

Monte Carlo Simulations of the English Premier League

Guided by:
Prof. Prasanna Chaporkar

Submitted By :
Deep Karman Pal Singh | 16D070063

1 Objectives

English Premier league is the top league of football in England. It is widely regarded as one of the most competitive leagues in the world. Of interest in this league are question like - how many points does a team need (on average) to win the league, how many to qualify for the UEFA Champions League, for the Europa League and how many points does a team need to avoid relegation. Through this project, I wish to answer the above questions. I also answer that looking at the data from the last 2 years, what are the odds for each team to win the league.

2 Why Monte Carlo Methods

A direct computation of the probabilities of the state of the table at the end of the league is not feasible due to the large number of possible states.

Consider, after each match, three cases - wins, loss or a draw. This gives that each state (of the points table) evolves with games, and three possibilities for the next state. Over the 380 games of the league, this means that there are 3^{380} possible states, which is too large to compute the probabilities over. Note that here we ignored goal difference to simplify analysis, which would increase the number of states by a large amount. Hence direct computation is not feasible, and we move to Monte Carlo methods.

3 Our Approach to the Problem

3.1 Simulating the League

We consider 20 teams in the league, as per the **2017-18 season**. In our model we would ignore the form of a team, hence the order in which matches are played is irrelevant. We create a simple schedule for these 20 teams and simulate the season match by match.

3.2 Simulating a match

In each match, we consider each team scores goals as per a **Poisson distribution**, whose parameter we estimate as below. Based on this we consider the evolution of the points table. Note that we would need to keep a track of goals to resolve differences in the points table (based on goal difference, and then goals for). We do this for all of the 380 games in the season.

3.3 Estimating the Parameters

We assumed that the distribution of goals by a team over many matches is Poisson. We note by empirically plotting the data that this is accurate. Now we wish to estimate these parameters.

We need to about a team's attack strength (at home and away) and defensive strength (at home and away), hence we shall need 4 parameters for each team. We consider the average number of goal a home team scores, and the average an away team scores over the entire league. These form our base parameters. We shall weigh these with respect to the attacks and defences of teams playing to get the Poisson parameters of the team.

To model the attacking strength of a home (away) team, we consider its weight as the factor by which the team scores more home (away) goals on average over the average number of home (away) goals scored in the league.

Similarly, the defensive strength (we actually calculate what could be considered an inverse of the strength, since smaller value implies stronger defense) of a home (away) team is the factor by which average number of goals the home (away) team concedes is more than the average number of goals conceded by home (away) teams.

We now expand on the above explanation with the equations used, say for Manchester City (ManC) versus Liverpool (Liv) match

$$\begin{aligned}
 \lambda_{av,H} &= \text{Average number of Home Goals scored in the league} \\
 \lambda_{\text{ManC}, \text{Manc v Liv}} &= \lambda_{av,H} \times \text{AttSrt}_{\text{ManC},H} \times \text{DefStr}_{\text{Liv},A} \\
 \lambda_{\text{Liv}, \text{Manc v Liv}} &= \lambda_{av,A} \times \text{AttSrt}_{\text{Liv},A} \times \text{DefStr}_{\text{ManC},H} \\
 \text{AttSrt}_{\text{ManC},H} &= \text{Attacking Strength of ManC at Home} \\
 \text{AttSrt}_{\text{ManC},H} &= \frac{\text{Average number of goals ManC scores when playing at home}}{\text{Average number of goals scored by home teams}} \\
 \text{DefStr}_{\text{Liv},A} &= \text{Defensive Strength of Liv at Away} \\
 \text{DefStr}_{\text{Liv},A} &= \frac{\text{Average number of goals Liv concedes when playing away}}{\text{Average number of goals conceded by away team}}
 \end{aligned}$$

Note that this gives that an average team shall score and concede an average number of goals. This is consistent since the average of a Poisson distribution is its parameter. Running our model, we see that the results appear to be distributed like a real table, hence lending credibility to our assumptions.

3.4 Data Acquisition and Cleanup

We took data from the entire 2017-18 season, as well as from the 2018-19 season. The source is linked [here](#). Merging of the data from two seasons led to some points we had to consider, listed below.

First, the 2018-19 season is not complete, but ongoing. So we just took the season so far, incorporating all the games. This is up to games 37 for all teams, only 1 matchday is missing for all.

Next, due to relegation and promotion there is a difference of teams in the two seasons. To solve this, we removed all matches with the three newly promoted teams in the 2018-19 season. Hence we only had 262 matches from the 2018-19 season in our dataset, and all 380 matches from 2017-18 season.

Since the data for the three teams that got relegated was missing from the 2018-19 season, we did not add any bias for the newer season over the old. Note that the data imbalance does not effect the calculations of the attacking and defensive strengths since we only consider averages over games played.

4 Observations & Results

We ran two simulations, one only with data from 2017-18 season and one with data from both the seasons. In each case, we ran multiple simulations, over 1000, 10,000, and 1,000,000 iterations.

4.1 Only 2017-18 Season's Data

We observe that (for the simulation with 1 million leagues)

- Due to Manchester City's amazing record in the 2017-18 season (total 100 points, highest ever in the league), they won 89.8% of all simulations. Next were Liverpool and Manchester United, winning 4.7% and 3.3% of the simulations respectively. Tottenham won 1.9% of the time, Chelsea 0.2% and Arsenal 0.01%.
- No team other than these (not in the colloquial "Big Six") won the league in any simulation. This only serves to show why Leicester City's 2015-16 season, where they won the league was so exceptional (according to this, rarer than one in a million).
- On average, the second placed team made 85.17 points. Hence, on average, **scoring 86 points would win you the league**.
- To qualify for UEFA Champions League, a team needs to be in the top 4. The 5th placed team made, on average, 69.5 points. Hence to **qualify for Champions League**, a team needs, on average, **70 points**.

- Similarly, **for Europa League**, 7th placed team made 57.4 points. Hence to be in the Europa League, a team needs **58 points**, on average.
- The 18th placed team made, on average, 32.89 points. Hence, **to avoid relegation**, a team need **33 points**, on average.

4.2 Both 2017-18 and 2018-19 Season's Data

We observe that (for the simulation with 1 million leagues)

- Due to Manchester City's amazing record in the both the seasons, they won 80.8% of all simulations. Next were Liverpool and Tottenham, winning 17.3% and 1.4% of the simulations respectively. Manchester United won 0.23% of the time, Chelsea 0.21% and Arsenal 0.02%.
- No team other than these (not in the colloquial "Big Six") won the league in any simulation. Liverpool won a much larger proportion due to their great record in the 2018-19 season, while the reverse occurred for Manchester United.
- On average, the second placed team made 86.35 points. Hence, on average, **scoring 87 points would win you the league**. This increased due to the great record of two teams, Manchester City and Liverpool over the two seasons.
- To qualify for UEFA Champions League, a team needs to be in the top 4. The 5th placed team made, on average, 68.15 points. Hence to **qualify for Champions League**, a team need, on average, **69 points**. This fell due to a weak performance by all three of Chelsea, Arsenal, and Manchester United in the 2018-19 season.
- Similarly, **for Europa League**, 7th placed team made 58.33 points. Hence to be in the Europa League, a team needs **59 points**, on average.
- The 18th placed team made, on average, 32.65 points. Hence, **to avoid relegation**, a team need **33 points**, on average.

A Codes

```
1 import os
2 import numpy as np
3 import sys
4 import scipy
5 import csv
6 import time
7 np.set_printoptions(threshold=sys.maxsize)
8 num_iter = 1000000
9
10 start = time.time()
11
12 teams_to_index = ["Arsenal",
13 "Bournemouth",
14 "Brighton",
15 "Burnley",
16 "Chelsea",
17 "Crystal Palace",
18 "Everton",
19 "Huddersfield",
20 "Leicester",
21 "Liverpool",
22 "Man City",
23 "Man United",
24 "Newcastle",
25 "Southampton",
26 "Stoke",
27 "Swansea",
28 "Tottenham",
29 "Watford",
30 "West Brom",
31 "West Ham"]
32
33 # Global var for the params
34 attack_home = np.zeros(20)
35 attack_away = np.zeros(20)
36 defense_away = np.zeros(20)
37 defense_home = np.zeros(20)
38 av_home_goals = 0.0
39 av_away_goals = 0.0
40
41 def get_params():
42     filename = "season-1718.csv"
43     filename2 = "season-1819.csv"
44
```

```

45     fields1 = []
46     fields2 = []
47     rows = []
48     rows2 = []
49
50     with open(filename, 'r') as csvfile: # 17-18 season
51         # creating a csv reader object
52         csvreader = csv.reader(csvfile)
53
54         # extracting field names through first row
55         fields1 = csvreader.next()
56
57         # extracting each data row one by one
58         for row in csvreader:
59             rows.append(row)
60
61     with open(filename2, 'r') as csvfile: # 18-19 season,
62         uptill now
63         csvreader = csv.reader(csvfile)
64
65         fields2 = csvreader.next()
66
67         for row in csvreader:
68             rows2.append(row)
69
70     fields = np.asarray(fields1)
71     rows = np.asarray(rows)
72     rows2 = np.asarray(rows2)
73
74     fields = fields[1:6]
75     rows = rows[:,1:6]
76     rows2 = rows2[:,1:6]
77     # HomeTeam, AwayTeam, HomeGoals, AwayGoals, Result
78
79     tot_games = 0
80     tot_home_goals = 0.0 # away conceded = home scored
81     tot_away_goals = 0.0 # home conceded = away scored
82     tot_games_home_perteam = np.zeros(20) # total number of
83         games each team played home
84     tot_games_away_perteam = np.zeros(20) # total number of
85         games each team played away
86     tot_home_perteam = np.zeros(20) # total number of goals
87         each scored at their home
88     tot_away_perteam = np.zeros(20) # total number of away
89         goals each team scored
90     tot_conceded_home_perteam = np.zeros(20) # total number of

```

```

    goals team conceded at home
86 tot_conceded_away_perteam = np.zeros(20) # total number of
    goals team conceded away
87 global av_home_goals
88 global av_away_goals
89
90 for i in rows:
91     tot_home_goals += int(i[2])
92     tot_away_goals += int(i[3])
93     tot_games += 1
94
95     home_index = teams_to_index.index(i[0])
96     away_index = teams_to_index.index(i[1])
97
98     tot_games_home_perteam[home_index] += 1
99     tot_games_away_perteam[away_index] += 1
100
101     tot_home_perteam[home_index] += int(i[2])
102     tot_away_perteam[away_index] += int(i[3])
103
104     tot_conceded_home_perteam[home_index] += int(i[3])
105     tot_conceded_away_perteam[away_index] += int(i[2])
106
107 for i in rows2:
108     tot_home_goals += int(i[2])
109     tot_away_goals += int(i[3])
110     tot_games += 1
111
112     home_index = teams_to_index.index(i[0])
113     away_index = teams_to_index.index(i[1])
114
115     tot_games_home_perteam[home_index] += 1
116     tot_games_away_perteam[away_index] += 1
117
118     tot_home_perteam[home_index] += int(i[2])
119     tot_away_perteam[away_index] += int(i[3])
120
121     tot_conceded_home_perteam[home_index] += int(i[3])
122     tot_conceded_away_perteam[away_index] += int(i[2])
123
124 av_home_goals = float(tot_home_goals/tot_games) # sum of
    all home goals/number of games
125 av_away_goals = float(tot_away_goals/tot_games) # sum of
    all away goals/number of games
126
127

```

```

128     for i in range(20):
129         attack_home[i] = (tot_home_perteam[i]/
130                             tot_games_home_perteam[i])/av_home_goals
131         attack_away[i] = (tot_away_perteam[i]/
132                             tot_games_away_perteam[i])/av_away_goals
133
134         defense_away[i] = (tot_conceded_away_perteam[i]/
135                             tot_games_away_perteam[i])/av_home_goals
136         defense_home[i] = (tot_conceded_home_perteam[i]/
137                             tot_games_home_perteam[i])/av_away_goals
138
139     return
140 get_params()
141
142 # print(av_home_goals)
143 # print(av_away_goals)
144 # print(attack_home)
145 # print(attack_away)
146 # print(defense_home)
147 # print(defense_away)
148 # 0 to 19
149
150 # Set up the schedule – form not considered, does not make a
151     difference
152     schedule = np.asarray([[0,1]])
153     for i in range(20):
154         for j in range(20):
155             # print(i,j)
156             if ((i == (j-1)) and (i == 0)):
157                 temp = np.asarray([[i,j]])
158             elif (i != j) :
159                 schedule = np.append(schedule,[[i,j]], axis = 0)
160     np.random.shuffle(schedule)
161
162 def league_sim():
163     # Set up data structure for the table
164     table = np.zeros((20,9)) # ordered as per teams to index
165     # Played, Wins, Draws, Losses, Goals For, Goals Against,
166     # Goal Diff, Points, index
167     for i in range(20):
168         (table[i])[8] = i
169
170     #print(schedule)
171
172     for match in schedule:
173         hometeam_stats = table[match[0]]

```



```

168 awayteam_stats = table[match[1]]
169
170 # both play a game
171 hometeam_stats[0] = hometeam_stats[0] + 1
172 awayteam_stats[0] = awayteam_stats[0] + 1
173
174 # Get the goals both teams scored
175 goals_home = np.random.poisson((av_home_goals*
    attack_home[match[0]]*defense_away[match[1]])) # TODO
    - how are goals generated
176 goals_away = np.random.poisson((av_away_goals*
    attack_away[match[1]]*defense_home[match[0]])) # TODO
    - how are goals generated
177
178 if (goals_home > goals_away):
179     hometeam_stats[1] += 1 # home won
180     awayteam_stats[3] += 1 # away lost
181
182     hometeam_stats[4] += goals_home # GF home
183     hometeam_stats[5] += goals_away # GA home
184
185     awayteam_stats[4] += goals_away # GF away
186     awayteam_stats[5] += goals_home # GA away
187
188     hometeam_stats[6] += (goals_home - goals_away) # GD
        home
189     awayteam_stats[6] -= (goals_home - goals_away) # GD
        away
190
191     hometeam_stats[7] += 3 # Winner gets points
192 elif (goals_home == goals_away):
193     hometeam_stats[2] += 1 # home drew
194     awayteam_stats[2] += 1 # away drew
195
196     hometeam_stats[4] += goals_home # GF home
197     hometeam_stats[5] += goals_away # GA home
198
199     awayteam_stats[4] += goals_away # GF away
200     awayteam_stats[5] += goals_home # GA away
201
202     hometeam_stats[7] += 1 # both get points
203     awayteam_stats[7] += 1
204 else : #(goals_home < goals_away)
205     hometeam_stats[3] += 1 # home lost
206     awayteam_stats[1] += 1 # away won
207

```

```

208     hometeam_stats[4] += goals_home # GF home
209     hometeam_stats[5] += goals_away # GA home
210
211     awayteam_stats[4] += goals_away # GF away
212     awayteam_stats[5] += goals_home # GA away
213
214     hometeam_stats[6] += (goals_home - goals_away) # GD
        home
215     awayteam_stats[6] -= (goals_home - goals_away) # GD
        away
216
217     awayteam_stats[7] += 3 # Winner gets points
218
219     table[match[0]] = hometeam_stats
220     table[match[1]] = awayteam_stats
221
222     #print(table)
223
224     # use mergesort as it is a stable sort
225     # the intial sorting remains if later sorting is not able
        to resolve
226     # last by Points, then GD, GF, GA
227     table = table[table[:,3].argsort(kind = 'mergesort')
        [::-1]] # sort by losses, reverse
228     table = table[table[:,1].argsort(kind = 'mergesort')] #
        sort by wins
229     table = table[table[:,4].argsort(kind = 'mergesort')] #
        sort by GF
230     table = table[table[:,6].argsort(kind = 'mergesort')] #
        sort by GD
231     table = table[table[:,7].argsort(kind = 'mergesort')
        [::-1]] # sort by points
232     # print(table) # if all above same, alphabetical
233     return table
234
235 num_wins = np.zeros(20)
236 av_points = np.zeros(20)
237
238 for i in range(num_iter):
239     table = league_sim()
240     num_wins[int(table[0][8])] += 1
241     for j in range(20):
242         av_points[j] += table[j][7]
243
244 av_points = av_points/num_iter
245

```

```

246 print(num_wins)
247 print(av_points)
248
249 with open("op1.txt","a") as file:
250     file.write("new sim starts, num iter is %d\n"% num_iter)
251     file.write("num wins is\n")
252     for i in range(20):
253         file.write("%s - %d \n"% (teams_to_index[i], num_wins[i]
254                                     ]))
255     file.write(" \n")
256     file.write("Average points total is\n")
257     for i in range(20):
258         file.write("%dth place - %6.2f \n"% (i, av_points[i]))
259     file.write(" \n")
260
261
262
263 end = time.time()
264 print(end-start)

```