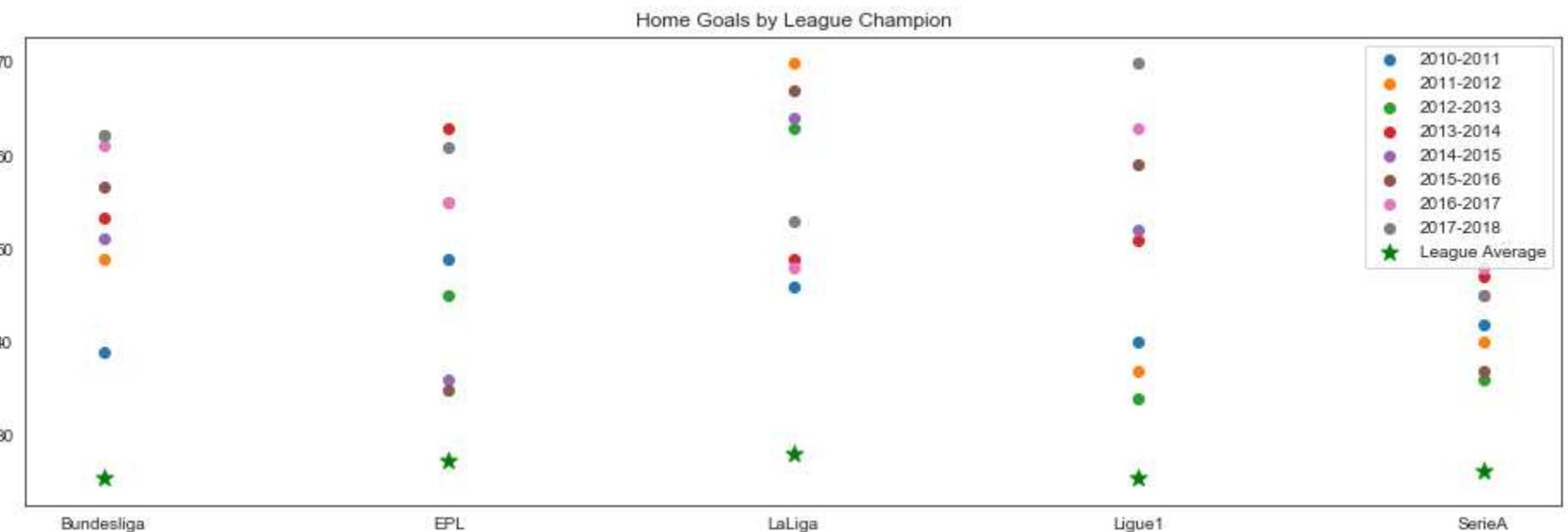
- ✓ Determine the Champion for each league for each season
 - ✓ Rules:
 - ✓ Win : 3 points
 - ✓ Draw : 1 point
 - ✓ Lose : 0 point
 - ✓ The team with highest accumulative points win the season
 - ✓ Juventus dominates Serie A wins 7 times since 2010.
 - ✓ Bayern Munich dominates Bundesliga wins 6 times.
 - ✓ Paris SG and Barcelona 5 times.
 - ✓ Only EPL doesn't have a team dominates the league.

Team	Count
Juventus	7
Bayern Munich	6
Paris SG	5
Barcelona	5
Man City	3
Real Madrid	2
Man United	2
Dortmund	2
Chelsea	2
Montpellier	1
Monaco	1
Milan	1
Lille	1
Leicester	1
Ath Madrid	1

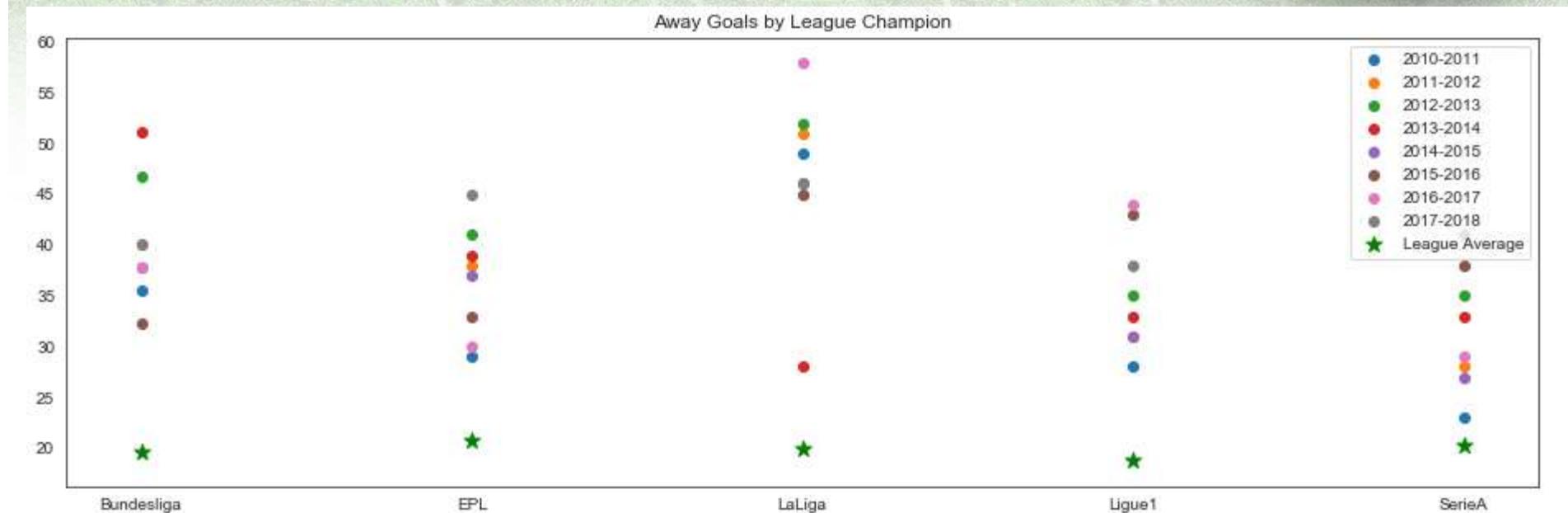
Match Outcome Feature Analysis



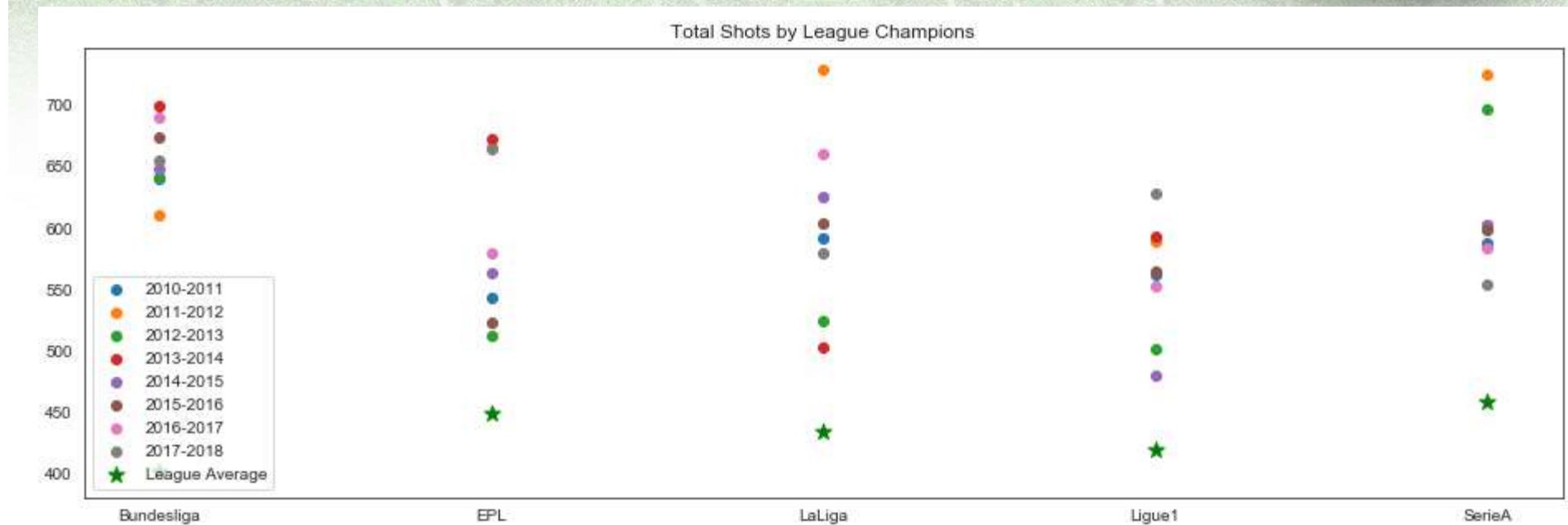
Match Outcome Feature Analysis



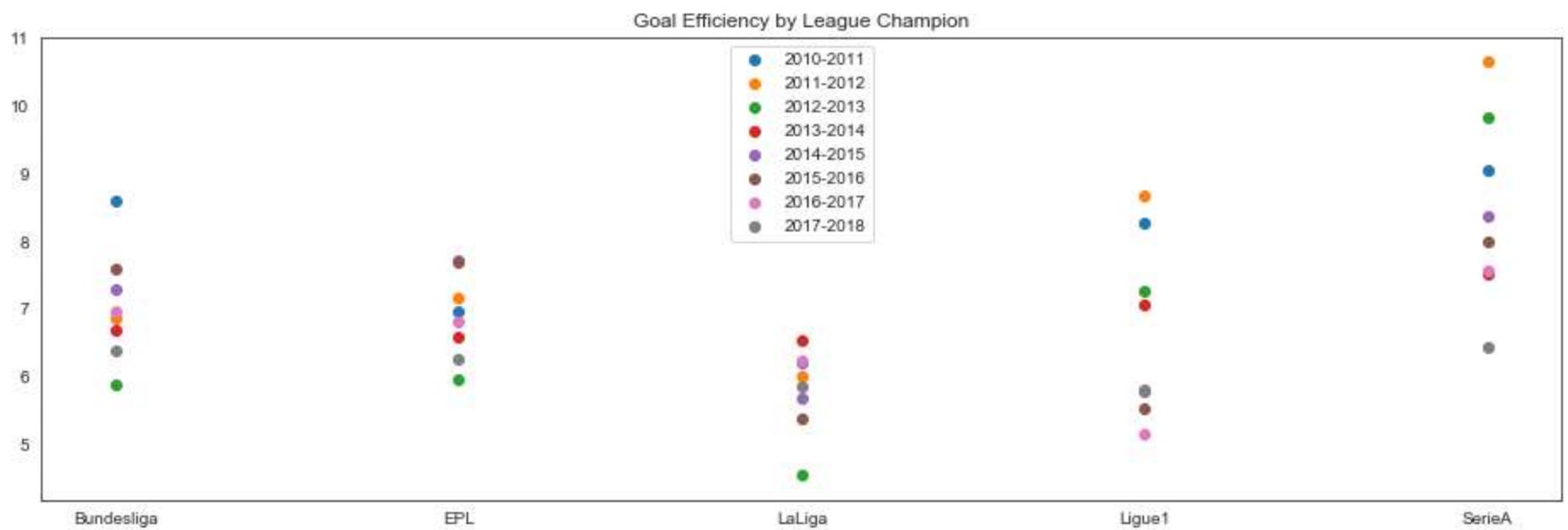
Match Outcome Feature Analysis



Match Outcome Feature Analysis

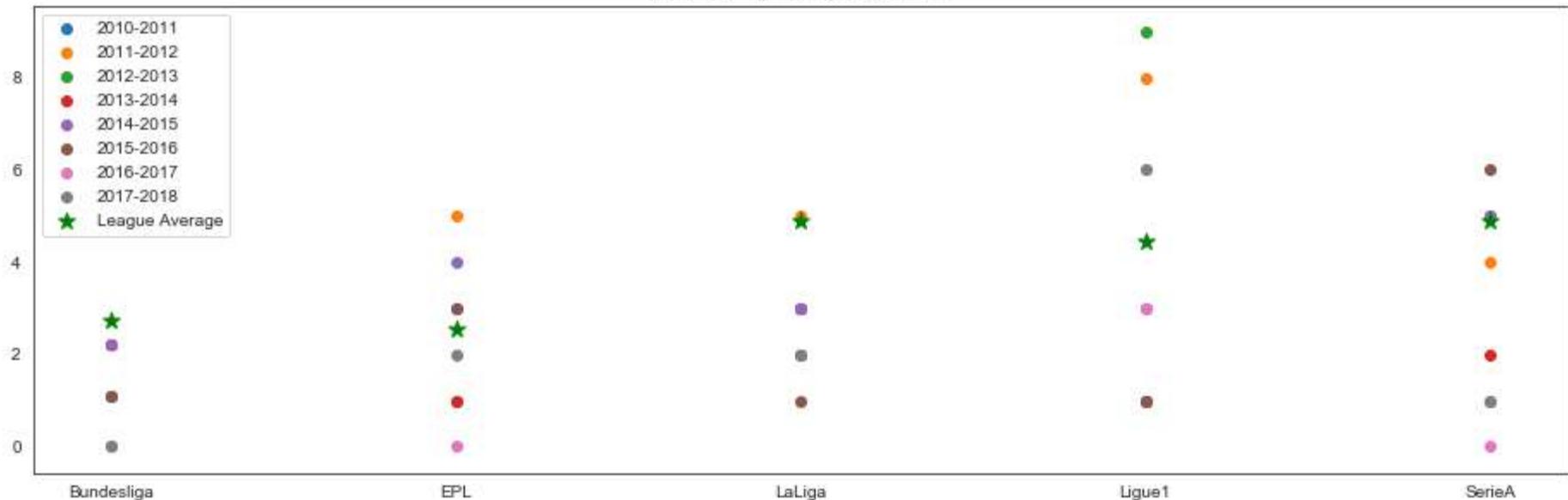


Match Outcome Feature Analysis

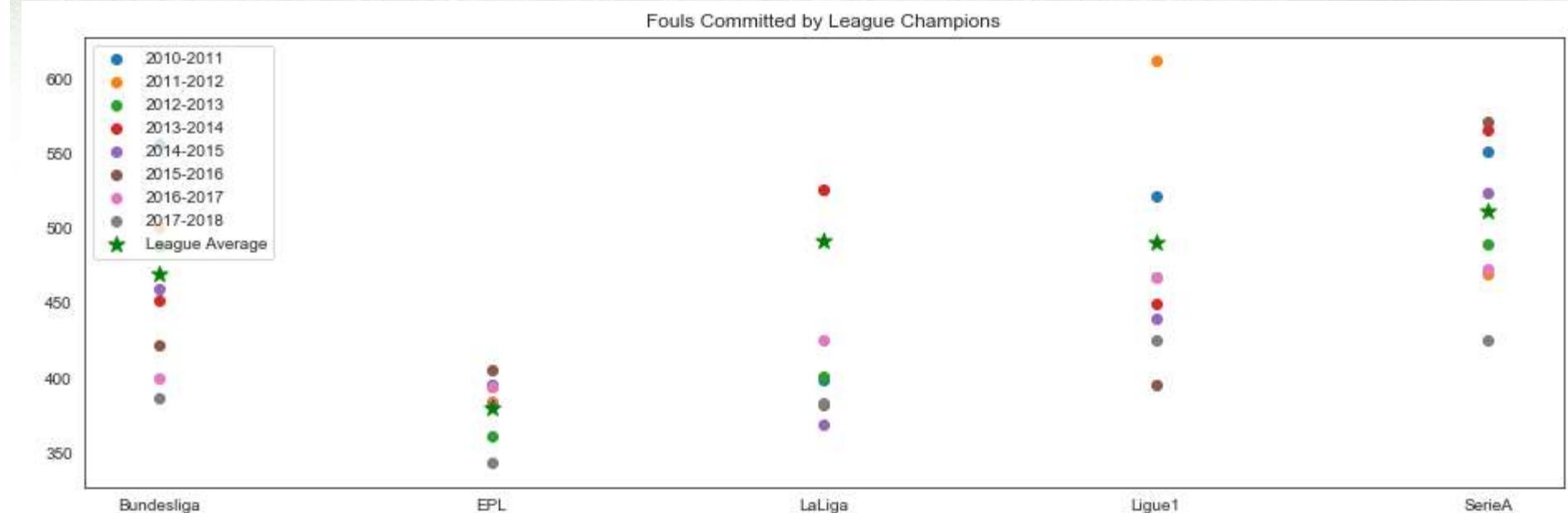


Match Outcome Feature Analysis

Red Card by League Champions



Match Outcome Feature Analysis

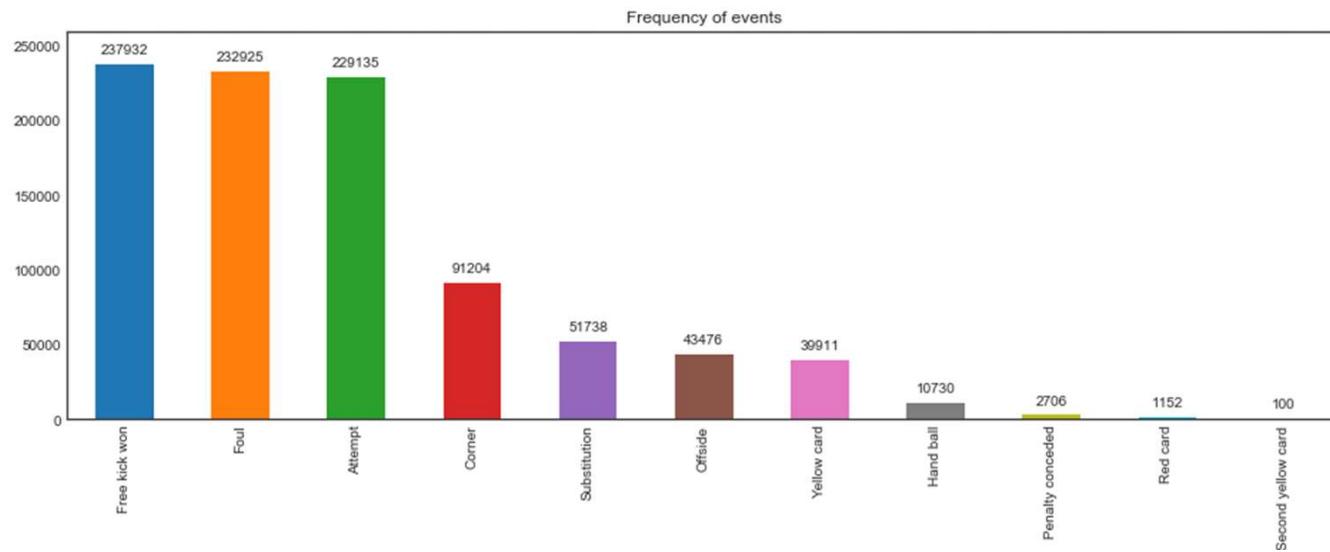


Match Outcome Feature Analysis Summary

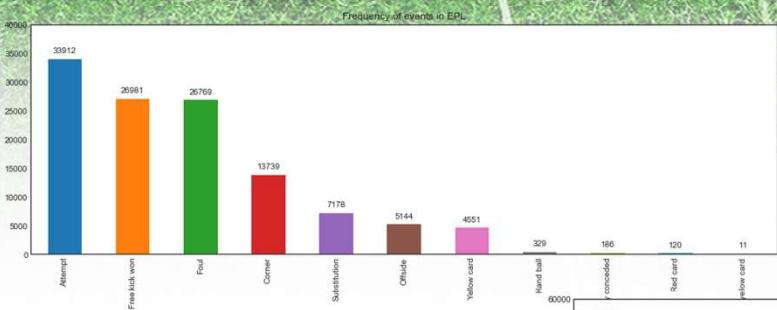
- ✓ Champions will have significantly more goals comparing the league average.
- ✓ Champion will have more goals at home games comparing to away games.
- ✓ Number of fouls and red card is similar to league champion and other teams.
- ✓ Comparison across leagues:
 - ✓ Serie A has the lease goal for champions and the lowest goal efficiency. Italian team tend to have better defence, this could be one of the reason.
 - ✓ La Liga has the largest gap between champion and league average on goals. La Liga also has the biggest dominance in the home game goal.
 - ✓ EPL's champion has the least fouls committed.

Event Analysis

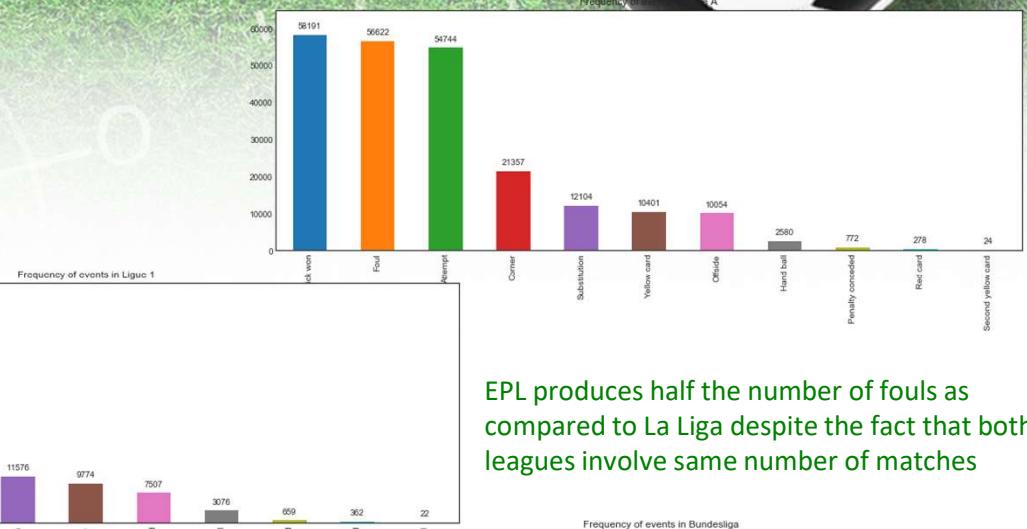
- We are using a secondary dataset to analyze events, individual performances and relationship between betting odds and actual match results
- Each row of this dataset corresponds to an event during a soccer match – overall the dataset has information for almost a million events
- Some examples of events are : Foul, Goal Scored, Goal Blocked, Substitution, Offside, Handball, Penalty etc.



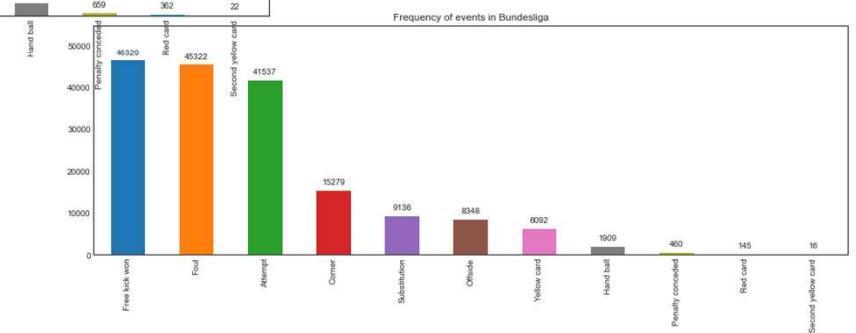
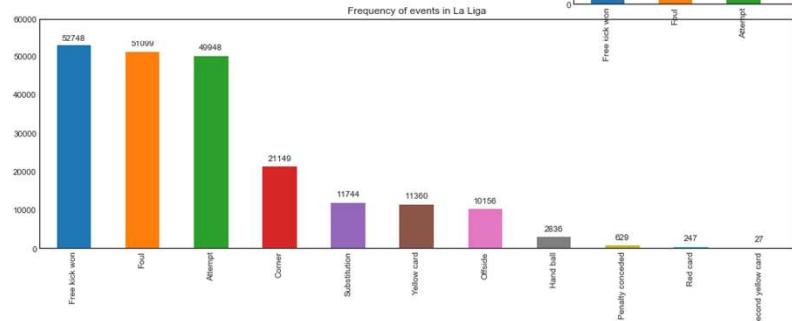
Event Analysis



So the top 3 events in every league are Free kick won, Fouls and Attempts

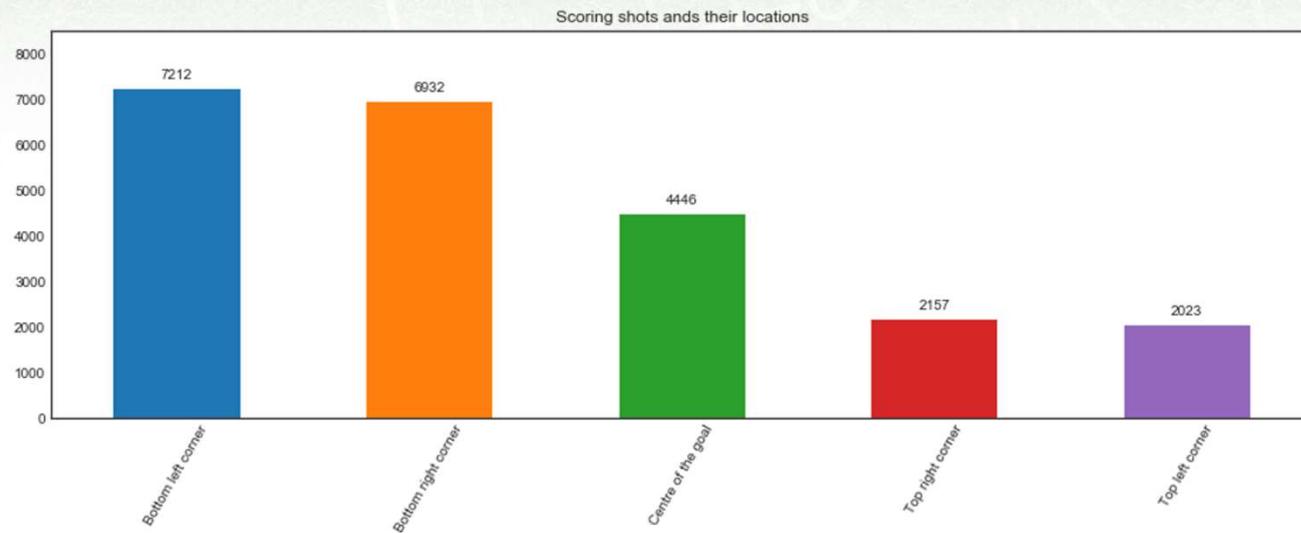


EPL produces half the number of fouls as compared to La Liga despite the fact that both leagues involve same number of matches



Event Analysis: Where to aim in order to score a goal ?

Here's what goal data for last 6 years reveal :



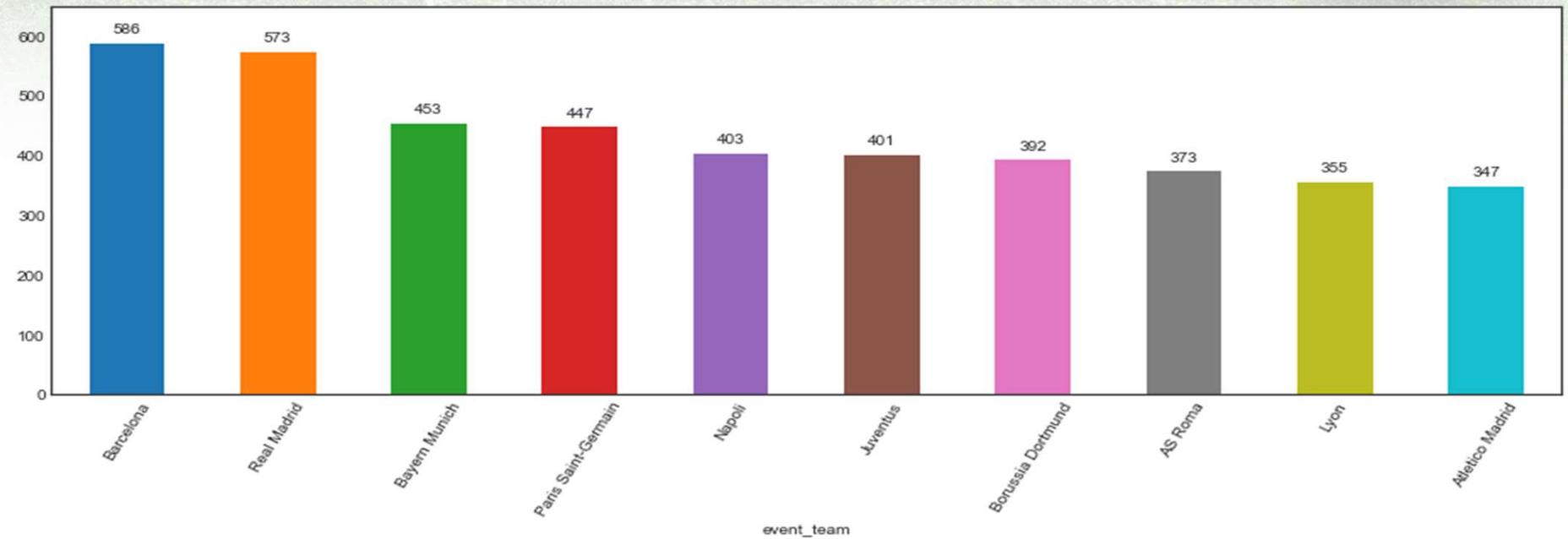
Data clearly suggests that most of the goals come off shots aimed at the bottom left or right corner

Which teams and individuals are the most prolific scorers ?



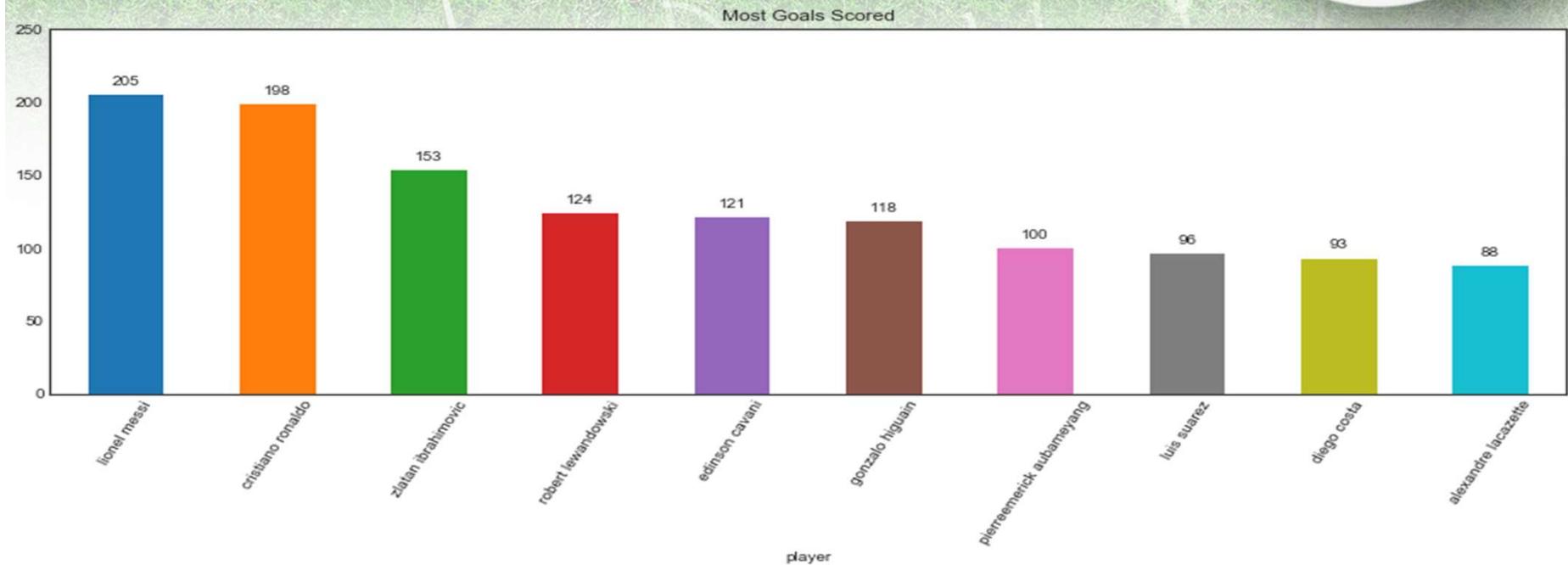
Lets look at the team level data first

Most Goals Scored



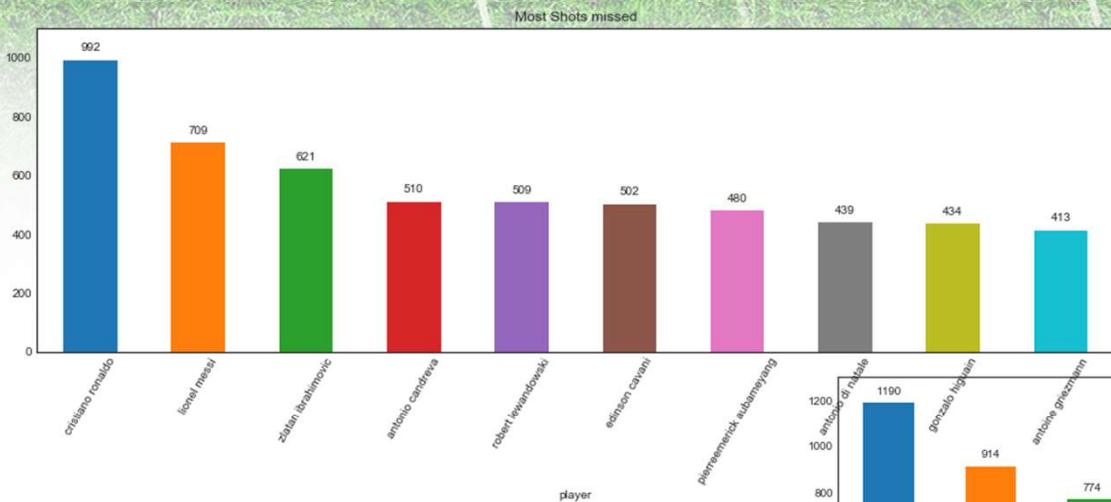
The La Liga giants – Barcelona and Real Madrid are at the top. Bayern Munich is at no 3 with 453 and even if we apply adjustments for La Liga, falls way short of top 2 at 506. Interestingly none of the EPL teams made it to the list of top 10 goal scorers

Which teams and individuals are the most prolific scorers ?



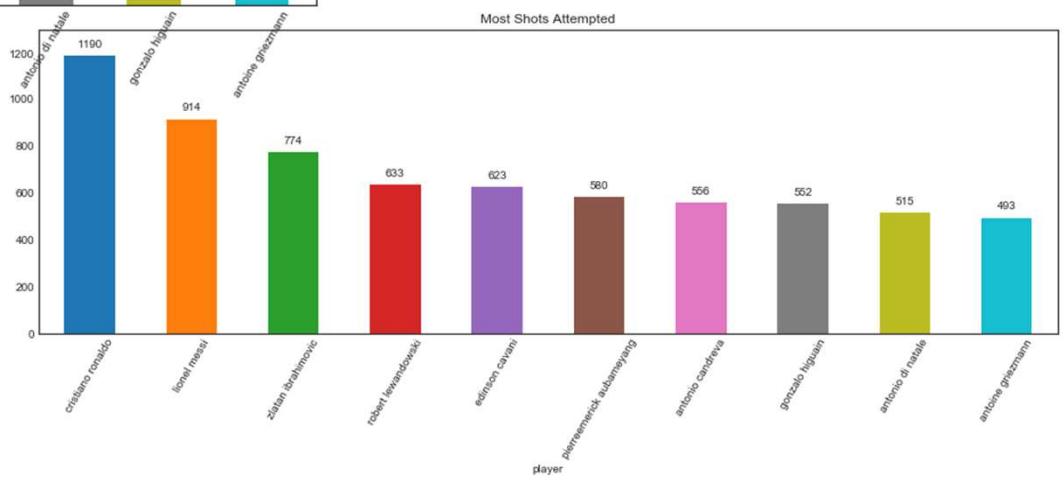
No surprises there at top 2. Zlatan Ibrahimović is the only player to make it to this list after playing in three different leagues (Serie A, Ligue One, EPL) during this period. Also since he plays for LA Galaxy now, his 2018-19 season data isn't included here. So he is much closer to no 2 than what the data suggests

But what about the number of goal attempts missed ?



Again the usual suspects at top 3 – clearly they create the most number of chances and hence score the most as well as miss the most

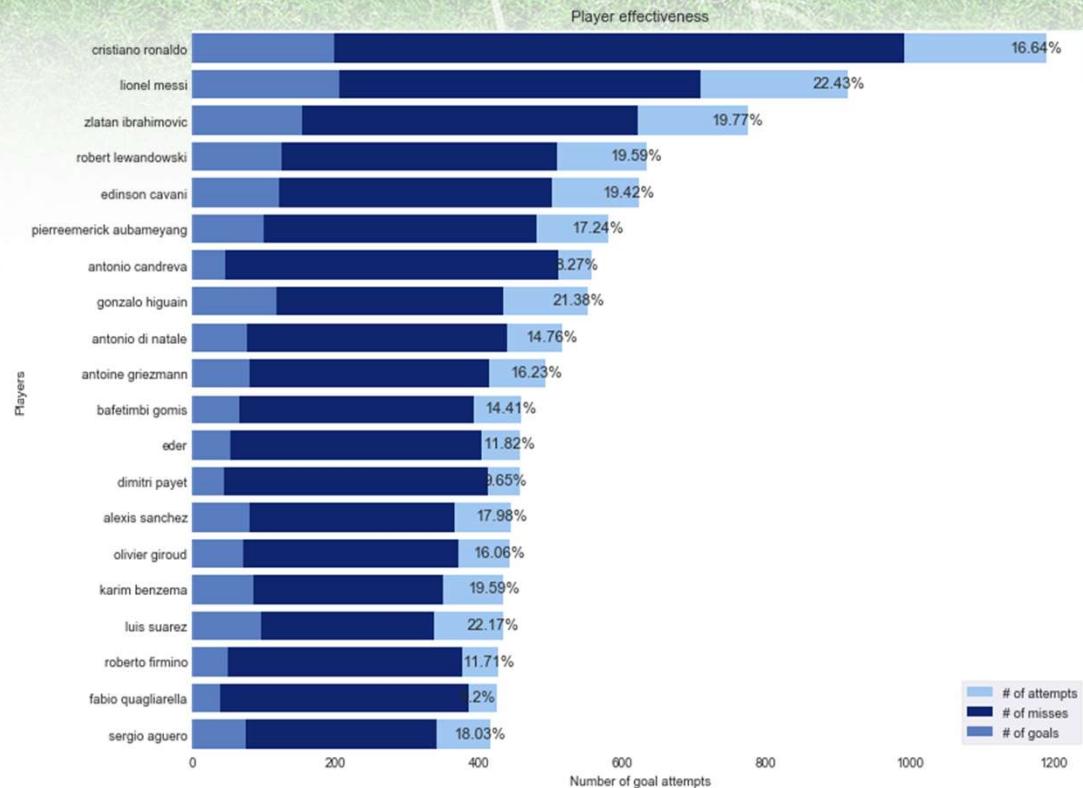
Cristiano Ronaldo has almost a 20 % lead over Lionel Messi in terms of number of terms



The top 3 in terms of number of attempts missed are also the same trio of Ronaldo, Messi and Zlatan



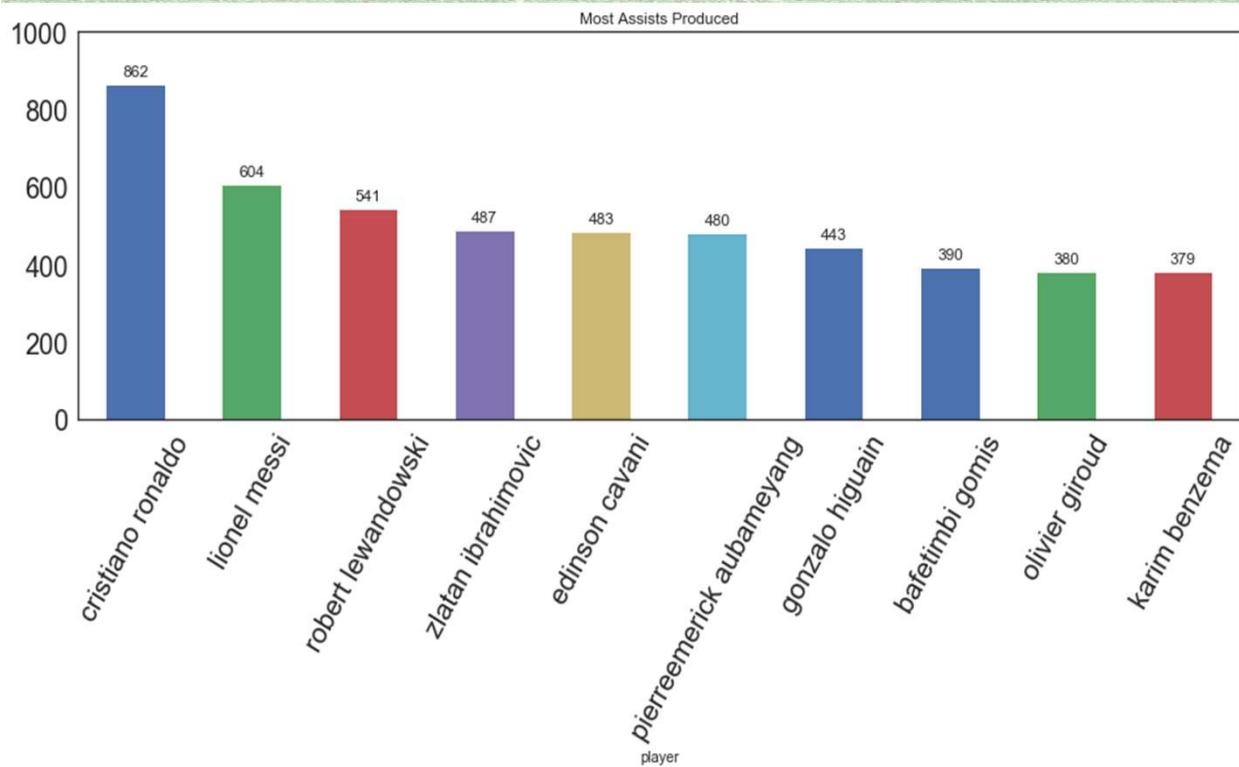
What about efficiency at scoring goals ?



Cristiano Ronaldo and Lionel Messi are considered among greatest of all times because they have been consistently creating chances.

Their teammates like Suarez and Higuain have better attempt-to-goal-conversion-rate but doesn't create half as many chances as Messi and Ronaldo

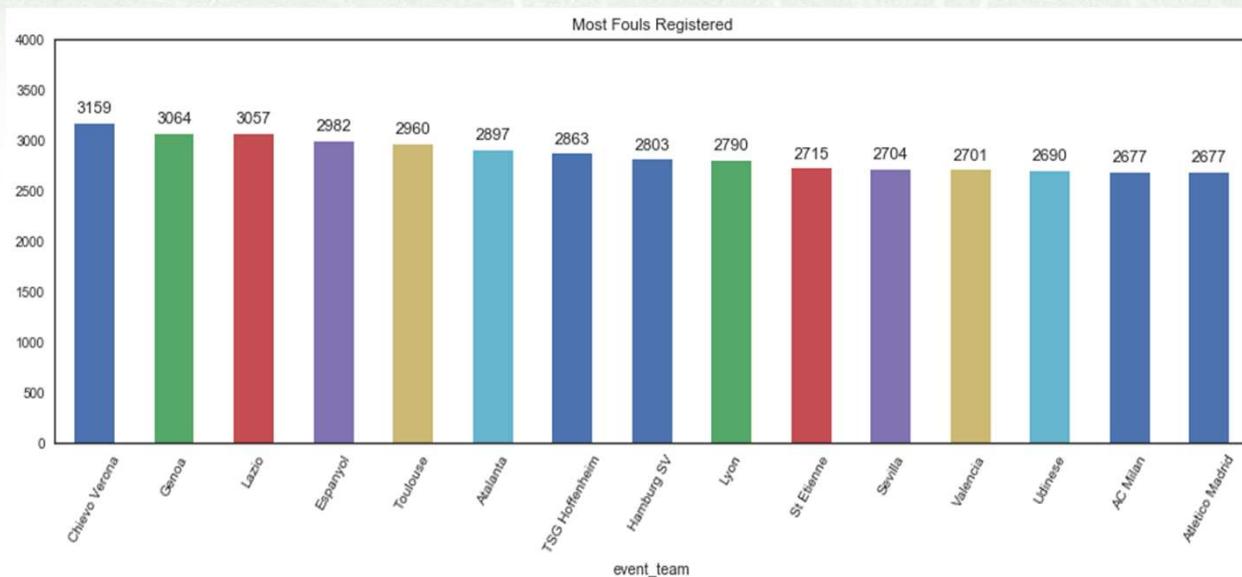
Which individuals produce the most number of goal assists ?



- Cristiano Ronaldo tops the number of assists as well (and is ahead of Messi at no 2 by nearly 30 %)
- Two of Ronaldo's teammates make it to the list as well – Benzema and Higuain. None of Messi's Barcelona teammates make it to this list
- Ronaldo, Messi and Zlatan are clearly the top 3 forwards from the last 6 years

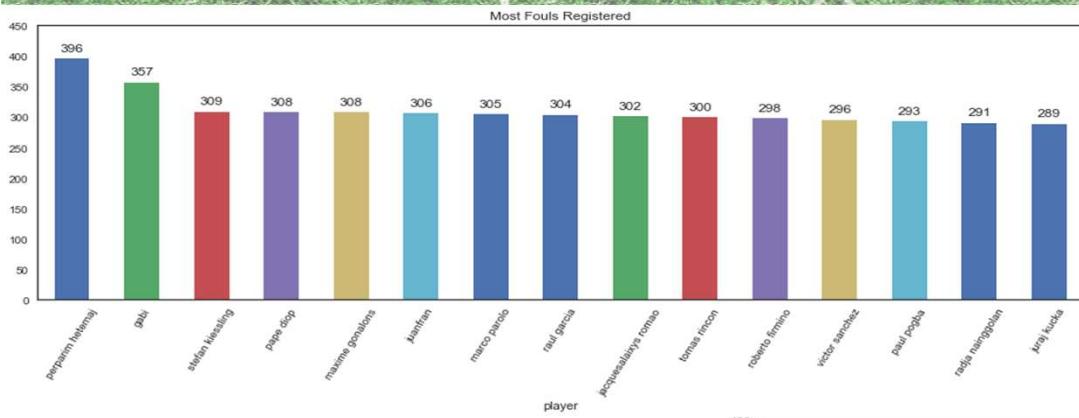


Defense yardsticks – number of fouls and cards



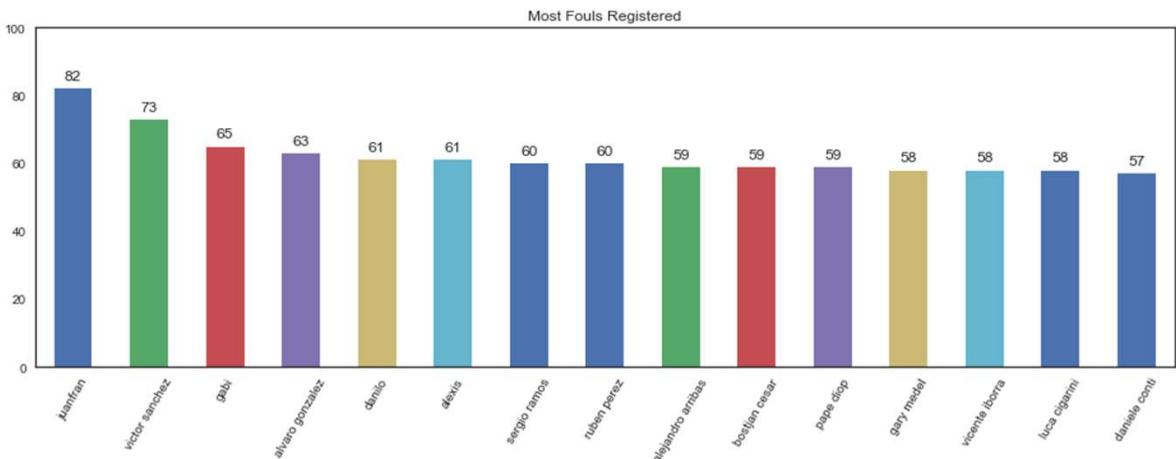
- More fouls are committed by lower ranked teams in the league.
- Lyon and Lazio are the only top tier teams with very high number of fouls

Defense yardsticks – number of fouls and cards



Surprising to see an attacking midfielder like Paul Pogba in this list. Equally surprising to see not see Sergio Ramos in this list

- Not many players made it to both the lists.
- Gabi from Athletico Madrid is among the worst offenders - he features in top 3 in terms of both the number of fouls committed as well as the number of cards



Building a player comparator function



```
1 def player_comparator(events, players):
2
3     data = DataFrame(columns=['Goals','Misses','Attempts','Assists','Cards','Fouls','Offsides'],index = players)
4
5     offsides_grouping_by_player = events[events['event_type']=='Offside'].groupby('player').count()
6     attempts_grouping_by_player = attempts.groupby('player').count()
7     misses_grouping_by_player = non_goals.groupby('player').count()
8     goals_grouping_by_player = goals.groupby('player').count()
9     assists_grouping_by_player = assists.groupby('player').count()
10    cards_grouping_by_player = cards.groupby('player').count()
11    fouls_grouping_by_player = fouls.groupby('player').count()
12
13    build_feature(data, goals_grouping_by_player, 'Goals', players )
14    build_feature(data, misses_grouping_by_player, 'Misses', players )
15    build_feature(data, attempts_grouping_by_player, 'Attempts', players )
16    build_feature(data, assists_grouping_by_player, 'Assists', players )
17    build_feature(data, cards_grouping_by_player, 'Cards', players )
18    build_feature(data, fouls_grouping_by_player, 'Fouls', players )
19    build_feature(data, offsides_grouping_by_player, 'Offsides', players )
20
21    return data
22
23
24 def build_feature(target, source, feature, players):
25     for p in players:
26         if (p in target.index):
27             row = target.loc[p]
28             row[feature]=source.loc[p]['id_odsp']
29
```

Building a player comparator – sample input and output



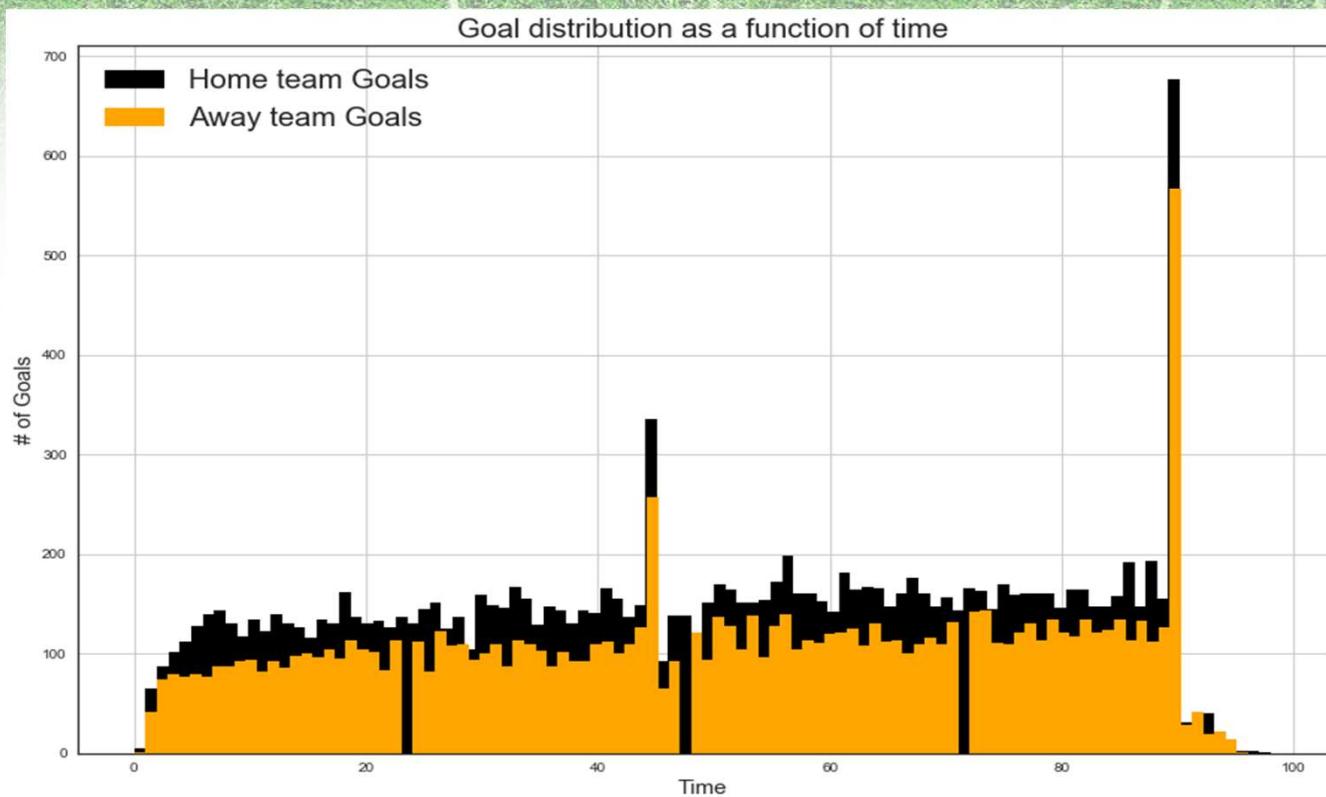
Sample Usage :

```
cmp = player_comparator(events, ['lionel messi','cristiano ronaldo','zlatan ibrahimovic','gonzalo  
higuain', 'luis suarez','karim benzema','gareth bale', 'harry kane','neymar'])  
display(HTML(cmp.to_html()))
```

Output:

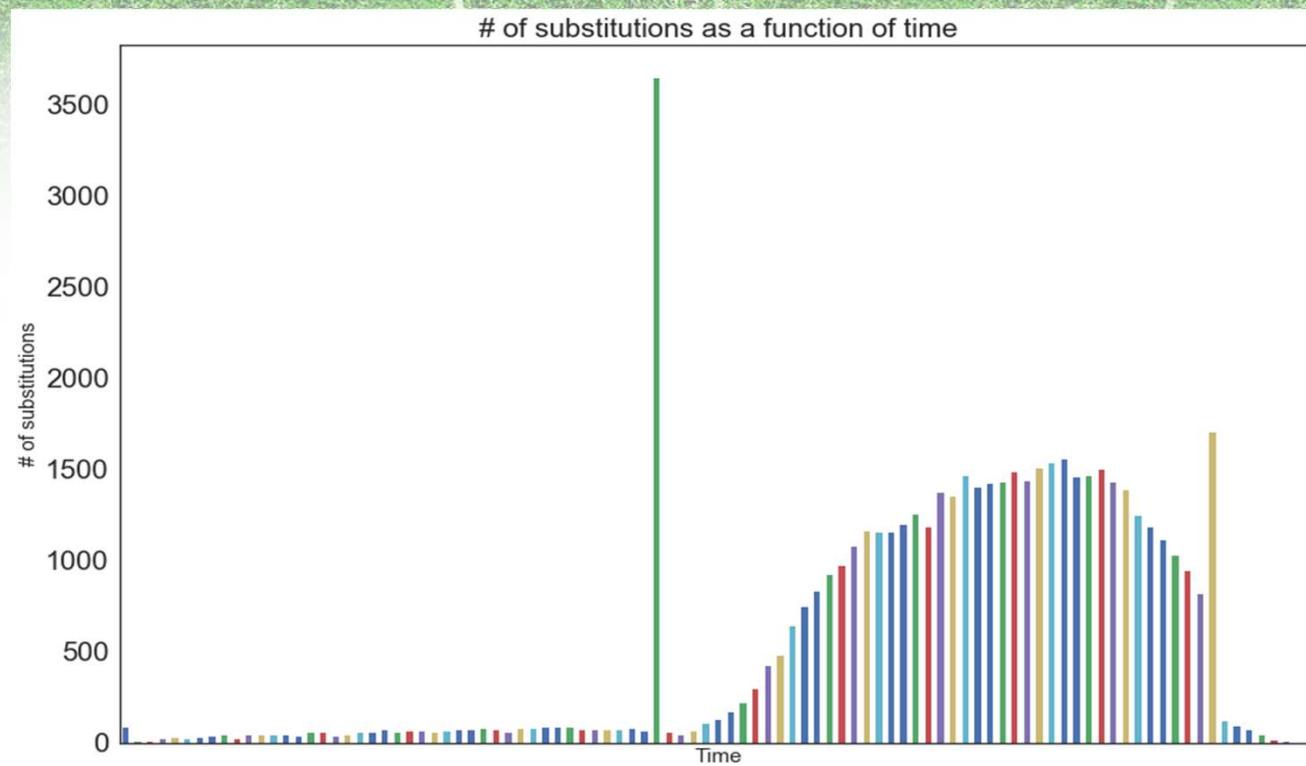
	Goals	Misses	Attempts	Assists	Cards	Fouls	Offsides
lionel messi	205	709	914	604	17	68	89
cristiano ronaldo	198	992	1190	862	29	100	206
zlatan ibrahimovic	153	621	774	487	29	248	225
gonzalo higuain	118	434	552	443	21	155	190
luis suarez	96	337	433	322	22	113	161
karim benzema	85	349	434	379	4	91	137
gareth bale	50	251	301	234	9	61	52
harry kane	65	288	353	278	9	85	78
neymar	58	260	318	248	23	106	69

When are the most goals scored during a football match ?



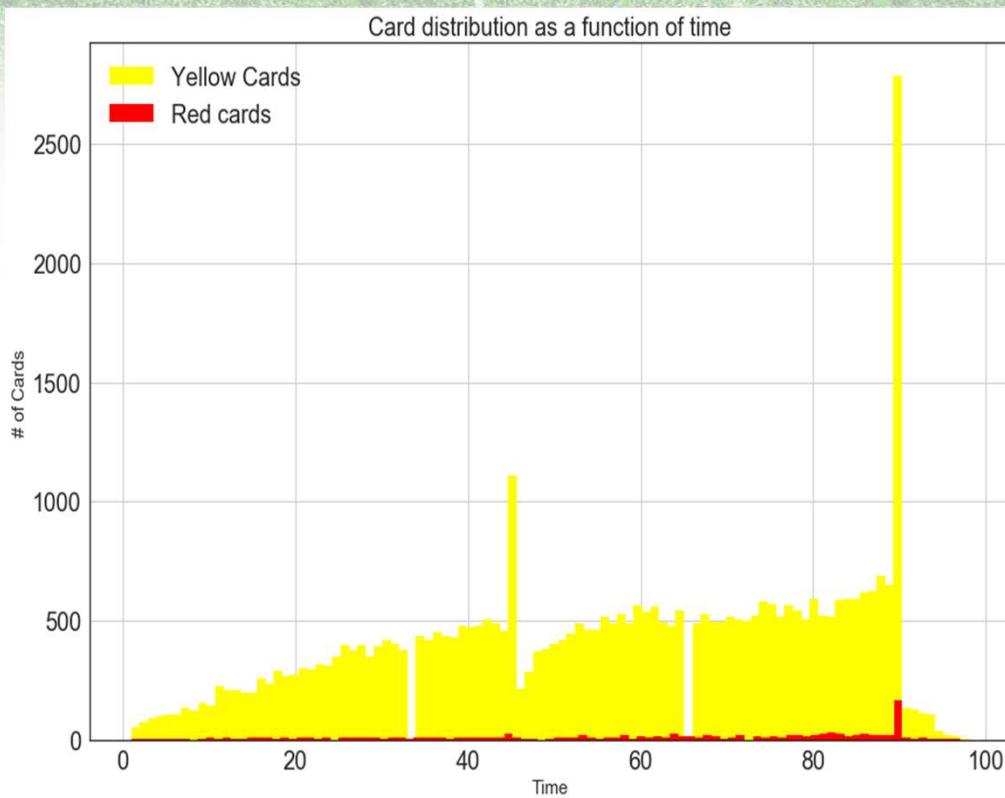
Interestingly there's a massive spike in number of goals just before half time and full time for both home and away teams

When do most player substitutions occur ?



Just like goal data, there's a massive spike in number of substitutions just around half time and full time.

When are the most number of red and yellow cards shown in a football match ?

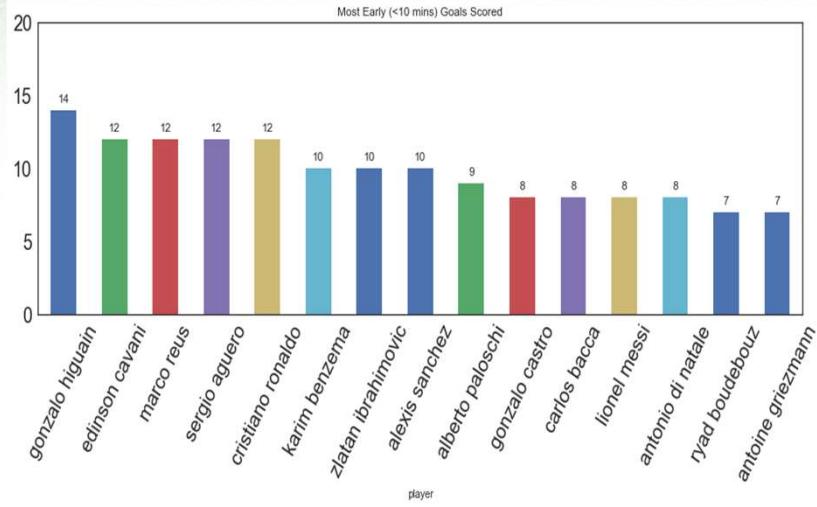


There's a spike in number of yellow cards just before half time and full time.

Most red cards are shown in last few minutes of a match (possibly when opponent is trailing by a goal and desperate to score and thus, the stakes are high)

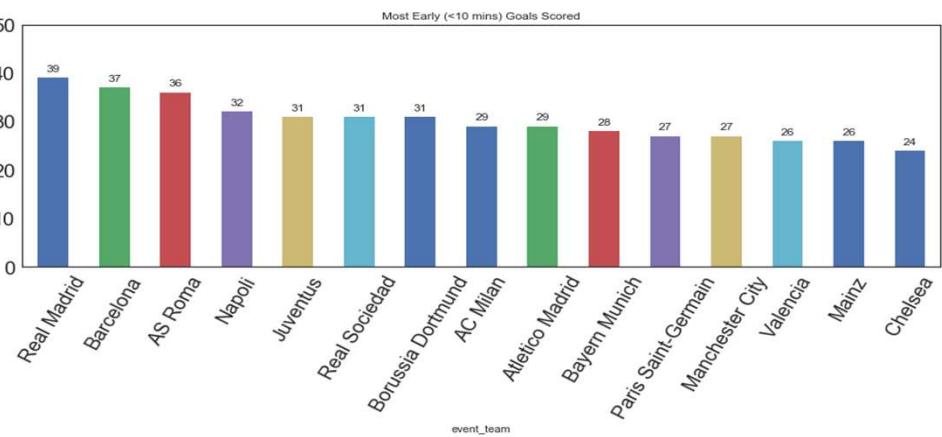
The early goal scorers

For our analysis, we are considering a goal scored during first 10 minutes of a game as an early goal



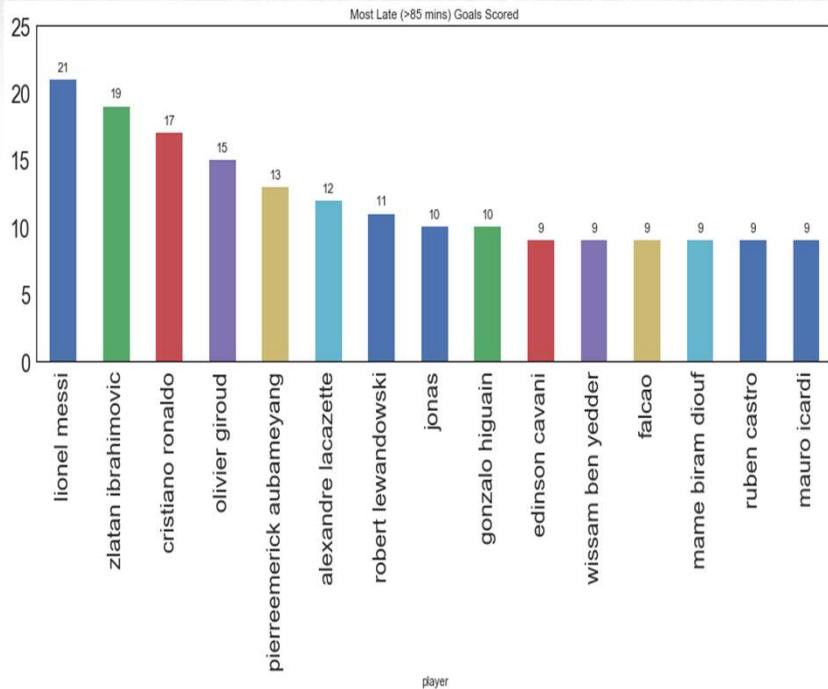
Real Madrid leads the early goal scorer table at club level

- There are very matches where a team scores a goal in first 10 minutes
- Finally we have some statistical yardstick for forwards which isn't led by Ronaldo or Messi or Zlatan
- Interestingly 3 out of top 15 early goal scorers are from Real Madrid

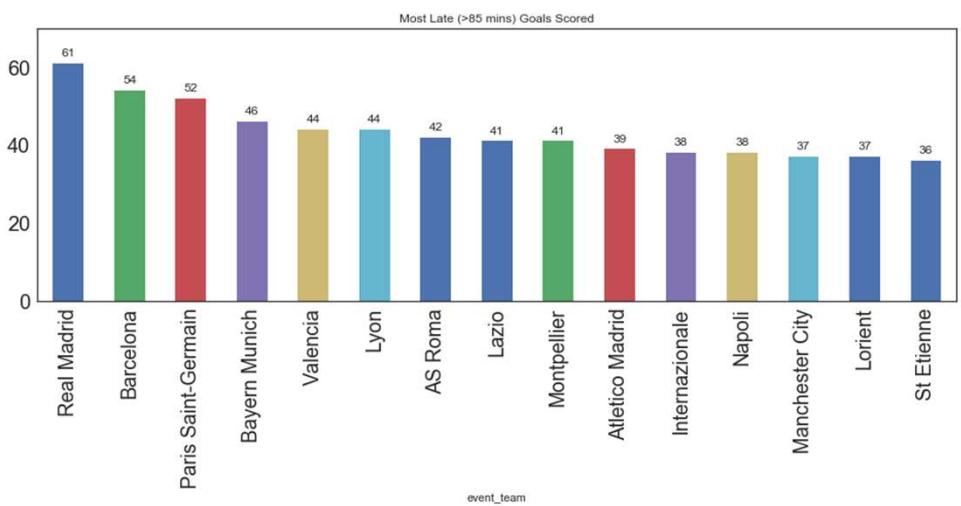


The last minute goal scorers

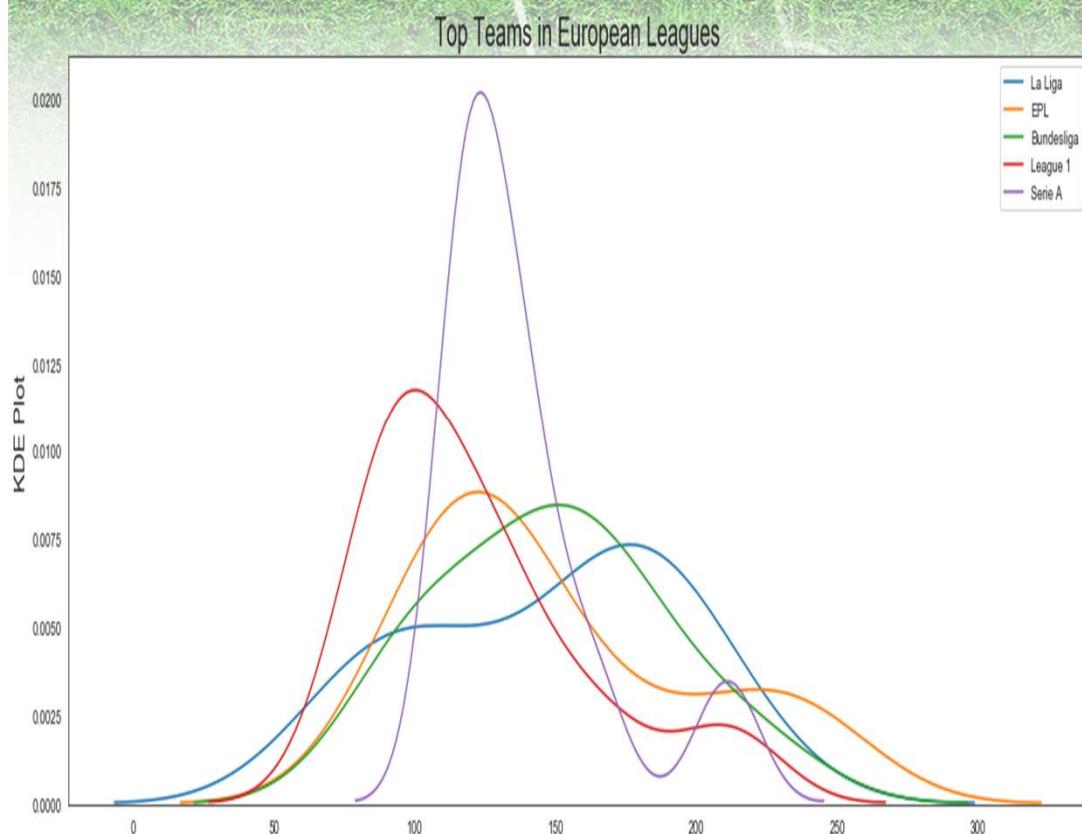
For our analysis, we are considering a goal scored after 85th minute of a game as a early goal



The forward trio are back in business again



Which league is the most competitive ?

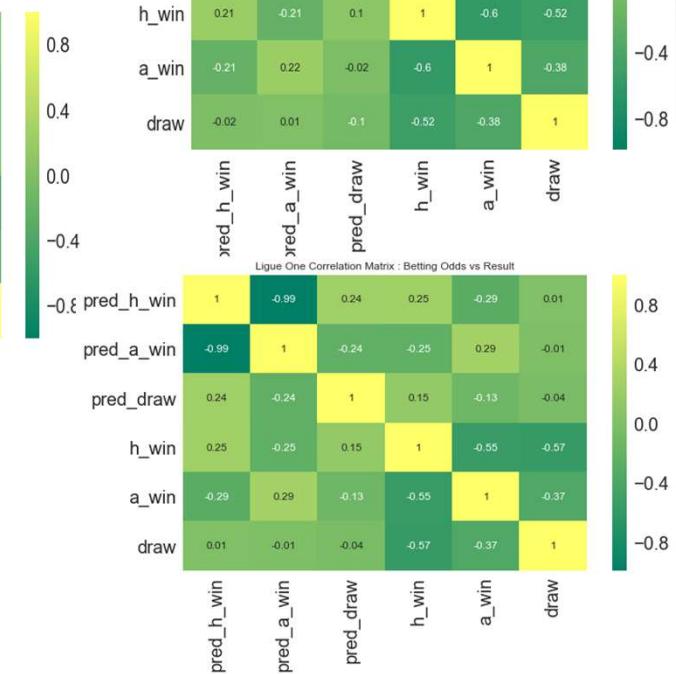
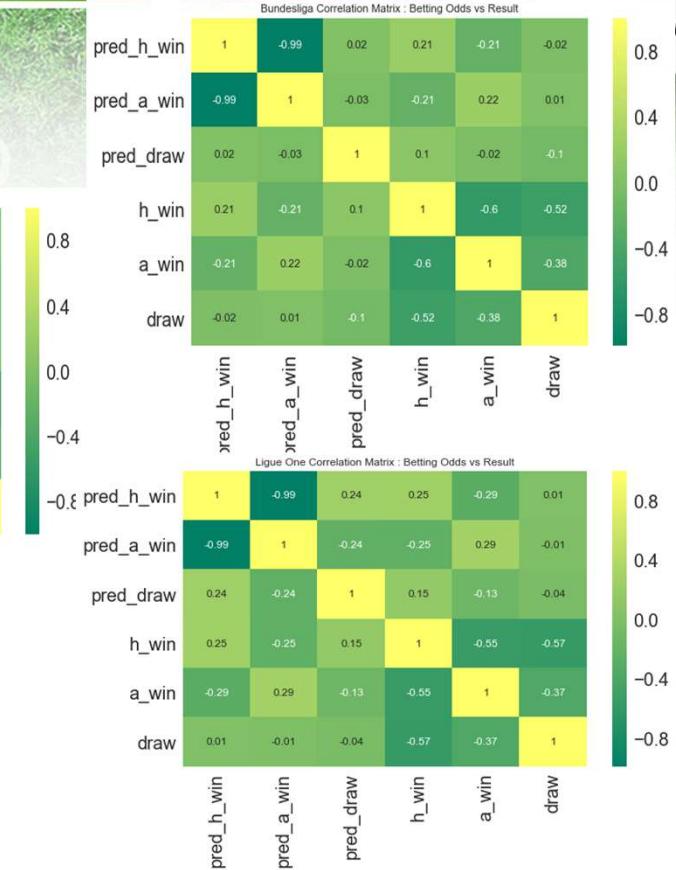
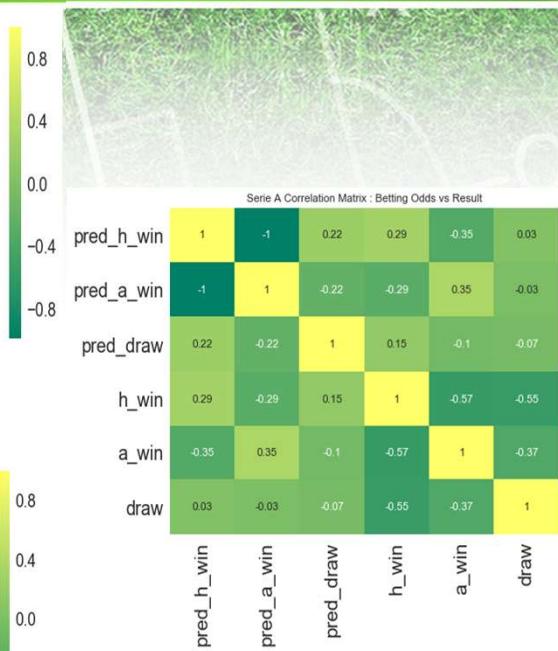
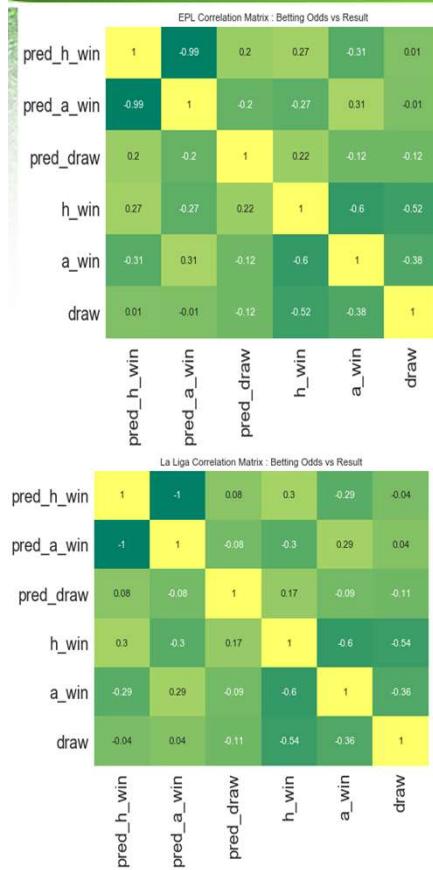


Leagues that are normally distributed are less competitive because the best and worst teams are equidistant from the mean

Left skewed distribution means only a few handful teams are close to the top whereas right skewed distribution means there's lot of competition at the top

From the plot it seems that La Liga and Bundesliga are the most competitive

Betting odds – are they indicative of results ?



Betting odds – are they indicative of results ?



- For every match we have 3 columns for betting odds : home_odds, away_odds and draw_odds
- The odds are expressed as numeric values (e.g a betting odd of 4 means, you'll win \$ 400 for every \$ 100 that you bet)
- From the odds, the probability of the event (as per bookies) can be calculated as :
$$P = 1/(1 + \text{odd})$$
So, a betting odd of 4 will translate to a probability of 0.2
- We converted the betting odd numbers into predicted outcome (lower odds = more likely outcome) and tried to find a correlation between the predicted and the actual outcome

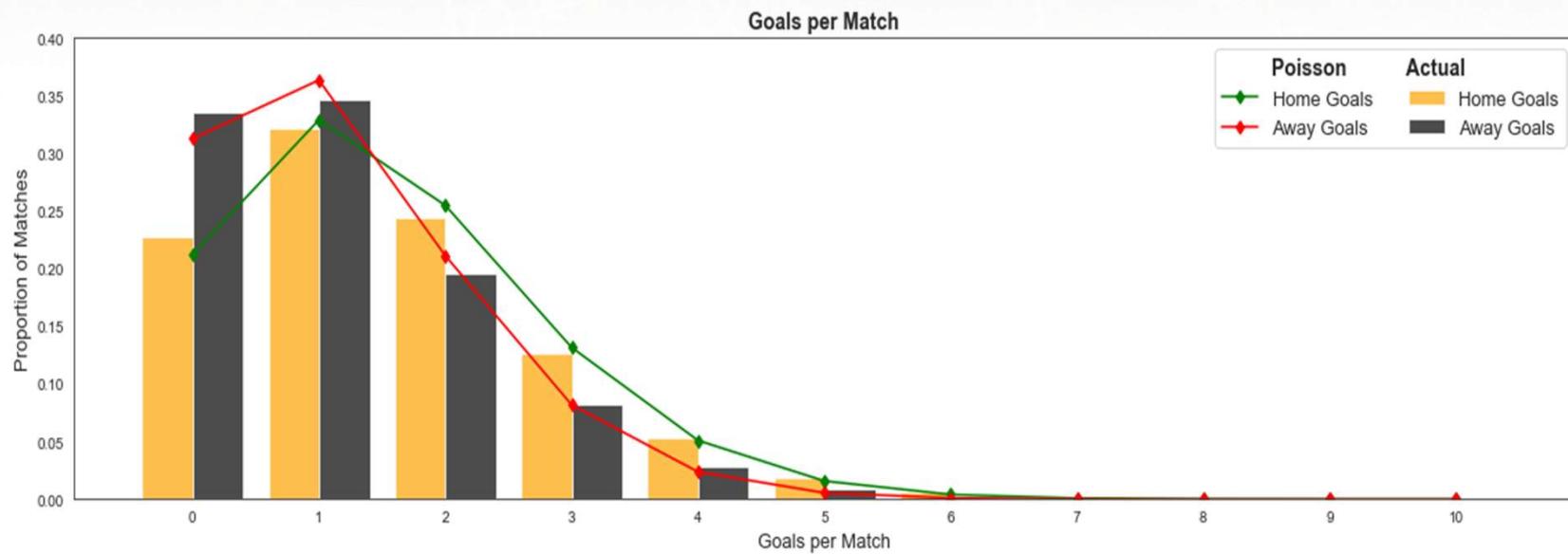
	EPL	LaLiga	Serie A	Bundesliga	Ligue One
Home Win	0.27	0.3	0.29	0.21	0.25
Away Win	0.31	0.29	0.35	0.22	0.29
Draw	-0.12	-0.11	-0.07	-0.1	-0.04

- Clearly betting odds have some correlation with actual results but aren't exactly indicative of results
- Moral of the story : It's probably not a sustainable idea to get rich by betting on football matches.

Predictive Modelling of Number of Goals



- Statistical studies on goal scoring pattern in football matches often suggest that goals in a football match follow Poisson distribution pattern.
- We continued with that assumption and use the mean expected number of goals for a match to form a probability distribution and express the number of goals as a function of average rate of scoring goals



Match simulation – Man City vs Man Utd at Etihad Stadium



Prediction by Model :

Probability Man City Win (Man City Home) : 0.5213758840709862
Probability Man U win (Man City Home) : 0.237715360715005
Probability Draw (Man City Home) 0.2409069842947822

Actual Results :



Match simulation –FC Barcelona vs Real Madrid at Camp Nou



Prediction by Model :

Probability Barcelona Win (Barcelona) : 0.6127844341417124
Probability Madrid win (Barcelona) : 0.20240895399328018
Probability Draw (Barcelona) 0.18475050565145482

Actual Results :

La Liga - 10-28

Full-time



5

1



Barcelona

Real Madrid

Philippe Coutinho 11'



Marcelo Vieira 50'

Luis Suárez 30' (P), 75', 83'

Arturo Vidal 87'

Football Data Analysis

Thank You