

1 Predicting the 2024 UEFA Euro Using Poisson Distribution

Football predictions have always been a fascinating topic for data scientists and enthusiasts alike. In this project, we leverage the power of the Poisson distribution to forecast the outcomes of the 2024 UEFA European Championship. This method is particularly well-suited for predicting football scores due to the nature of the game, where goals are relatively rare events.

```
In [1]: 1 #installing the necessary libraries
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 from scipy.stats import poisson
6 from sklearn.metrics import accuracy_score
7 from sklearn.metrics import roc_curve, roc_auc_score
8 import matplotlib.colors as mcolors
9 import warnings
10 pd.set_option('display.max_rows', 10000)
11 pd.set_option('display.max_columns', 500)
12 warnings.filterwarnings('ignore')
```

2 Data Collection

The first step in our analysis involved gathering historical match data. The dataset sourced from Kaggle contains the results of international football matches starting from the very first official match in 1872 up to 2024.

```
In [2]: 1 #Reading the
2 df = pd.read_csv('results.csv', sep=',',)
3 #Dropping matches without final result
4 df=df[~df['home_score'].isna()]
```

```
In [3]: 1 df
```

Out[3]:

	date	home_team	away_team	home_score	away_score	tournament	city	cour
0	1872-11-30	Scotland	England	0.0	0.0	Friendly	Glasgow	Scotl
1	1873-03-08	England	Scotland	4.0	2.0	Friendly	London	Engl
2	1874-03-07	Scotland	England	2.0	1.0	Friendly	Glasgow	Scotl
3	1875-03-06	England	Scotland	2.0	2.0	Friendly	London	Engl
4	1876-03-04	Scotland	England	3.0	0.0	Friendly	Glasgow	Scotl
...
47365	2024-07-02	Costa Rica	Paraguay	2.0	1.0	Copa América	Austin	Un St
47366	2024-07-04	Argentina	Ecuador	1.0	1.0	Copa América	Houston	Un St
47367	2024-07-05	Germany	Spain	1.0	2.0	UEFA Euro	Stuttgart	Germ
47368	2024-07-05	Portugal	France	0.0	0.0	UEFA Euro	Hamburg	Germ
47369	2024-07-05	Venezuela	Canada	1.0	1.0	Copa América	Arlington	Un St

47370 rows × 9 columns

3 ELO Ranking System

The Elo ranking system is a method for calculating the relative skill levels of teams. It is useful for ranking teams based on their performance in previous matches. The Elo rating is dynamic, meaning it changes based on the results of matches, rewarding teams for wins and penalizing them for losses, while also considering the strength of the opponents.

To create our dataset, we need to define some parameters, such as the importance of each tournament in the points exchanged by the Elo ranking system.

```

In [4]: 1 #Auxiliary funtions for update the ratings
2 confederation_tournaments=['AFC Asian Cup','African Cup of Nations','U
3 confederation_clasification=['UEFA Euro qualification','African Cup o
4 def k_value(tournament):
5     k=5
6     if tournament == 'Friendly':
7         k=10
8     elif tournament in confederation_clasification:
9         k=20
10    elif tournament == 'FIFA World Cup qualification':
11        k=25
12    elif tournament in confederation_tournaments:
13        k=40
14    elif tournament == 'FIFA World Cup':
15        k=50
16    return k
17
18 def expected_result(loc,aw):
19     dr=loc-aw
20     we=(1/(10**(-dr/400)+1))
21     return [np.round(we,3),1-np.round(we,3)]
22
23 def actual_result(loc,aw):
24     if loc<aw:
25         wa=1
26         wl=0
27     elif loc>aw:
28         wa=0
29         wl=1
30     elif loc==aw:
31         wa=0.5
32         wl=0.5
33     return [wl,wa]
34
35 def calculate_elo(elo_l,elo_v,local_goals,away_goals,tournament):
36
37     k=k_value(tournament)
38     wl,wv=actual_result(local_goals,away_goals)
39     wel,wev=expected_result(elo_l,elo_v)
40
41     elo_ln=elo_l+k*(wl-wel)
42     elo_vn=elo_v+k*(wv-wev)
43
44     return elo_ln,elo_vn,wel,wev

```

```

In [5]: 1 #Calculating the elo points for all historical international matches
        2 current_elo={}
        3 for idx,row in df.iterrows():
        4
        5     local=row['home_team']
        6     away=row['away_team']
        7     local_goals=row['home_score']
        8     away_goals=row['away_score']
        9     tournament = row['tournament']
        10
        11
        12     if local not in current_elo.keys():
        13         current_elo[local]=1300
        14
        15     if away not in current_elo.keys():
        16         current_elo[away]=1300
        17
        18     elo_l=current_elo[local]
        19     elo_v=current_elo[away]
        20     elo_ln,elo_vn, wel,wev=calculate_elo(elo_l,elo_v,local_goals,away_
        21
        22     current_elo[local]=elo_ln
        23     current_elo[away]=elo_vn
        24
        25     df.loc[idx,'Elo_h_after']=elo_ln
        26     df.loc[idx,'Elo_a_after']=elo_vn
        27     df.loc[idx,'Elo_h_before']=elo_l
        28     df.loc[idx,'Elo_a_before']=elo_v
        29     df.loc[idx,'probH']=wel
        30     df.loc[idx,'probA']=wev

```

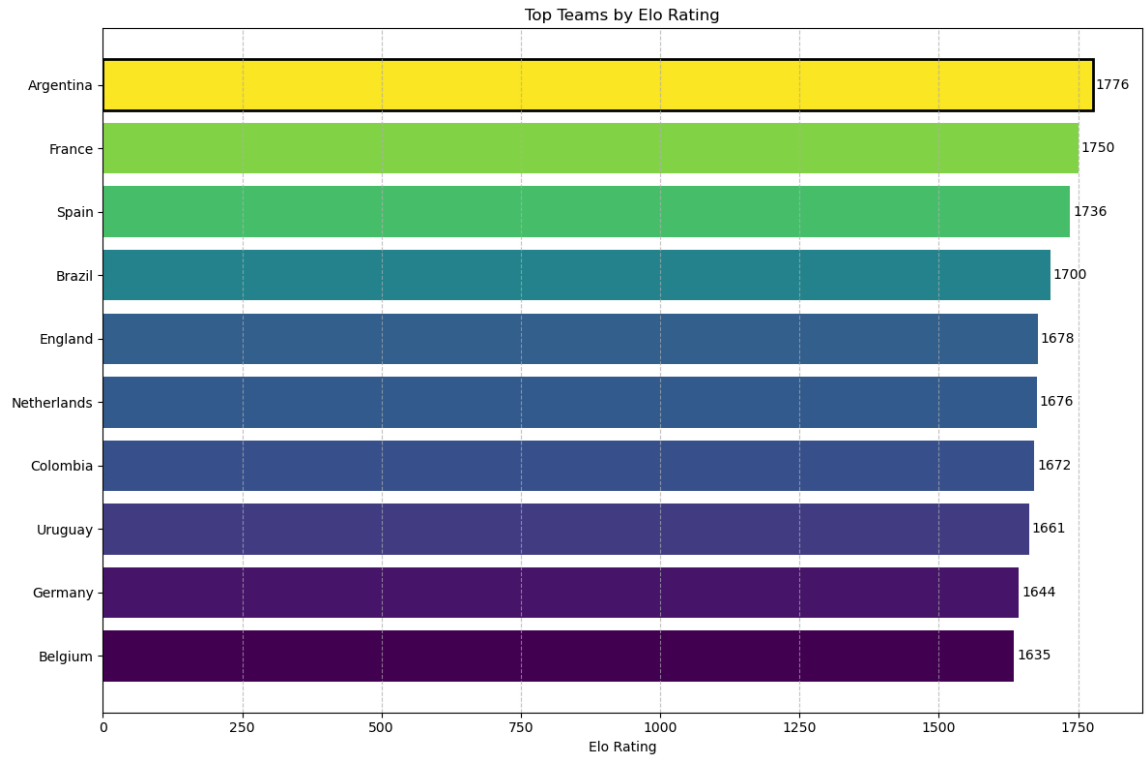
Once we have processed the historical data of international matches, we can calculate the final Elo points for each team. This gives us the current Elo ranking, which reflects the latest performance levels of the teams.

```

In [6]: 1 elos=pd.concat([df[['date','home_team','Elo_h_after']].rename(columns:
        2 elos.sort_values(by='date', ascending=False,inplace=True)
        3 elos.drop_duplicates('Team',inplace=True)
        4 elos.sort_values(by='Elo', ascending=False, inplace=True)
        5 elos.reset_index(drop=True, inplace=True)
        6 elos['position']=elos.index+1

```

```
In [7]: 1 def plot_top_elo(elos, n=10):
2         """
3         Plots the top n teams by Elo rating in a horizontal bar chart.
4
5         Parameters:
6         elos (pd.DataFrame): DataFrame containing 'Team' and 'Elo' columns
7         n (int): Number of top teams to display.
8         """
9         # Sort elos DataFrame by Elo rating in descending order
10        elos.sort_values(by='Elo', ascending=False, inplace=True)
11
12        # Select only the top n teams
13        top_teams = elos.head(n)
14
15        # Create a color map
16        cmap = plt.get_cmap('viridis')
17        norm = mcolors.Normalize(vmin=min(top_teams['Elo']), vmax=max(top_teams['Elo']))
18        colors = cmap(norm(top_teams['Elo']))
19
20        # Plotting
21        plt.figure(figsize=(12, 8))
22        bars = plt.barh(top_teams['Team'], top_teams['Elo'], color=colors)
23
24        # Add text annotations
25        for bar in bars:
26            plt.text(
27                bar.get_width() + 5,
28                bar.get_y() + bar.get_height() / 2,
29                f'{bar.get_width():.0f}',
30                ha='left',
31                va='center',
32                fontsize=10
33            )
34
35        # Highlighting the top team
36        bars[0].set_edgecolor('black')
37        bars[0].set_linewidth(2)
38
39        # Customizing the plot
40        plt.xlabel('Elo Rating')
41        plt.title('Top Teams by Elo Rating')
42        plt.gca().invert_yaxis() # Invert y-axis to display the highest Elo at the top
43        plt.grid(axis='x', linestyle='--', alpha=0.7)
44
45        # Show plot
46        plt.tight_layout()
47        plt.show()
48
49        # Call the function
50        plot_top_elo(elos, 10)
51
```

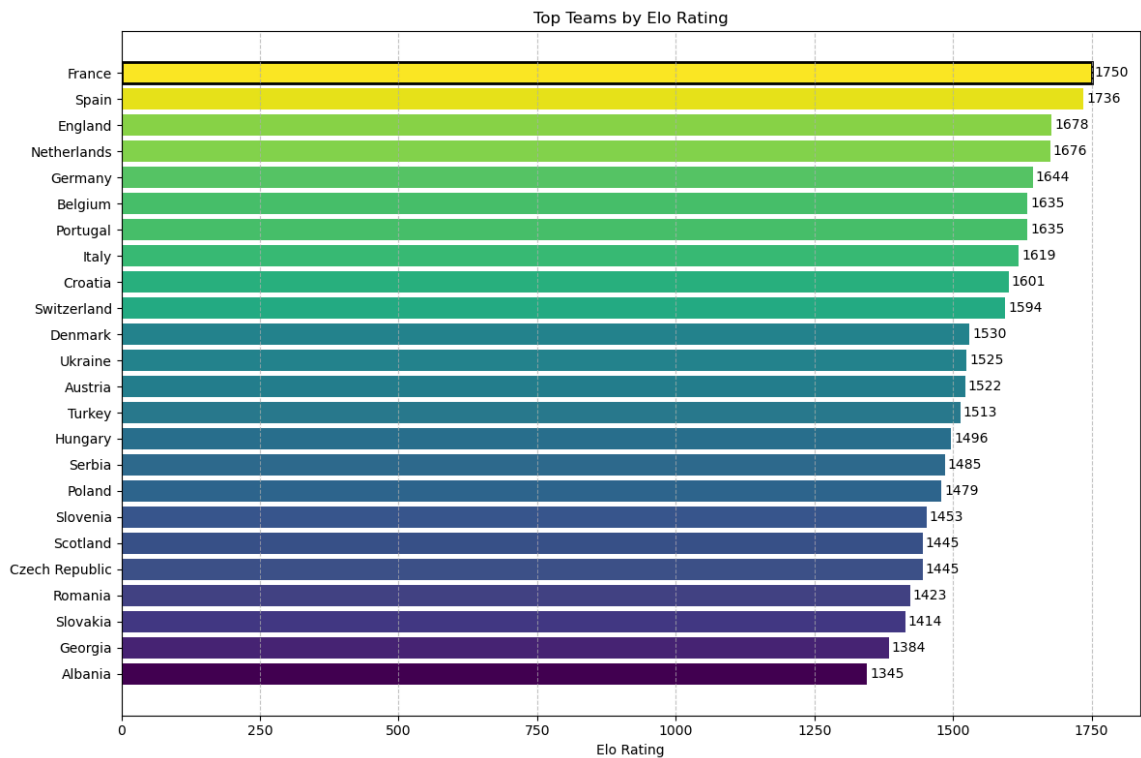


We will focus only on the teams participating in the 2024 Euro. Therefore, the ranking would be as follows:

```

In [8]: 1 eurocup_teams =[
2         "Albania", "Scotland", "Hungary", "Romania",
3         "Germany", "Slovakia", "England", "Czech Republic",
4         "Austria", "Slovenia", "Italy", "Serbia",
5         "Belgium", "Spain", "Netherlands", "Switzerland",
6         "Croatia", "France", "Poland", "Turkey",
7         "Denmark", "Georgia", "Portugal", "Ukraine"
8     ]
9
10 elos_ranking_euro = elos[elos.Team.isin(eurocup_teams)].copy()
11 plot_top_elo(elos_ranking_euro,24)

```



Based on this ranking, we can say that France is the favorite to win the Euro, while Albania is the weakest team in the tournament.

One of the problems with international matches is that we don't have the statistics of a league, making it difficult to determine which team has a better attack because it depends on the quality of the opponent. To solve this, and using the same idea as the Elo ranking system above, we can create a ranking for attacking power and defensive power.

In [9]:

1 elos_ranking_euro

Out[9]:

	date	Team	Elo	position
1	2024-07-05	France	1750.370	2
2	2024-07-05	Spain	1735.715	3
4	2024-06-30	England	1677.695	5
5	2024-07-02	Netherlands	1675.520	6
8	2024-07-05	Germany	1643.990	9
9	2024-07-01	Belgium	1634.765	10
10	2024-07-05	Portugal	1634.685	11
11	2024-06-29	Italy	1619.060	12
13	2024-06-24	Croatia	1600.850	14
14	2024-06-29	Switzerland	1594.235	15
22	2024-06-29	Denmark	1529.990	23
23	2024-06-26	Ukraine	1524.660	24
24	2024-07-02	Austria	1521.590	25
28	2024-07-02	Turkey	1513.125	29
32	2024-06-23	Hungary	1496.465	33
35	2024-06-25	Serbia	1485.295	36
37	2024-06-25	Poland	1479.350	38
51	2024-07-01	Slovenia	1452.545	52
55	2024-06-23	Scotland	1445.070	56
56	2024-06-26	Czech Republic	1444.845	57
62	2024-07-02	Romania	1423.000	63
64	2024-06-30	Slovakia	1413.625	65
72	2024-06-30	Georgia	1384.090	73
82	2024-06-24	Albania	1345.260	83

In [10]:

```

1 def calculate_attdef(Attack_l,Attack_v,Deffend_l,Deffend_v,localGoals
2     k=k_value(tournament)
3     if neutral==False:
4         c=1.27
5     else:
6         c=1
7     ehg=Attack_l*Deffend_v*c
8     eag=Attack_v*Deffend_l/c
9     Attack_ln=Attack_l+(k/2000)*(localGoals-ehg)
10    Attack_vn=Attack_v+(k/2000)*(awayGoals-eag)
11    Deffend_ln=Deffend_l+(k/2000)*(awayGoals-eag)
12    Deffend_vn=Deffend_v+(k/2000)*(localGoals-ehg)
13
14    return Attack_ln,Attack_vn,Deffend_ln,Deffend_vn, ehg, eag

```



```

In [11]: 1 current_att={}
          2 current_def={}
          3 for idx,row in df.iterrows():
          4
          5     local=row['home_team']
          6     away=row['away_team']
          7     local_goals=row['home_score']
          8     away_goals=row['away_score']
          9     tournament=row['tournament']
         10     neutral=row['neutral']
         11
         12     if local not in current_att.keys():
         13         current_att[local]=1.3
         14     if away not in current_att.keys():
         15         current_att[away]=1.3
         16     if local not in current_def.keys():
         17         current_def[local]=1.3
         18     if away not in current_def.keys():
         19         current_def[away]=1.3
         20
         21     att_l=current_att[local]
         22     att_v=current_att[away]
         23     def_l=current_def[local]
         24     def_v=current_def[away]
         25
         26     att_ln,att_vn, def_ln, def_vn, ehg, eag=calculate_attdef(att_l,att_v,def_l,def_v,ehg,eag)
         27
         28     current_att[local]=att_ln
         29     current_att[away]=att_vn
         30     current_def[local]=def_ln
         31     current_def[away]=def_vn
         32
         33     df.loc[idx,'att_h_after']=att_ln
         34     df.loc[idx,'att_a_after']=att_vn
         35     df.loc[idx,'att_h_before']=att_l
         36     df.loc[idx,'att_a_before']=att_v
         37     df.loc[idx,'def_h_after']=def_ln
         38     df.loc[idx,'def_a_after']=def_vn
         39     df.loc[idx,'def_h_before']=def_l
         40     df.loc[idx,'def_a_before']=def_v
         41     df.loc[idx,'XGhome']=ehg
         42     df.loc[idx,'XGaway']=eag

```

```

In [12]: 1 attdef=df.copy()
          2 attdef=pd.concat([df[['date','home_team','att_h_after','def_h_after']]
          3 attdef.sort_values(by='date', ascending=False,inplace=True)
          4 attdef.drop_duplicates('Team',inplace=True)
          5 attdef.sort_values(by='def', ascending=True,inplace=True)
          6 attdef.reset_index(drop=True, inplace=True)
          7 attdef['def_position']=attdef.index+1
          8 attdef.sort_values(by='att', ascending=False,inplace=True)
          9 attdef.reset_index(drop=True, inplace=True)
         10 attdef['att_position']=attdef.index+1

```


In [13]:

```

1 import matplotlib.pyplot as plt
2 import pandas as pd
3
4 def plot_top_att_def(elos, n=10):
5     """
6     Plot the top teams by attacking and defensive power.
7
8     Parameters:
9     elos (pd.DataFrame): DataFrame containing team ratings.
10    n (int): Number of top teams to display.
11    """
12    # Sort elos DataFrame by attacking and defensive ratings
13    elos_sorted_att = elos.sort_values(by='att', ascending=False).head(n)
14    elos_sorted_def = elos.sort_values(by='def', ascending=True).head(n)
15
16    # Create a color map for both plots
17    colors_att = plt.cm.viridis(elos_sorted_att['att'] / max(elos_sorted_att['att']))
18    colors_def = plt.cm.viridis(elos_sorted_def['def'] / max(elos_sorted_def['def']))
19
20    # Create subplots
21    fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(20, 10))
22
23    # Plot for attacking ratings
24    bars_att = ax1.barh(elos_sorted_att['Team'], elos_sorted_att['att'], color=colors_att)
25    ax1.set_xlabel('Attacking Power', fontsize=14)
26    ax1.set_title('Top Teams by Attacking Power', fontsize=16, weight='bold')
27    ax1.invert_yaxis() # Invert y-axis to display the highest rating at the top
28    ax1.grid(axis='x', linestyle='--', alpha=0.7)
29
30    # Add text annotations for attacking ratings
31    for bar in bars_att:
32        ax1.text(
33            bar.get_width() + 0.01,
34            bar.get_y() + bar.get_height() / 2,
35            f'{bar.get_width():.2f}',
36            ha='left',
37            va='center',
38            fontsize=12,
39            color='black'
40        )
41
42    # Plot for defensive ratings
43    bars_def = ax2.barh(elos_sorted_def['Team'], elos_sorted_def['def'], color=colors_def)
44    ax2.set_xlabel('Defensive Power', fontsize=14)
45    ax2.set_title('Top Teams by Defensive Power', fontsize=16, weight='bold')
46    ax2.invert_yaxis() # Invert y-axis to display the highest rating at the top
47    ax2.grid(axis='x', linestyle='--', alpha=0.7)
48
49    # Add text annotations for defensive ratings
50    for bar in bars_def:
51        ax2.text(
52            bar.get_width() + 0.01,
53            bar.get_y() + bar.get_height() / 2,
54            f'{bar.get_width():.2f}',
55            ha='left',
56            va='center',
57            fontsize=12,
58            color='black'
59        )
60
61    # Adjust Layout

```

```

62 plt.tight_layout()
63 plt.show()
64

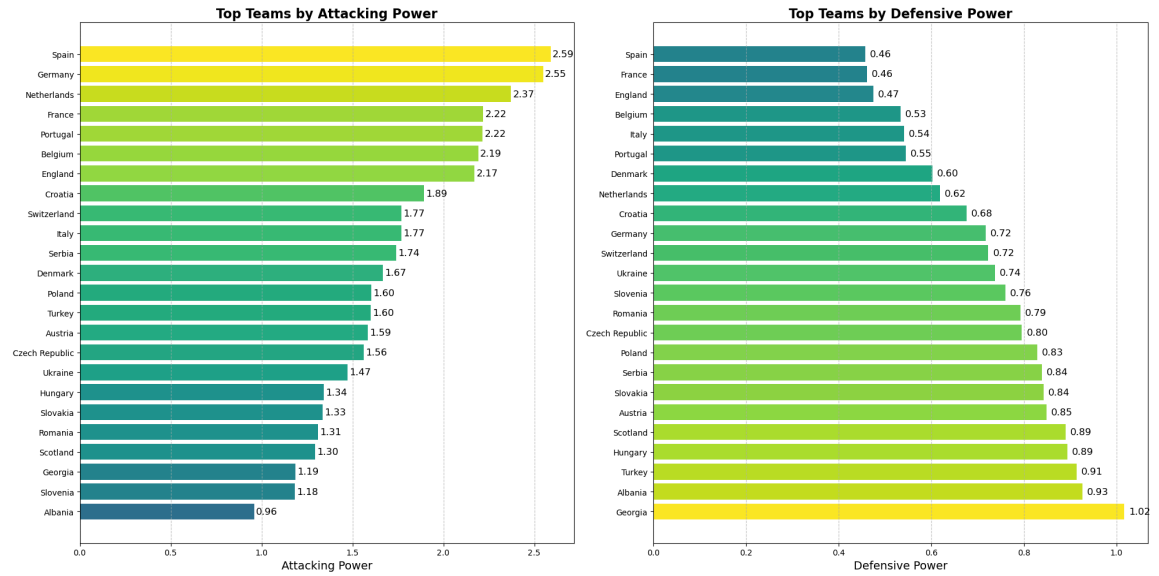
```

In [14]:

```

1 attdef_ranking_europe = attdef[attdef.Team.isin(eurocup_teams)]
2 plot_top_att_def(attdef_ranking_europe,24)

```



In [15]:

1 attdef_ranking_europe

Out[15]:

	date	Team	att	def	def_position	att_position
0	2024-07-05	Spain	2.589579	0.457570	3	1
1	2024-07-05	Germany	2.548609	0.717336	36	2
3	2024-07-02	Netherlands	2.371421	0.618433	18	4
6	2024-07-05	France	2.219333	0.461540	4	7
7	2024-07-05	Portugal	2.216051	0.545238	11	8
8	2024-07-01	Belgium	2.191970	0.533718	9	9
9	2024-06-30	England	2.171419	0.474628	6	10
12	2024-06-24	Croatia	1.892096	0.676163	26	13
16	2024-06-29	Switzerland	1.769522	0.723214	38	17
17	2024-06-29	Italy	1.769313	0.540849	10	18
18	2024-06-25	Serbia	1.740022	0.839128	63	19
25	2024-06-29	Denmark	1.665533	0.602726	17	26
33	2024-06-25	Poland	1.604044	0.828754	59	34
34	2024-07-02	Turkey	1.599998	0.913647	77	35
35	2024-07-02	Austria	1.585274	0.848737	67	36
39	2024-06-26	Czech Republic	1.561265	0.795722	51	40
48	2024-06-26	Ukraine	1.474243	0.737801	41	49
79	2024-06-23	Hungary	1.342773	0.893565	75	80
83	2024-06-30	Slovakia	1.334461	0.842814	64	84
111	2024-07-02	Romania	1.310627	0.792474	49	112
133	2024-06-23	Scotland	1.295363	0.889733	74	134
208	2024-06-30	Georgia	1.187143	1.016745	92	209
210	2024-07-01	Slovenia	1.182072	0.760566	45	211
272	2024-06-24	Albania	0.958692	0.926319	79	273

Based on this ranking, we can conclude that the teams with the best attack are Germany and Spain. The teams with the best defense are England, Spain, and Italy.

4 Goals Expected and Poisson Distribution

The best aspect of this ranking is that we can use it to calculate the expected goals in a match. The expected goals for a team can be calculated as:

$$ExpectedGoals = AttackingPowerTeam \times DefensivePowerOpponent$$

Additionally, in non-neutral matches, we have to set a parameter for the advantage of the home team. This parameter adjusts the expected goals to account for the home field advantage.

Let see if the results adjust with Expected goals we predicted, for doing that we consider the results of European tournaments of the last 14 years.

```
In [16]: 1 #Creating the reference dataset
2 df2 = df[(df.date>"2010-09-01") & (df.tournament.isin(["UEFA Euro","U
3
4 # Define the bins
5 bins = np.array([1.5,2.25, 2.5, 2.75, 3,10]) # adjust this to fit you
6 df2['total_score'] = df2['away_score'] + df2['home_score']
7 df2['total_xG']=df2['XGhome']+df2['XGaway']
8 df2['total_xG_bin'] = pd.cut(df2['total_xG'], bins)
9
10 # Group by the bin and calculate the mean of the 'Goals' column
11 binned_means = df2.groupby('total_xG_bin')['total_score'].agg(["mean"
12
13 print(binned_means)
```

	mean	count
total_xG_bin		
(1.5, 2.25]	2.249322	369
(2.25, 2.5]	2.366234	385
(2.5, 2.75]	2.515406	357
(2.75, 3.0]	2.923767	223
(3.0, 10.0]	3.505988	334

```
In [17]: 1 df2
```

```
Out[17]:
```

	date	home_team	away_team	home_score	away_score	tournament	
33959	2010-09-02	Israel	Malta	3.0	1.0	UEFA Euro qualification	Ram
33961	2010-09-03	Andorra	Russia	0.0	2.0	UEFA Euro qualification	Sant J
33962	2010-09-03	Armenia	Republic of Ireland	0.0	1.0	UEFA Euro qualification	Ye
33963	2010-09-03	Belgium	Germany	0.0	1.0	UEFA Euro qualification	Br
33965	2010-09-03	England	Bulgaria	4.0	0.0	UEFA Euro qualification	L
33966	2010-09-03	Estonia	Italy	1.0	2.0	UEFA Euro qualification	
33967	2010-09-03	Faroe Islands	Serbia	0.0	3.0	UEFA Euro qualification	Tól

```
In [18]: 1 df2.columns
```

```
Out[18]: Index(['date', 'home_team', 'away_team', 'home_score', 'away_score',
'tournament', 'city', 'country', 'neutral', 'Elo_h_after',
'Elo_a_after', 'Elo_h_before', 'Elo_a_before', 'probH', 'probA',
'att_h_after', 'att_a_after', 'att_h_before', 'att_a_before',
'def_h_after', 'def_a_after', 'def_h_before', 'def_a_before', 'XGho
me',
'XGaway', 'total_score', 'total_xG', 'total_xG_bin'],
dtype='object')
```

In [19]:

```
1 df2['year']=df2['date'].apply(lambda x: int(x[:4]))
2 df2
```

Out[19]:

	date	home_team	away_team	home_score	away_score	tournament	
33959	2010-09-02	Israel	Malta	3.0	1.0	UEFA Euro qualification	Ram
33961	2010-09-03	Andorra	Russia	0.0	2.0	UEFA Euro qualification	Sant J
33962	2010-09-03	Armenia	Republic of Ireland	0.0	1.0	UEFA Euro qualification	Ye
33963	2010-09-03	Belgium	Germany	0.0	1.0	UEFA Euro qualification	Br
33965	2010-09-03	England	Bulgaria	4.0	0.0	UEFA Euro qualification	L
33966	2010-09-03	Estonia	Italy	1.0	2.0	UEFA Euro qualification	
33967	2010-09-03	Faroe Islands	Serbia	0.0	3.0	UEFA Euro qualification	Tól

In [20]:

1 df2.info()

```

<class 'pandas.core.frame.DataFrame'>
Index: 1668 entries, 33959 to 47368
Data columns (total 29 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   date                  1668 non-null   object
 1   home_team             1668 non-null   object
 2   away_team             1668 non-null   object
 3   home_score            1668 non-null   float64
 4   away_score            1668 non-null   float64
 5   tournament            1668 non-null   object
 6   city                  1668 non-null   object
 7   country               1668 non-null   object
 8   neutral               1668 non-null   bool
 9   Elo_h_after           1668 non-null   float64
10   Elo_a_after           1668 non-null   float64
11   Elo_h_before          1668 non-null   float64
12   Elo_a_before          1668 non-null   float64
13   probH                 1668 non-null   float64
14   probA                 1668 non-null   float64
15   att_h_after           1668 non-null   float64
16   att_a_after           1668 non-null   float64
17   att_h_before          1668 non-null   float64
18   att_a_before          1668 non-null   float64
19   def_h_after           1668 non-null   float64
20   def_a_after           1668 non-null   float64
21   def_h_before          1668 non-null   float64
22   def_a_before          1668 non-null   float64
23   XGhome               1668 non-null   float64
24   XGaway               1668 non-null   float64
25   total_score           1668 non-null   float64
26   total_xG              1668 non-null   float64
27   total_xG_bin          1668 non-null   category
28   year                 1668 non-null   int64
dtypes: bool(1), category(1), float64(20), int64(1), object(6)
memory usage: 368.4+ KB

```



```
In [21]: 1 df2['home_diff_elo']=df2['Elo_h_after']-df2['Elo_a_after']
        2 df2['away_diff_elo']=df2['Elo_a_after']-df2['Elo_h_after']
        3 df2
```

Out[21]:

	date	home_team	away_team	home_score	away_score	tournament	
33959	2010-09-02	Israel	Malta	3.0	1.0	UEFA Euro qualification	Ram
33961	2010-09-03	Andorra	Russia	0.0	2.0	UEFA Euro qualification	Sant J
33962	2010-09-03	Armenia	Republic of Ireland	0.0	1.0	UEFA Euro qualification	Ye
33963	2010-09-03	Belgium	Germany	0.0	1.0	UEFA Euro qualification	Br
33965	2010-09-03	England	Bulgaria	4.0	0.0	UEFA Euro qualification	L
33966	2010-09-03	Estonia	Italy	1.0	2.0	UEFA Euro qualification	
33967	2010-09-03	Faroe Islands	Serbia	0.0	3.0	UEFA Euro qualification	Tól

```

In [22]: 1 aP_teams = {}
        2 for idx, match in df2.iterrows():
        3     home_team = match["home_team"]
        4     away_team = match["away_team"]
        5     home_score = match["home_score"]
        6     away_score = match["away_score"]
        7
        8     match_year = int(match['date'][:4])
        9     if match_year == 2024:
       10         match_year = 2023
       11
       12     #     if home_team == 'Northern Ireland':
       13     #         home_team = 'Republic of Ireland'
       14
       15     #     if away_team == 'Northern Ireland':
       16     #         away_team = 'Republic of Ireland'
       17
       18     home_elo_series = df2[(df2['year'] == match_year) & (df2['home_team'] == home_team)]
       19     away_elo_series = df2[(df2['year'] == match_year) & (df2['away_team'] == away_team)]
       20
       21     try:
       22         home_elo = home_elo_series.tolist()[0]
       23         away_elo = away_elo_series.tolist()[0]
       24     except:
       25         continue
       26
       27     home_diff = home_elo - away_elo
       28     away_diff = away_elo - home_elo
       29
       30     if home_team in eurocup_teams:
       31         if home_team not in aP_teams:
       32             aP_teams[home_team] = [home_score, 1, [home_score], away_score, home_diff]
       33         else:
       34             aP_teams[home_team][0] += home_score
       35             aP_teams[home_team][1] += 1
       36             aP_teams[home_team][2].append(home_score)
       37             aP_teams[home_team][3] += away_score
       38             aP_teams[home_team][4].append(away_score)
       39             aP_teams[home_team][5].append(home_diff)
       40
       41     if away_team in eurocup_teams:
       42         if away_team not in aP_teams:
       43             aP_teams[away_team] = [away_score, 1, [away_score], home_score, away_diff]
       44         else:
       45             aP_teams[away_team][0] += away_score
       46             aP_teams[away_team][1] += 1
       47             aP_teams[away_team][2].append(away_score)
       48             aP_teams[away_team][3] += home_score
       49             aP_teams[away_team][4].append(home_score)
       50             aP_teams[away_team][5].append(away_diff)

```

In [23]:

```

1 # Summary of aP and Code Function
2 # aP: This dictionary ends up containing the average number of goals s
3 # Print Output: For each team, it prints the team name, the average go
4 import statistics
5 aP = {}
6 for i in aP_teams.keys():
7     aP[i] = aP_teams[i][0] / aP_teams[i][1]
8
9     print(f"{i},{aP[i]},{statistics.stdev(aP_teams[i][2]))}")

```

```

Belgium,2.3714285714285714,1.7870016063789165
Germany,2.288135593220339,1.848119710423572
England,2.0985915492957745,1.8059815659355425
Italy,1.7625,1.4861704248320815
Serbia,1.5789473684210527,1.223977246827685
France,1.8333333333333333,1.9890727127006333
Georgia,1.2580645161290323,1.213793455719087
Turkey,1.4696969696969697,1.3268783141667742
Croatia,1.6428571428571428,1.3623904214488174
Spain,2.3947368421052633,1.7669480080517068
Scotland,1.576271186440678,1.499464155215507
Portugal,2.0666666666666667,1.833128571693857
Romania,1.303030303030303,1.2021736335577509
Albania,1.1071428571428572,1.1548879745526117
Netherlands,2.208955223880597,1.8630157486674404
Slovakia,1.1940298507462686,1.1446688811967767
Slovenia,1.359375,1.337697470958783
Hungary,1.492537313432836,1.3070106341683057
Austria,1.6153846153846154,1.3996908999434894
Czech Republic,1.373134328358209,1.056521176708736
Denmark,1.671875,1.4148272371958037
Switzerland,1.878787878787879,1.6503160677709963
Poland,1.5833333333333333,1.5762180296000163
Ukraine,1.3166666666666667,1.2418093233182887

```

```
In [24]: 1 # Structure of aP_teams
2 # Key: Team name.
3 # Value: A list containing:
4 # Total goals scored by the team ([0]).
5 # Total number of matches ([1]).
6 # A list of goals scored in each match ([2]).
7 # Total goals conceded ([3]).
8 # A list of goals conceded in each match ([4]).
9 # A list of ELO rating differences for each match ([5]).
10
11
12 import pprint # Will be used later for some nicer prints
13 pprint.pprint(aP_teams["England"])
```

```
[149.0,
 71,
 [4.0,
  3.0,
  0.0,
  2.0,
  2.0,
  3.0,
  1.0,
  2.0,
  1.0,
  3.0,
  1.0,
  0.0,
  2.0,
  5.0,
  1.0,
  3.0,
  4.0,
  ~ ~]
```

```
In [25]: 1 mega = [[], []]
2 for idx, i in enumerate(list(aP_teams.values())):
3     mega[1] += (np.array(i[2]) - np.array(len(i[2]))*[aP[list(aP_teams
4     mega[0] += i[5]
```

```
In [26]: 1 from sklearn.linear_model import LinearRegression
2
3 X, y = np.array(mega[0]).reshape(-1, 1), np.array(mega[1]).reshape(-1, 1)
4
5
6 reg = LinearRegression()
7 reg.fit(X, y)
8
9
10 print(reg.score(X, y))
11
12 print(reg.coef_, reg.intercept_)
```

```
0.1575846986689724
[[0.00279906]] [-0.28575836]
```

```
In [27]: 1 max(X)
```

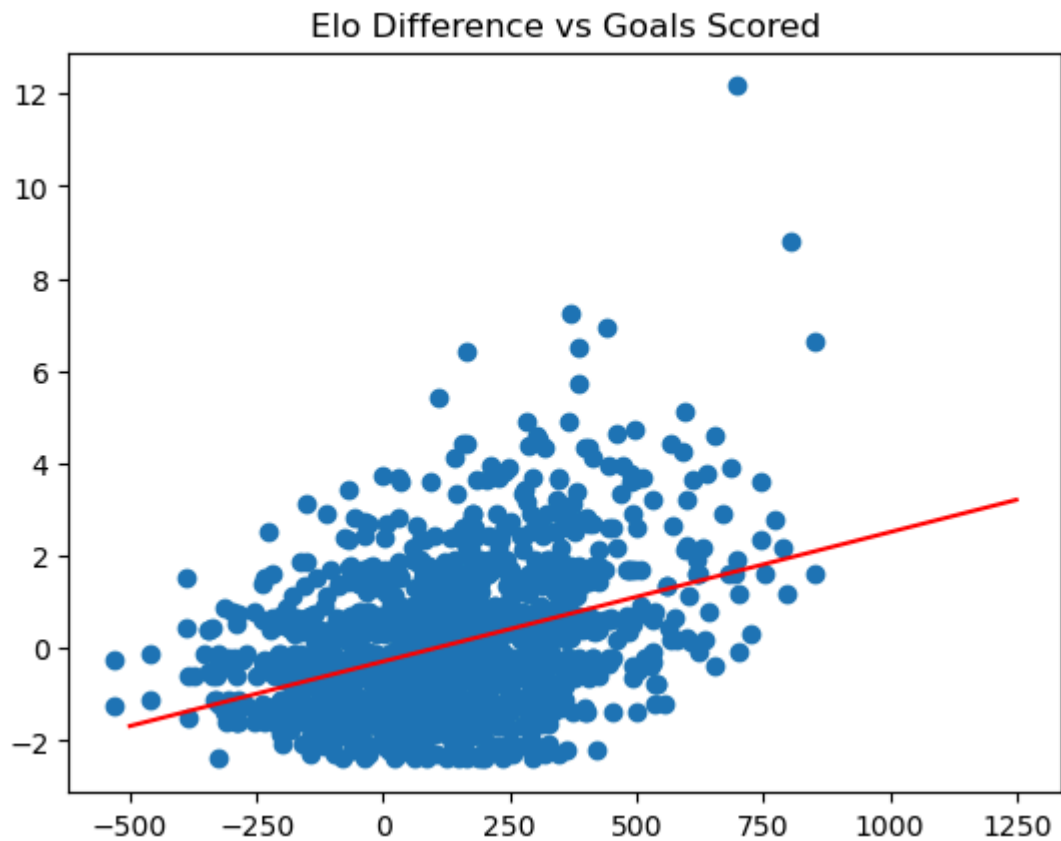
```
Out[27]: array([852.45])
```

In [28]: 1 `min(X)`

Out[28]: `array([-531.88])`

In [29]: 1 `plt.title('Elo Difference vs Goals Scored')`
 2 `plt.scatter(mega[0], mega[1])`
 3 `x = np.arange(-500, 1250)`
 4 `y = x*reg.coef_[0] + reg.intercept_[0]`
 5 `plt.plot(x, y, color='red')`

Out[29]: `[<matplotlib.lines.Line2D at 0x2711fbe0610>]`



We can see that the expected goals we created adjust very well to the real number of goals in the matches. So we can use the poisson distribution to predict the results of the matches

The Poisson distribution is defined by the probability mass function:

$$P(X=k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

where:

$P(X=k)$ is the probability of k goals being scored, λ is the average number of goals scored in a match.

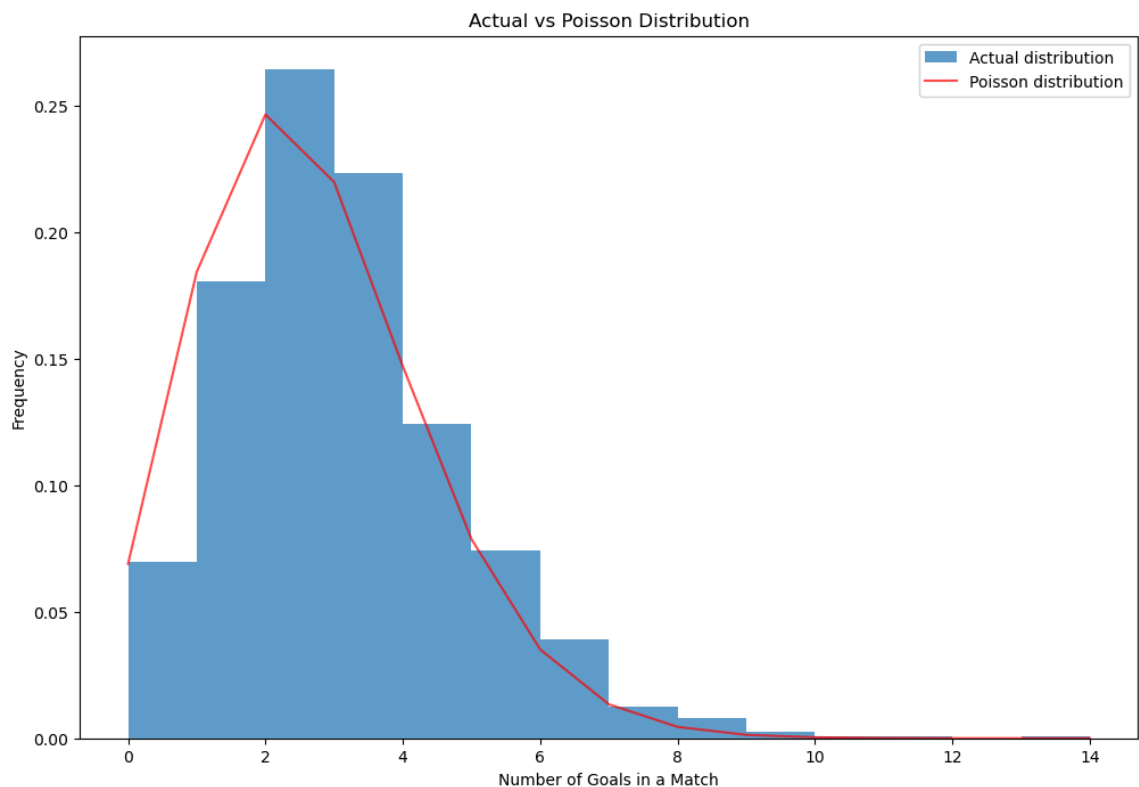
In our context, λ represents the average number of goals a team is expected to score.

In [30]:

```

1 data = df2['total_score']
2
3 # Calculate the mean of the data as it will be the Lambda (rate parameter)
4 mu = data.mean()
5
6 # Create a range of numbers from 0 to the maximum number of goals in a match
7 k = np.arange(0, data.max()+1)
8
9 # Create a Poisson distribution with the mean obtained
10 poisson_pmf = poisson.pmf(k, mu)
11
12 plt.figure(figsize=(12,8))
13
14 # Plot the actual distribution of goals
15 plt.hist(data, bins=k, density=True, alpha=0.7, label='Actual distribution')
16
17 # Plot the Poisson distribution
18 plt.plot(k, poisson_pmf, 'r-', alpha=0.7, label='Poisson distribution')
19
20 plt.title('Actual vs Poisson Distribution')
21 plt.xlabel('Number of Goals in a Match')
22 plt.ylabel('Frequency')
23 plt.legend()
24
25 plt.show()

```



With the utilization of the given distribution, we can use the probabilities to predict the matches.

```

In [31]: 1 def calculate_match_probabilities(team_a_lambda, team_b_lambda, max_g
2         # Create a matrix of zeros
3         matrix = np.zeros((max_goals + 1, max_goals + 1))
4
5         # Populate the matrix with probabilities
6         for i in range(max_goals + 1):
7             for j in range(max_goals + 1):
8                 matrix[i, j] = poisson.pmf(i, team_a_lambda) * poisson.pmf(j, team_b_lambda)
9
10        # Calculate the total goals probabilities
11        total_goals_prob = [np.sum(np.diag(matrix[:, -1], k)) for k in range(-1, max_goals + 1)]
12
13        # Calculate the probability for the specific result
14        draw_prob = np.sum(np.diag(matrix))
15        away_prob = np.sum(np.triu(matrix, 1))
16        local_prob = np.sum(np.tril(matrix, -1))
17
18        df_match_goals = pd.DataFrame(matrix)
19        highest_value = df_match_goals.max().max()
20        max_position = df_match_goals.stack().idxmax()
21        h_goals = max_position[0]
22        a_goals = max_position[1]
23
24
25        return total_goals_prob, draw_prob, away_prob, local_prob, h_goals, a_goals
26

```

For instance, let's consider the inaugural match of the 2024 Euro between Germany and Scotland. Germany boasts an offensive power of 2.58, while Scotland 1.31.

```

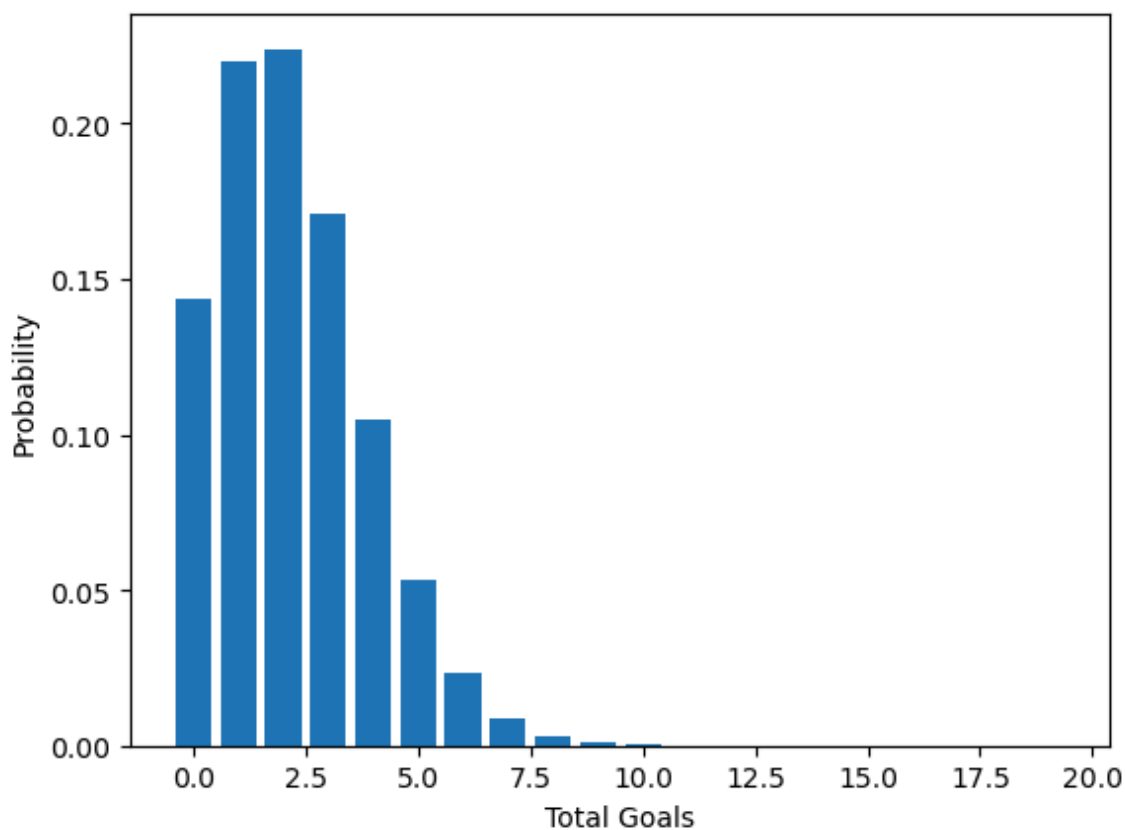
In [32]: 1 # Example usage
2 # Albania 1342.140
3 # France 1752.570
4 # Georgia 1.177952 0.987056
5 # Spain 2.548556 0.458797
6
7 team_a_lambda = 2.548556 * 0.987056
8 team_b_lambda = 1.177952 * 0.458797
9 total_goals_prob, draw_prob, away_prob, local_prob, h_goals, a_goals =
10
11 # Print the results
12 '''print("Total goals probabilities:")
13 for goals, prob in enumerate(total_goals_prob):
14     print(f'Probability of total {goals} goals: {prob}''')
15
16
17 print(f"Probability for France: {round(local_prob,2)*100} %")
18 print(f"Probability for Albania: {round(away_prob,2)*100} %")
19 print(f"Probability for draw: {round(draw_prob,2)*100} %")
20
21 import matplotlib.pyplot as plt
22 # Plot the distribution of possible total goals
23 plt.bar(range(len(total_goals_prob)), total_goals_prob)
24 plt.xlabel('Total Goals')
25 plt.ylabel('Probability')
26 plt.show()

```

Probability for France: 80.0 %

Probability for Albania: 6.0 %

Probability for draw: 14.000000000000002 %



In [33]: 1 h_goals

Out[33]: 2

In [34]: 1 a_goals

Out[34]: 0

In [35]: 1 max(total_goals_prob)

Out[35]: 0.22392614125034904

5 Adjusted Goals Expected with ELO form and Poisson Distribution

In [36]: 1 # phase1=pd.read_excel('Prediction_template.xlsx',sheet_name='Match pr
2 # phase1=phase1[['Home Team','Away Team','Home_goal ','Away_Goal ']]
3 # phase1=phase1.iloc[:36,:]
4 # phase1

In [37]: 1 # home_goals = phase1.groupby('Home Team').agg(goals_scored=('Home_goa
2 # home_goals.rename(columns={'Home Team': 'Team'}, inplace=True)
3 # away_goals = phase1.groupby('Away Team').agg(goals_scored=('Away_Goa
4 # away_goals.rename(columns={'Away Team': 'Team'}, inplace=True)
5 # teams_pred = pd.concat([home_goals, away_goals]).groupby('Team').sum
6 # teams_pred

In [38]: 1 # teams_pred.columns

In []: 1

5.1 Adding euro form

In [39]:

```
1 import pandas as pd
2 euros24_group=pd.read_excel('euroresults.xlsx')
3 euros24_group
```

Out[39]:

	Team	goals_scored	goals_against	points
0	Spain	5	0	9
1	Germany	8	2	7
2	Austria	6	4	6
3	Turkey	5	5	6
4	Portugal	5	3	6
5	England	2	1	5
6	France	2	1	5
7	Switzerland	5	3	5
8	Italy	3	3	4
9	Romania	4	3	4
10	Belgium	2	1	4
11	Netherlands	4	4	4
12	Slovakia	3	3	4
13	Georgia	4	4	4
14	Ukraine	2	4	4
15	Hungary	2	5	3
16	Slovenia	2	2	3
17	Denmark	2	2	3
18	Serbia	1	2	2
19	Croatia	3	6	2
20	Poland	3	6	1
21	Scotland	2	7	1
22	Albania	3	5	1
23	Czech Republic	3	5	1

```

In [40]: 1 group16=[
          2     'Austria',
          3     'England',
          4     'Germany',
          5     'Portugal',
          6     'Romania',
          7     'Spain',
          8     'Belgium',
          9     'Denmark',
         10     'France',
         11     'Italy',
         12     'Switzerland',
         13     'Turkey',
         14     'Georgia',
         15     'Netherlands',
         16     'Slovakia',
         17     'Slovenia']
         18
         19 group16_euroteams=euros24_group[euros24_group['Team'].isin(group16)]
         20 group16_euroteams_sorted = group16_euroteams.sort_values(by=['points',
         21 group16_euroteams_sorted['euro_form'] = (group16_euroteams_sorted['points']
         22 group16_euroteams_sorted
         23

```

Out[40]:

	Team	goals_scored	goals_against	points	euro_form
0	Spain	5	0	9	2.000000
1	Germany	8	2	7	1.555556
2	Austria	6	4	6	1.333333
3	Turkey	5	5	6	1.333333
4	Portugal	5	3	6	1.333333
5	England	2	1	5	1.111111
6	France	2	1	5	1.111111
7	Switzerland	5	3	5	1.111111
8	Italy	3	3	4	0.888889
9	Romania	4	3	4	0.888889
10	Belgium	2	1	4	0.888889
11	Netherlands	4	4	4	0.888889
12	Slovakia	3	3	4	0.888889
13	Georgia	4	4	4	0.888889
16	Slovenia	2	2	3	0.666667
17	Denmark	2	2	3	0.666667

```

In [41]: 1 group16_euroteams_sorted['points'].median()

```

Out[41]: 4.5

In []:

1

In [42]:

```
1 # group16_euroteams_sorted_predict = pd.merge(group16_euroteams_sort
2 # group16_euroteams_sorted_predict['diff_goals_scored']=group16_euro
3 # group16_euroteams_sorted_predict['diff_goals_agaisnt']=group16_euro
4 # group16_euroteams_sorted_predict['avg_goals_scored']=group16_eurost
5 # group16_euroteams_sorted_predict['avg_goals_agaisnt']=group16_eurost
6 # group16_euroteams_sorted_predict['euro_form']=group16_euroteams_sc
7 # group16_euroteams_sorted_predict['form_goals_scored']=group16_eurost
8 # group16_euroteams_sorted_predict['form_goals_agaisnt']=group16_euro
9 # group16_euroteams_sorted_predict
10
```

In []:

1

In []:

1

In []:

1

```

In [43]: 1 def calculate_match_probabilities_elo(team_a_lambda, team_b_lambda, home_advantage):
2         # Create a matrix of zeros
3         matrix = np.zeros((max_goals + 1, max_goals + 1))
4
5         # Populate the matrix with probabilities
6         for i in range(max_goals + 1):
7             for j in range(max_goals + 1):
8                 matrix[i, j] = poisson.pmf(i, team_a_lambda) * poisson.pmf(j, team_b_lambda)
9
10        # Calculate the total goals probabilities
11        total_goals_prob = [np.sum(np.diag(matrix[:, -1], k)) for k in range(-1, max_goals + 1)]
12
13        # Calculate the probability for the specific result
14        draw_prob = np.sum(np.diag(matrix))
15        away_prob = np.sum(np.triu(matrix, 1))
16        local_prob = np.sum(np.tril(matrix, -1))
17
18        df_match_goals = pd.DataFrame(matrix)
19        highest_value = df_match_goals.max().max()
20        max_position = df_match_goals.stack().idxmax()
21        #array results so need to convert to int
22        h_goals_adjust = int(reg.predict(np.array(home_diff).reshape(-1, 1)))
23        a_goals_adjust = int(reg.predict(np.array(away_diff).reshape(-1, 1)))
24        h_goals = max_position[0]
25        a_goals = max_position[1]
26
27        fin_h_goals = h_goals + h_goals_adjust
28        fin_a_goals = a_goals + a_goals_adjust
29        if (fin_a_goals == -1):
30            adj_sum = abs(h_goals_adjust) + abs(a_goals_adjust)
31            fin_h_goals += adj_sum
32            fin_a_goals += 1
33        elif (fin_h_goals == -1):
34            adj_sum = abs(h_goals_adjust) + abs(a_goals_adjust)
35            fin_a_goals += adj_sum
36            fin_h_goals += 1
37
38
39        return total_goals_prob, draw_prob, away_prob, local_prob, h_goals, a_goals

```

```

In [44]: 1 float(elos_ranking_euro[elos_ranking_euro['Team']=='Germany']['Elo'])

```

Out[44]: 1643.99000000000023

```

In [45]: 1
2         # home_diff = float(elos_ranking_euro[elos_ranking_euro['Team']=='Germany']['home_diff'])
3         # away_diff = float(elos_ranking_euro[elos_ranking_euro['Team']=='Scotland']['away_diff'])
4         # print(f'home_diff : {home_diff}')
5         # print(f'away_diff : {away_diff}')

```

```

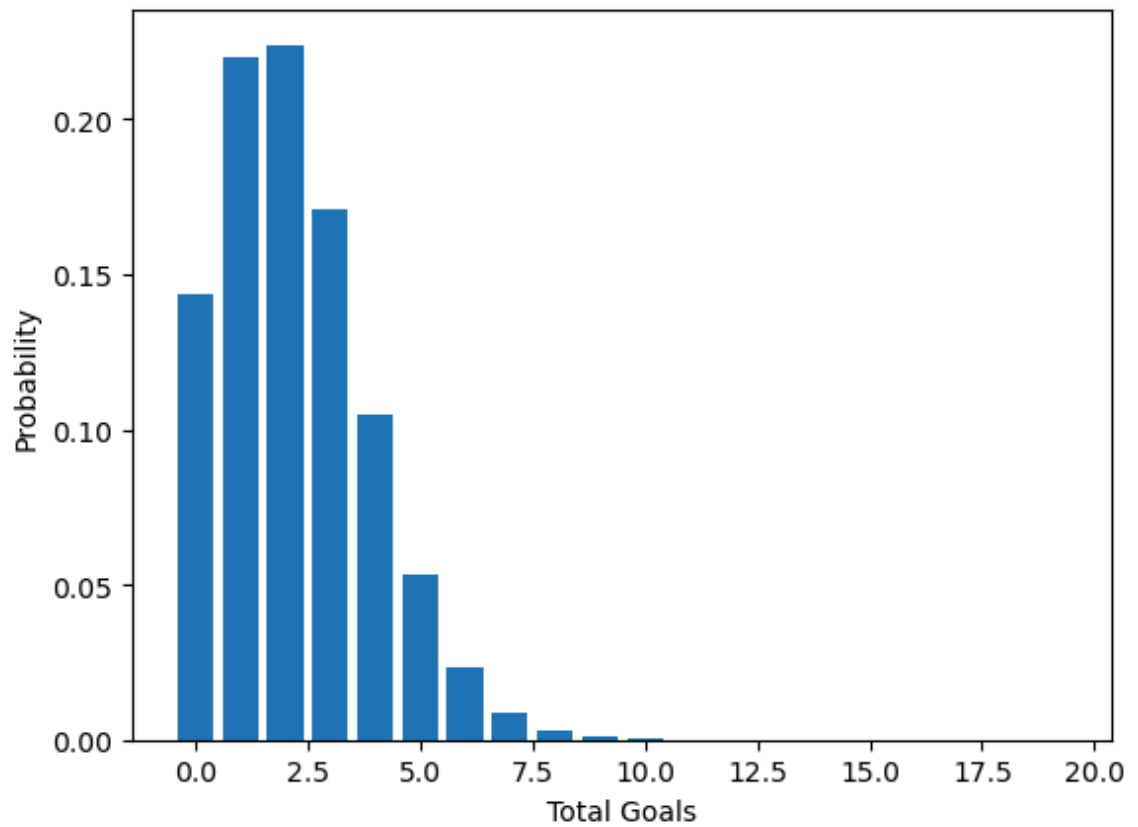
In [46]: 1 # Example usage
2 # Albania 1342.140
3 # France 1752.570
4 # Georgia 1.177952 0.987056
5 # Spain 2.548556 0.458797
6
7 team_a_lambda = 2.548556 * 0.987056
8 team_b_lambda = 1.177952 * 0.458797
9 home_diff = float(elos_ranking_euro[elos_ranking_euro['Team']=='Spain']
10 away_diff = float(elos_ranking_euro[elos_ranking_euro['Team']=='Georgia']
11 total_goals_prob, draw_prob, away_prob, local_prob, h_goals, a_goals, fi
12
13
14
15 # Print the results
16 '''print("Total goals probabilities:")
17 for goals, prob in enumerate(total_goals_prob):
18     print(f'Probability of total {goals} goals: {prob}''')
19
20
21 print(f"Probability for France: {round(local_prob,2)*100} %")
22 print(f"Probability for Albania: {round(away_prob,2)*100} %")
23 print(f"Probability for draw: {round(draw_prob,2)*100} %")
24
25 import matplotlib.pyplot as plt
26 # Plot the distribution of possible total goals
27 plt.bar(range(len(total_goals_prob)), total_goals_prob)
28 plt.xlabel('Total Goals')
29 plt.ylabel('Probability')
30 plt.show()

```

Probability for France: 80.0 %

Probability for Albania: 6.0 %

Probability for draw: 14.000000000000002 %



In [47]: 1 h_goals,a_goals

Out[47]: (2, 0)

In [48]: 1 fin_h_goals,fin_a_goals

Out[48]: (3, 0)

```

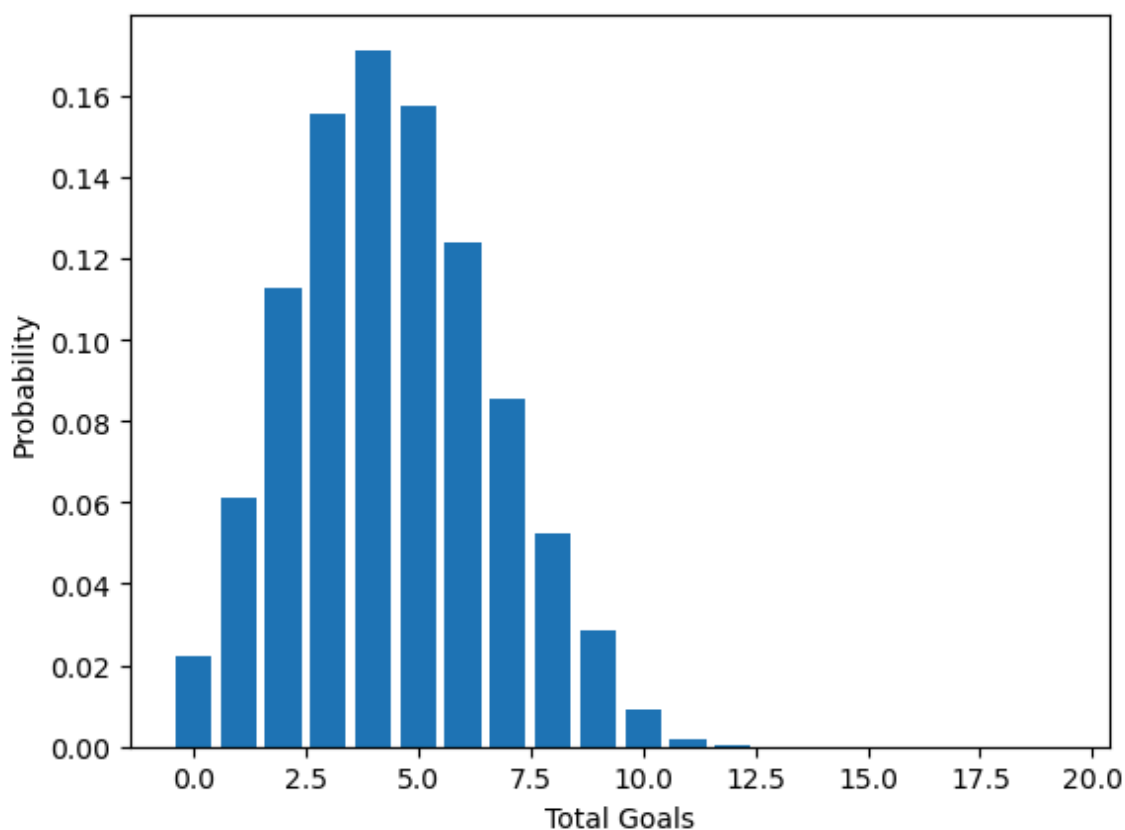
In [49]: 1 # Example usage
2 # Georgia 1389.450
3 # Spain 1713.675
4 # Georgia 1.177952 0.987056
5 # Spain 2.548556 0.458797
6
7 team_a_lambda = 2.548556 * 0.987056*2.000000
8 team_b_lambda = 1.177952 * 0.458797*0.888889
9 home_diff = float(elos_ranking_euro[elos_ranking_euro['Team']=='Spain
10 away_diff = float(elos_ranking_euro[elos_ranking_euro['Team']=='Georg
11 total_goals_prob, draw_prob, away_prob, local_prob,h_goals,a_goals,fi
12
13
14 # Print the results
15 '''print("Total goals probabilities:")
16 for goals, prob in enumerate(total_goals_prob):
17     print(f'Probability of total {goals} goals: {prob}''')
18
19
20 print(f"Probability for France: {round(local_prob,2)*100} %")
21 print(f"Probability for Albania: {round(away_prob,2)*100} %")
22 print(f"Probability for draw: {round(draw_prob,2)*100} %")
23
24 import matplotlib.pyplot as plt
25 # Plot the distribution of possible total goals
26 plt.bar(range(len(total_goals_prob)), total_goals_prob)
27 plt.xlabel('Total Goals')
28 plt.ylabel('Probability')
29 plt.show()

```

Probability for France: 96.0 %

Probability for Albania: 1.0 %

Probability for draw: 2.0 %



In [50]: 1 h_goals,a_goals

Out[50]: (5, 0)

In [51]: 1 fin_h_goals,fin_a_goals

Out[51]: (6, 0)

In []: 1

In []: 1

In [52]: 1 for x in range(-532,853):
2 h_goals_adjust=int(reg.predict(np.array(x).reshape(-1, 1)))
3 print(h_goals_adjust,x)

-1 -532

-1 -531

-1 -530

-1 -529

-1 -528

-1 -527

-1 -526

-1 -525

-1 -524

-1 -523

-1 -522

-1 -521

-1 -520

-1 -519

-1 -518

-1 -517

-1 -516

-1 -515

-1 -514

-1 -513

```

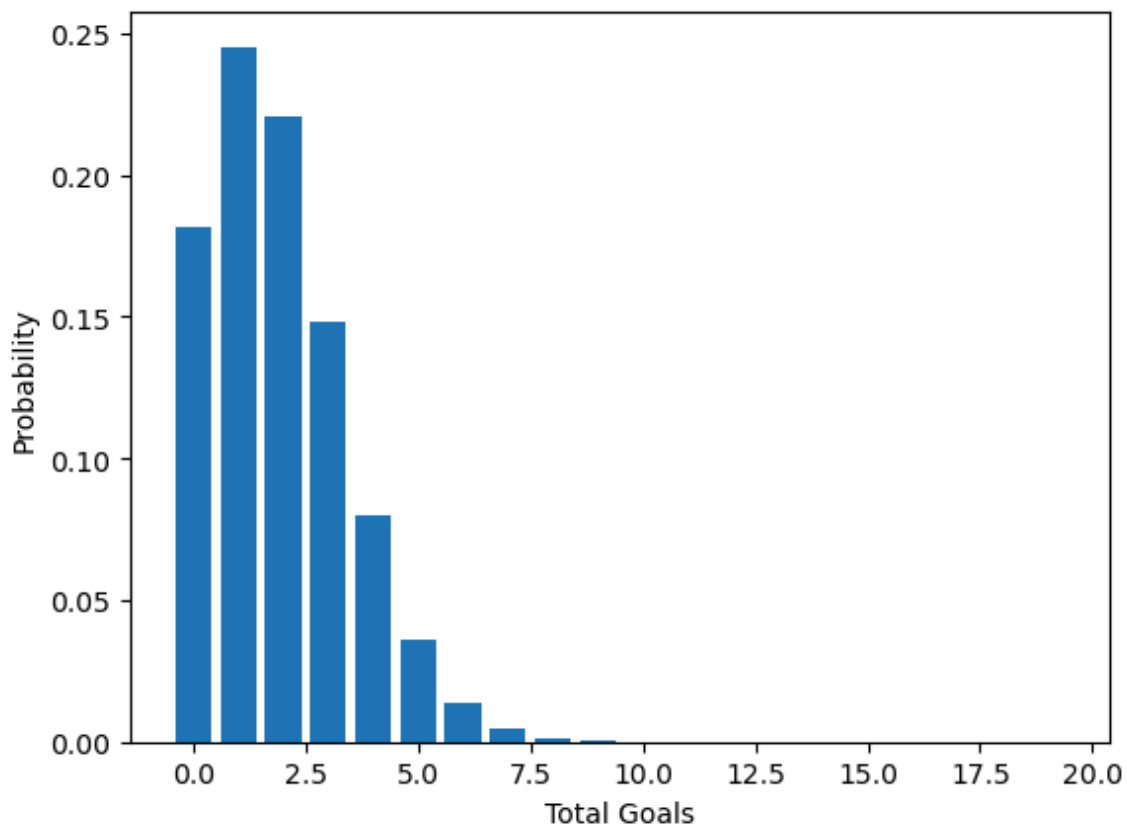
In [53]: 1 # Example usage
2 # Albania 1342.140
3 # France 1752.570
4 # Albania 0.929605 0.945796
5 # France 2.316062 0.543255
6
7
8 team_b_lambda = 2.316062 * 0.945796
9 team_a_lambda = 0.543255 * 0.929605
10 away_diff = float(elos_ranking_euro[elos_ranking_euro['Team']=='France'])
11 home_diff = float(elos_ranking_euro[elos_ranking_euro['Team']=='Albania'])
12 total_goals_prob, draw_prob, away_prob, local_prob, h_goals, a_goals, fi
13
14
15
16 # Print the results
17 '''print("Total goals probabilities:")
18 for goals, prob in enumerate(total_goals_prob):
19     print(f'Probability of total {goals} goals: {prob}''')
20
21
22 print(f"Probability for France: {round(local_prob,2)*100} %")
23 print(f"Probability for Albania: {round(away_prob,2)*100} %")
24 print(f"Probability for draw: {round(draw_prob,2)*100} %")
25
26 import matplotlib.pyplot as plt
27 # Plot the distribution of possible total goals
28 plt.bar(range(len(total_goals_prob)), total_goals_prob)
29 plt.xlabel('Total Goals')
30 plt.ylabel('Probability')
31 plt.show()

```

Probability for France: 7.000000000000001 %

Probability for Albania: 76.0 %

Probability for draw: 17.0 %



In [54]: 1 h_goals,a_goals

Out[54]: (0, 2)

In [55]: 1 fin_h_goals,fin_a_goals

Out[55]: (0, 3)

In [56]: 1 team_a_lambda = 2.316062 * 0.945796
 2 team_b_lambda = 0.543255 * 0.929605
 3 total_goals_prob, draw_prob, away_prob, local_prob,h_goals,a_goals,fi
 4

In [57]: 1 h_goals,a_goals

Out[57]: (2, 0)

In [58]: 1 fin_h_goals,fin_a_goals

Out[58]: (5, 0)

In [59]: 1 team_a_lambda = 2.316062 * 0.945796
 2 team_b_lambda = 0.543255 * 0.929605
 3 total_goals_prob, draw_prob, away_prob, local_prob,h_goals,a_goals,fi
 4

In [60]: 1 h_goals,a_goals

Out[60]: (0, 2)

In [61]: 1 fin_h_goals,fin_a_goals

Out[61]: (0, 5)

Now it is time to validate if the probabilities given by the poisson distribution are accurate or not

In [62]: 1 df2.info()

```
<class 'pandas.core.frame.DataFrame'>
Index: 1668 entries, 33959 to 47368
Data columns (total 31 columns):
#   Column                Non-Null Count  Dtype
---  -
0   date                  1668 non-null   object
1   home_team             1668 non-null   object
2   away_team             1668 non-null   object
3   home_score            1668 non-null   float64
4   away_score            1668 non-null   float64
5   tournament            1668 non-null   object
6   city                  1668 non-null   object
7   country               1668 non-null   object
8   neutral               1668 non-null   bool
9   Elo_h_after           1668 non-null   float64
10  Elo_a_after           1668 non-null   float64
11  Elo_h_before          1668 non-null   float64
12  Elo_a_before          1668 non-null   float64
13  probH                 1668 non-null   float64
14  probA                 1668 non-null   float64
15  att_h_after           1668 non-null   float64
16  att_a_after           1668 non-null   float64
17  att_h_before          1668 non-null   float64
18  att_a_before          1668 non-null   float64
19  def_h_after           1668 non-null   float64
20  def_a_after           1668 non-null   float64
21  def_h_before          1668 non-null   float64
22  def_a_before          1668 non-null   float64
23  XGhome                1668 non-null   float64
24  XGaway                1668 non-null   float64
25  total_score           1668 non-null   float64
26  total_xG              1668 non-null   float64
27  total_xG_bin          1668 non-null   category
28  year                  1668 non-null   int64
29  home_diff_elo         1668 non-null   float64
30  away_diff_elo         1668 non-null   float64
dtypes: bool(1), category(1), float64(22), int64(1), object(6)
memory usage: 394.4+ KB
```

```

In [63]: 1 df_euro=df2.copy()
          2 # Function to apply to each row of the DataFrame
          3 def apply_calculation_elo(row):
          4 #     home_diff = float(elos_ranking_euro[elos_ranking_euro['Team']==
          5 #     away_diff = float(elos_ranking_euro[elos_ranking_euro['Team']==
          6     total_goals_prob, draw_prob, away_prob, local_prob,h_goals,a_goals
          7     return pd.Series({
          8         'total_goals_prob': total_goals_prob,
          9         'draw_prob': draw_prob,
         10         'away_prob': away_prob,
         11         'local_prob': local_prob,
         12         'home_goals': h_goals,
         13         'away_goals': a_goals,
         14         'home_goals_elo': fin_h_goals,
         15         'away_goals_elo': fin_a_goals,
         16     })
         17
         18 # Apply the function to each row
         19 calculated_variables = df_euro.apply(apply_calculation_elo, axis=1)
         20
         21 # Concatenate the result with the original DataFrame
         22 matches_df = pd.concat([df_euro, calculated_variables], axis=1)
         23
         24 df_euro[['total_goals_prob', 'draw_prob', 'away_prob', 'local_prob', 'l
         25
         26 # Iterate over each row in the DataFrame
         27 for index, row in df_euro.iterrows():
         28     # Iterate over each element in the 'total_goals_prob' list
         29     for i, prob in enumerate(row['total_goals_prob']):
         30         # Create new column with the name 'prob_index_goals_i' and ass
         31         df_euro.at[index, f'prob_index_goals_{i}'] = prob

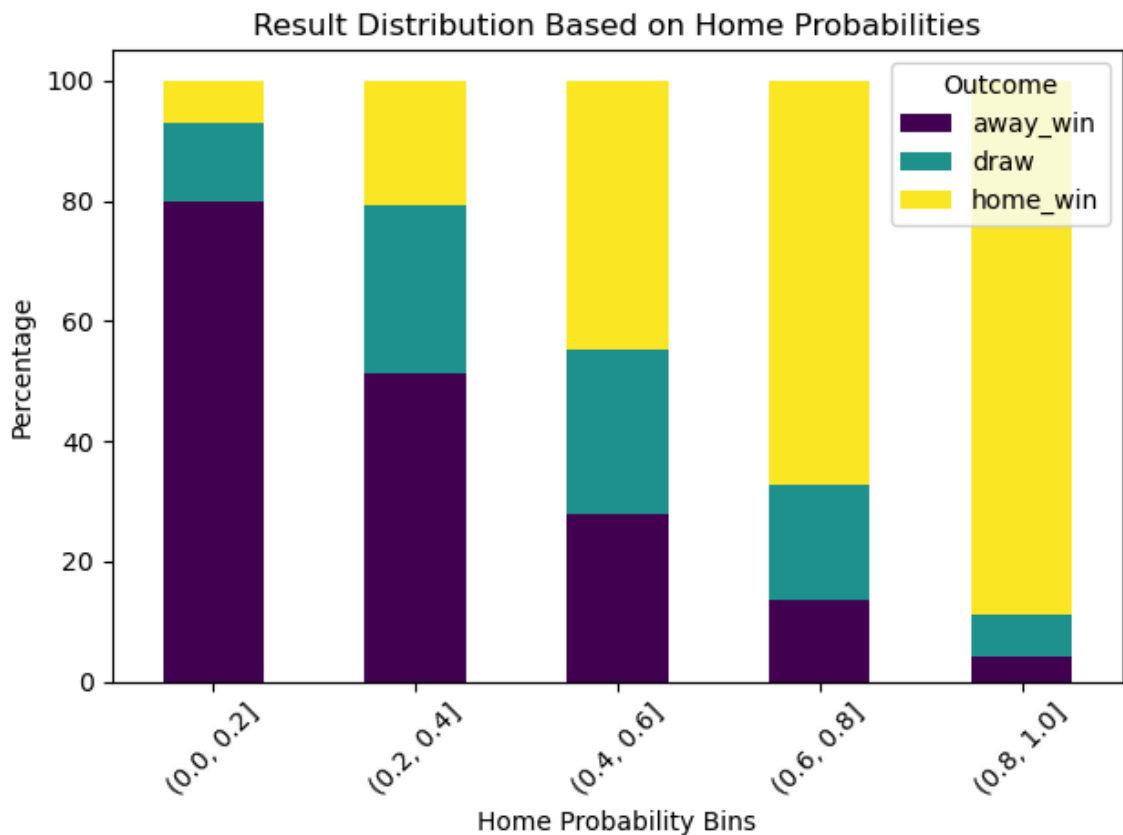
```

```

In [64]: 1 # Compute the actual outcome based on the home and away scores
          2 df_euro['result'] = np.where(df_euro['home_score'] > df_euro['away_score'], 'home_win',
          3                                     np.where(df_euro['home_score'] == df_euro['away_score'], 'draw', 'away_win'))
          4
          5 df_euro['home_win_flag'] = np.where(df_euro['result']=='home_win', 1, 0)
          6 df_euro['draw_flag'] = np.where(df_euro['result']=='draw', 1, 0)
          7 df_euro['away_win_flag'] = np.where(df_euro['result']=='away_win', 1, 0)
          8 # Define bins for grouping based on predicted probabilities (adjust as needed)
          9 bins = [0, 0.2, 0.4, 0.6, 0.8, 1]
         10
         11 # Create a new column indicating the bin for each match based on predicted probabilities
         12 df_euro['probability_bin'] = pd.cut(df_euro['local_prob'], bins=bins)
         13
         14 # Group the DataFrame by the probability bins and compute the percentage distribution of results
         15 result_distribution = df_euro.groupby('probability_bin')['result'].value_counts().unstack()
         16
         17 # Plotting the result distribution based on home probabilities
         18 plt.figure(figsize=(12, 6))
         19 result_distribution.plot(kind='bar', stacked=True, cmap='viridis')
         20 plt.title('Result Distribution Based on Home Probabilities')
         21 plt.xlabel('Home Probability Bins')
         22 plt.ylabel('Percentage')
         23 plt.xticks(rotation=45)
         24 plt.legend(title='Outcome')
         25 plt.tight_layout()

```

<Figure size 1200x600 with 0 Axes>

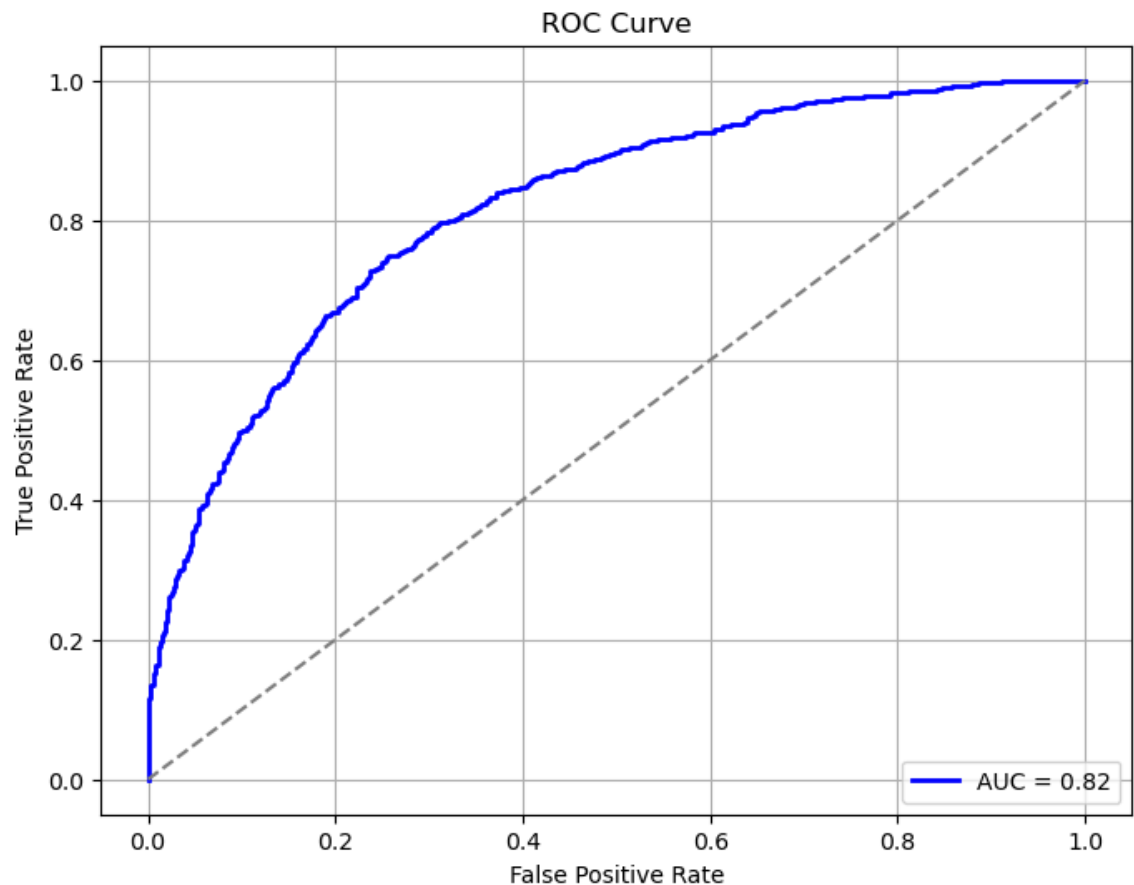


In [65]:

```

1  from sklearn.metrics import roc_curve, roc_auc_score
2  import matplotlib.pyplot as plt
3
4  # Compute ROC curve
5  fpr, tpr, thresholds = roc_curve(df_euro['home_win_flag'], df_euro['local_prob'])
6
7  # Compute AUC
8  auc = roc_auc_score(df_euro['home_win_flag'], df_euro['local_prob'])
9
10 # Plot ROC curve
11 plt.figure(figsize=(8, 6))
12 plt.plot(fpr, tpr, color='blue', lw=2, label=f'AUC = {auc:.2f}')
13 plt.plot([0, 1], [0, 1], color='gray', linestyle='--')
14 plt.xlabel('False Positive Rate')
15 plt.ylabel('True Positive Rate')
16 plt.title('ROC Curve')
17 plt.legend(loc='lower right')
18 plt.grid(True)
19 plt.show()
20
21 print(f"AUC: {auc:.2f}")

```



AUC: 0.82

It can be seen that the probabilities for home winning are working good.

```
In [66]: 1 # Create a new column indicating the bin for each match based on pred
2 bins = [0, 0.15, 0.3, 0.6, 1]
3 df_euro['probability_bin'] = pd.cut(df_euro['away_prob'], bins=bins)
4
5 # Group the DataFrame by the probability bins and compute the percenta
6 result_distribution = df_euro.groupby('probability_bin')['result'].va
7 print(result_distribution)
8 # Compute AUC
9 auc_away = roc_auc_score(df_euro['away_win_flag'], df_euro['away_prob
10
11 print(f"AUC: {auc_away:.2f}")
```

result	away_win	draw	home_win
probability_bin			
(0.0, 0.15]	8.007449	14.152700	77.839851
(0.15, 0.3]	25.223214	27.455357	47.321429
(0.3, 0.6]	53.053435	26.717557	20.229008
(0.6, 1.0]	87.421384	7.547170	5.031447

AUC: 0.81

It can be said the same for the away probabilities, so we can conclude that this predictions are relative good for the matches in Euro.

6 2024 Euro Prediction

Once we have the probabilities, we can finally use it to predict the results of the 2024 Euro. Let see it.

In [88]:

```
1 matches_data = [  
2     # Grupo 16  
3     {"Grupo": "phase3", "Home_Team": "Spain", "Away_Team": "France"},  
4     {"Grupo": "phase3", "Home_Team": "Netherlands", "Away_Team": "England"},  
5     {"Grupo": "phase3", "Home_Team": "Spain", "Away_Team": "England"},  
6     {"Grupo": "16", "Home_Team": "Spain", "Away_Team": "Georgia"},  
7     {"Grupo": "16", "Home_Team": "France", "Away_Team": "Belgium"},  
8     {"Grupo": "16", "Home_Team": "Portugal", "Away_Team": "Slovenia"},  
9     {"Grupo": "16", "Home_Team": "Romania", "Away_Team": "Netherlands"},  
10    {"Grupo": "16", "Home_Team": "Austria", "Away_Team": "Turkey"},  
11  
12    # Grupo 8  
13    {"Grupo": "8", "Home_Team": "Italy", "Away_Team": "England"},  
14    {"Grupo": "8", "Home_Team": "Germany", "Away_Team": "Spain"},  
15    {"Grupo": "8", "Home_Team": "France", "Away_Team": "Portugal"},  
16    {"Grupo": "8", "Home_Team": "Netherlands", "Away_Team": "Austria"},  
17  
18  
19  
20    {"Grupo": "16", "Home_Team": "Slovenia", "Away_Team": "Denmark"},  
21    {"Grupo": "16", "Home_Team": "Serbia", "Away_Team": "England"},  
22    {"Grupo": "16", "Home_Team": "Slovenia", "Away_Team": "Serbia"},  
23    {"Grupo": "16", "Home_Team": "Denmark", "Away_Team": "England"},  
24    {"Grupo": "16", "Home_Team": "England", "Away_Team": "Slovenia"},  
25    {"Grupo": "16", "Home_Team": "Denmark", "Away_Team": "Serbia"},  
26    # Grupo D  
27    {"Grupo": "16", "Home_Team": "Poland", "Away_Team": "Netherlands"},  
28    {"Grupo": "16", "Home_Team": "Austria", "Away_Team": "France"},  
29    {"Grupo": "16", "Home_Team": "Poland", "Away_Team": "Austria"},  
30    {"Grupo": "16", "Home_Team": "Netherlands", "Away_Team": "France"},  
31    {"Grupo": "D", "Home_Team": "Netherlands", "Away_Team": "Austria"},  
32    {"Grupo": "D", "Home_Team": "France", "Away_Team": "Poland"},  
33    # Grupo E  
34    {"Grupo": "E", "Home_Team": "Romania", "Away_Team": "Ukraine"},  
35    {"Grupo": "E", "Home_Team": "Belgium", "Away_Team": "Slovakia"},  
36    {"Grupo": "E", "Home_Team": "Slovakia", "Away_Team": "Ukraine"},  
37    {"Grupo": "E", "Home_Team": "Belgium", "Away_Team": "Romania"},  
38    {"Grupo": "E", "Home_Team": "Slovakia", "Away_Team": "Romania"},  
39    {"Grupo": "E", "Home_Team": "Ukraine", "Away_Team": "Belgium"},  
40    # Grupo F  
41    {"Grupo": "F", "Home_Team": "Turkey", "Away_Team": "Georgia"},  
42    {"Grupo": "F", "Home_Team": "Portugal", "Away_Team": "Czech Republic"},  
43    {"Grupo": "F", "Home_Team": "Georgia", "Away_Team": "Czech Republic"},  
44    {"Grupo": "F", "Home_Team": "Turkey", "Away_Team": "Portugal"},  
45    {"Grupo": "F", "Home_Team": "Georgia", "Away_Team": "Portugal"},  
46    {"Grupo": "F", "Home_Team": "Czech Republic", "Away_Team": "Turkey"},  
47 ]  
48  
49 # Create a dataframe  
50 matches_df = pd.DataFrame(matches_data)
```

In [89]:

```
1 matches_df
```

Out[89]:

	Grupo	Home_Team	Away_Team
0	phase3	Spain	France
1	phase3	Netherlands	England
2	phase3	Spain	England

In [90]:

1

elos

Out[90]:

	date	Team	Elo	position
0	2024-07-04	Argentina	1776.345	1
1	2024-07-05	France	1750.370	2
2	2024-07-05	Spain	1735.715	3
3	2024-07-02	Brazil	1699.620	4
4	2024-06-30	England	1677.695	5
5	2024-07-02	Netherlands	1675.520	6
6	2024-07-02	Colombia	1671.785	7
7	2024-07-01	Uruguay	1661.435	8
8	2024-07-05	Germany	1643.990	9
9	2024-07-01	Belgium	1634.765	10
10	2024-07-05	Portugal	1634.685	11

In [91]:

1

```
# merged_df = pd.merge(matches_df, attdef[["Team", "att", "def"]].add_prefix('att_'))
# merged_df_elo = pd.merge(merged_df, elos[["Team", "Elo"]].add_prefix('elo_'))
# merged_df_elo_form = pd.merge(merged_df_elo, group16_euroteams_sorted[["Team", "form", "Home_att", "Home_def", "Away_att", "Away_def"]].add_prefix('form_'))
# group_matches_elo = pd.merge(merged_df_elo_form, attdef[["Team", "att", "def"]].add_prefix('att_'))
# group_matches_form = pd.merge(group_matches_elo, elos[["Team", "Elo"]].add_prefix('elo_'))
# group_matches = pd.merge(group_matches_form, group16_euroteams_sorted[["Team", "form", "Home_att", "Home_def", "Away_att", "Away_def"]].add_prefix('form_'))
# # group_matches['form_Home_att']=group_matches['Home_att']+group_matches['form_Home_att']
# # group_matches['form_Home_def']=group_matches['Home_def']+group_matches['form_Home_def']
# # group_matches['form_Away_att']=group_matches['Away_att']+group_matches['form_Away_att']
# # group_matches['form_Away_def']=group_matches['Away_def']+group_matches['form_Away_def']
# group_matches
```

In [92]:

1

```
merged_df = pd.merge(matches_df, attdef[["Team", "att", "def"]].add_prefix('att_'))
merged_df_elo = pd.merge(merged_df, elos[["Team", "Elo"]].add_prefix('elo_'))
merged_df_elo_form = pd.merge(merged_df_elo, group16_euroteams_sorted[["Team", "form", "Home_att", "Home_def", "Away_att", "Away_def"]].add_prefix('form_'))
group_matches_elo = pd.merge(merged_df_elo_form, attdef[["Team", "att", "def"]].add_prefix('att_'))
group_matches_form = pd.merge(group_matches_elo, elos[["Team", "Elo"]].add_prefix('elo_'))
group_matches = pd.merge(group_matches_form, group16_euroteams_sorted[["Team", "form", "Home_att", "Home_def", "Away_att", "Away_def"]].add_prefix('form_'))
group_matches
```

Out[92]:

	Grupo	Home_Team	Away_Team	Home_att	Home_def	Home_Elo	Home_euro_form	Away_Elo
0	phase3	Spain	France	2.589579	0.457570	1735.715	2.000000	2.21
1	phase3	Netherlands	England	2.371421	0.618433	1675.520	0.888889	2.17
2	phase3	Spain	England	2.589579	0.457570	1735.715	2.000000	2.17

In [93]:

```
1 #Computing the expected goals
2 group_matches['XGhome']=group_matches['Home_att']*group_matches['Away_
3 group_matches['XGaway']=group_matches['Home_def']*group_matches['Away_
4 # group_matches['XGhome_form']=group_matches['Home_att']*group_matches
5 # group_matches['XGaway_form']=group_matches['Home_def']*group_matches
6 group_matches['home_diff_elo']=group_matches['Home_Elo']-group_matches
7 group_matches['away_diff_elo']=group_matches['Away_Elo']-group_matches
8
9 group_matches[['total_goals_prob', 'draw_prob', 'away_prob', 'local_p
10 group_matches.drop("total_goals_prob",axis=1, inplace=True)
```

In [94]:

```
1
2 group_matches
```

Out[94]:

	o	Away_euro_form	XGhome	XGaway	home_diff_elo	away_diff_elo	draw_prob	away_prob	I
	0	1.111111	2.390389	1.128333	-14.655	14.655	0.182667	0.164327	
	5	1.111111	1.000482	1.492086	-2.175	2.175	0.260620	0.485769	
	5	1.111111	2.458173	1.103973	58.020	-58.020	0.176375	0.153105	

```

In [74]: 1 # Define a function to determine the result based on probabilities
2 def determine_result_random(row):
3     # np.random.seed(1) # Note: Setting the seed inside the function
4     random_number = np.random.random()
5
6     if random_number < row["draw_prob"]:
7         return pd.Series(["Draw", random_number])
8     elif random_number < row["draw_prob"] + row["away_prob"]:
9         return pd.Series([row["Away_Team"], random_number])
10    else:
11        return pd.Series([row["Home_Team"], random_number])
12
13 def determine_result(row):
14     # max_prob = max(row["draw_prob"], row["away_prob"], row["local_p
15
16     if row["home_goals_elo"] == row["away_goals_elo"]:
17         return "Draw"
18     elif row["away_goals_elo"] > row["away_goals_elo"]:
19         return row["Away_Team"]
20     else:
21         return row["Home_Team"]
22
23 def determine_result_knockout(row):
24     if row["away_prob"] > row["home_prob"]:
25         return row["Away_Team"]
26     else:
27         return row["Home_Team"]
28
29 def determine_result_knockout_random(row):
30     random_number = np.random.random()
31     if random_number < row["draw_prob"]:
32         random_number2 = np.random.random()
33         if random_number < 0.5:
34             return row["Away_Team"]
35         else:
36             return row["Home_Team"]
37     elif random_number < row["draw_prob"] + row["away_prob"]:
38         return row["Away_Team"]
39     else:
40         return row["Home_Team"]
41
42 # Apply the function to create the "Results" column
43 group_matches["Result"] = group_matches.apply(determine_result, axis=
44

```

We can see the predictions for each match, based on the expected goals, we simulate the number of goals scored by each team in a match.

```
In [75]: 1 for x in range(-532,853):
          2     h_goals_adjust=int(reg.predict(np.array(x).reshape(-1, 1)))
          3     print(h_goals_adjust,x)
```

-1 -532
-1 -531
-1 -530
-1 -529
-1 -528
-1 -527
-1 -526
-1 -525
-1 -524
-1 -523
-1 -522
-1 -521
-1 -520
-1 -519
-1 -518
-1 -517
-1 -516
-1 -515
-1 -514
-1 -513

```
In [76]: 1 group_matches[group_matches['away_diff_elo']<=-243]
```

Out[76]:

	Grupo	Home_Team	Away_Team	Home_att	Home_def	Home_Elo	Home_euro_form	Away
1	16	England	Slovakia	2.171419	0.474628	1677.695	1.111111	1.33
3	16	Spain	Georgia	2.589579	0.457570	1735.715	2.000000	1.18

```
In [77]: 1 group_matches[group_matches['home_diff_elo']<=-243]
```

Out[77]:

	Grupo	Home_Team	Away_Team	Home_att	Home_def	Home_Elo	Home_euro_form	Away
6	16	Romania	Netherlands	1.310627	0.792474	1423.0	0.888889	2.37

```
In [78]: 1 group_matches[group_matches.Grupo=="A"][["Home_Team", "Away_Team", "local_prob", "draw_prob", "away_prob", "home_goals", "away_goals", "home_goals_adj", "away_goals_adj"]]
```

Out[78]:

	Home_Team	Away_Team	local_prob	draw_prob	away_prob	home_goals	away_goals	home_goals_adj	away_goals_adj
1	England	Slovakia	0.474628	0.474628	0.474628	1677.695	1.111111	1.33	1.33

```
In [79]: 1 group_matches[group_matches.Grupo=="B"][["Home_Team", "Away_Team", "local_prob", "draw_prob", "away_prob", "home_goals", "away_goals", "home_goals_adj", "away_goals_adj"]]
```

Out[79]:

	Home_Team	Away_Team	local_prob	draw_prob	away_prob	home_goals	away_goals	home_goals_adj	away_goals_adj
1	Romania	Netherlands	0.792474	0.792474	0.792474	1423.0	0.888889	2.37	2.37

In [80]: 1 group_matches[group_matches.Grupo=="C"](["Home_Team", "Away_Team", "local_prob", "draw_prob", "away_prob", "home_goals", "away_goals", "home_goals_elo", "away_goals_elo", "Result"])

Out[80]:

	Home_Team	Away_Team	local_prob	draw_prob	away_prob	home_goals	away_goals	home_goals_elo	away_goals_elo	Result
0	Switzerland	Italy	1.769522	0.723214	1594.235	1.111111	1.769522	1.111111	1.769522	Draw

In [81]: 1 group_matches[group_matches.Grupo=="D"](["Home_Team", "Away_Team", "local_prob", "draw_prob", "away_prob", "home_goals", "away_goals", "home_goals_elo", "away_goals_elo", "Result"])

Out[81]:

	Home_Team	Away_Team	local_prob	draw_prob	away_prob	home_goals	away_goals	home_goals_elo	away_goals_elo	Result
7	Austria	Turkey	1.585274	0.848737	1521.590	1.333333	1.585274	1.333333	1.585274	Draw

In [82]: 1 group_matches[group_matches.Grupo=="E"](["Home_Team", "Away_Team", "local_prob", "draw_prob", "away_prob", "home_goals", "away_goals", "home_goals_elo", "away_goals_elo", "Result"])

Out[82]:

	Home_Team	Away_Team	local_prob	draw_prob	away_prob	home_goals	away_goals	home_goals_elo	away_goals_elo	Result
10	France	Portugal	2.219333	0.461540	1750.370	1.111111	2.219333	1.111111	2.219333	Draw

In [83]: 1 group_matches[group_matches.Grupo=="F"](["Home_Team", "Away_Team", "local_prob", "draw_prob", "away_prob", "home_goals", "away_goals", "home_goals_elo", "away_goals_elo", "Result"])

Out[83]:

	Home_Team	Away_Team	local_prob	draw_prob	away_prob	home_goals	away_goals	home_goals_elo	away_goals_elo	Result
11	Netherlands	Austria	2.371421	0.618433	1675.520	0.888889	2.371421	0.888889	2.371421	Draw

By running numerous simulations, we can estimate the probabilities of different outcomes for each match. We simulate each match multiple times to predict the likely outcomes and the points each team might accumulate. The table presented below illustrates the predicted probabilities for each team to conclude the group stage in a particular position.

In [84]: 1 group_matches[group_matches['Result']=='Draw']

Out[84]:

	Grupo	Home_Team	Away_Team	Home_att	Home_def	Home_Elo	Home_euro_form	Away_att	Away_def	Away_Elo	Away_euro_form
0	16	Switzerland	Italy	1.769522	0.723214	1594.235	1.111111	1.769522	0.723214	1594.235	1.111111
7	16	Austria	Turkey	1.585274	0.848737	1521.590	1.333333	1.585274	0.848737	1521.590	1.333333
10	8	France	Portugal	2.219333	0.461540	1750.370	1.111111	2.219333	0.461540	1750.370	1.111111
11	8	Netherlands	Austria	2.371421	0.618433	1675.520	0.888889	2.371421	0.618433	1675.520	0.888889

In [85]: 1 group_matches.columns

Out[85]: Index(['Grupo', 'Home_Team', 'Away_Team', 'Home_att', 'Home_def', 'Home_Elo', 'Home_euro_form', 'Away_att', 'Away_def', 'Away_Elo', 'Away_euro_form', 'XGhome', 'XGaway', 'home_diff_elo', 'away_diff_elo', 'draw_prob', 'away_prob', 'local_prob', 'home_goals', 'away_goals', 'home_goals_elo', 'away_goals_elo', 'Result'], dtype='object')

```
In [86]: 1 def calculate_points(result, home_team, team):  
2     if result == "Draw":  
3         return 1  
4     elif result == home_team:  
5         return 3  
6     elif result == team:  
7         return 0
```

```

In [87]: 1 # Simulate match results and calculate standings
2 num_simulations = 10000 # Number of simulations
3 standings_list = [] # List to store standings for each simulation
4
5 for _ in range(num_simulations):
6     # Apply the function to create the "Results" column
7     group_matches[["Random_Result", 'Random_number']] = group_matches.
8     group_matches["Result"] = group_matches.apply(determine_result, a
9
10     # Apply the function to create the "Points" column for both Home
11     group_matches["Points_Home"] = group_matches.apply(lambda row: ca
12     group_matches["Points_Away"] = group_matches.apply(lambda row: ca
13     # Create a dataframe to store points for each teams
14     standings = pd.DataFrame(columns=["Grupo", "Team", "Points"])
15     # Concatenate Home and Away points by group
16     for group_name, group_data in group_matches.groupby("Grupo"):
17         group_points_df = pd.DataFrame(columns=["Team", "Points"])
18         for team in set(group_data["Home_Team"]).union(set(group_data
19             home_points = group_data.loc[group_data["Home_Team"] == t
20             away_points = group_data.loc[group_data["Away_Team"] == t
21             total_points = home_points + away_points
22             group_points_df = group_points_df._append({"Team": team,
23             group_points_df = group_points_df.sort_values(by=["Points"],
24             group_points_df["Grupo"] = group_name
25             group_points_df["Rank"] = range(1, len(group_points_df) + 1)
26             group_points_df["Qualified_1"] = np.where(group_points_df["Rank
27             group_points_df["Qualified_2"] = np.where(group_points_df["Rank
28             group_points_df["Qualified_3"] = np.where(group_points_df["Rank
29             group_points_df["Qualified_4"] = np.where(group_points_df["Rank
30             standings = pd.concat([standings, group_points_df], ignore_in
31             standings_list.append(standings)

```



```

-----
-
KeyError                                Traceback (most recent call last)
File C:\ProgramData\anaconda3\Lib\site-packages\pandas\core\indexes\base.py:3653, in Index.get_loc(self, key)
    3652 try:
-> 3653     return self._engine.get_loc(casted_key)
    3654 except KeyError as err:

File C:\ProgramData\anaconda3\Lib\site-packages\pandas\_libs\index.pyx:147, in pandas._libs.index.IndexEngine.get_loc()

File C:\ProgramData\anaconda3\Lib\site-packages\pandas\_libs\index.pyx:176, in pandas._libs.index.IndexEngine.get_loc()

File pandas\_libs\hashtable_class_helper.pxi:7080, in pandas._libs.hashtable.PyObjectHashTable.get_item()

File pandas\_libs\hashtable_class_helper.pxi:7088, in pandas._libs.hashtable.PyObjectHashTable.get_item()

```

KeyError: 'Rank'

The above exception was the direct cause of the following exception:

```

KeyError                                Traceback (most recent call last)
File C:\ProgramData\anaconda3\Lib\site-packages\pandas\core\frame.py:4110, in DataFrame._set_item_mgr(self, key, value)
    4109 try:
-> 4110     loc = self._info_axis.get_loc(key)
    4111 except KeyError:
    4112     # This item wasn't present, just insert at end

File C:\ProgramData\anaconda3\Lib\site-packages\pandas\core\indexes\base.py:3655, in Index.get_loc(self, key)
    3654 except KeyError as err:
-> 3655     raise KeyError(key) from err
    3656 except TypeError:
    3657     # If we have a listlike key, _check_indexing_error will raise
    3658     # InvalidIndexError. Otherwise we fall through and re-raise
    3659     # the TypeError.

```

KeyError: 'Rank'

During handling of the above exception, another exception occurred:

```

KeyboardInterrupt                        Traceback (most recent call last)
Cell In[87], line 25
    23 group_points_df = group_points_df.sort_values(by=["Points"], ascending=False).reset_index(drop=True)
    24 group_points_df["Grupo"] = group_name
--> 25 group_points_df["Rank"] = range(1, len(group_points_df) + 1)
    26 group_points_df["Qualified_1"] = np.where(group_points_df["Rank"] == 1, 1, 0)
    27 group_points_df["Qualified_2"] = np.where(group_points_df["Rank"] == 2, 1, 0)

```

File C:\ProgramData\anaconda3\Lib\site-packages\pandas\core\frame.py:3950,

```

in DataFrame.__setitem__(self, key, value)
3947     self._setitem_array([key], value)
3948 else:
3949     # set column
-> 3950     self._set_item(key, value)

File C:\ProgramData\anaconda3\Lib\site-packages\pandas\core\frame.py:4156,
in DataFrame._set_item(self, key, value)
4153     if isinstance(existing_piece, DataFrame):
4154         value = np.tile(value, (len(existing_piece.columns),
1)).T
-> 4156 self._set_item_mgr(key, value)

File C:\ProgramData\anaconda3\Lib\site-packages\pandas\core\frame.py:4113,
in DataFrame._set_item_mgr(self, key, value)
4110     loc = self._info_axis.get_loc(key)
4111 except KeyError:
4112     # This item wasn't present, just insert at end
-> 4113     self._mgr.insert(len(self._info_axis), key, value)
4114 else:
4115     self._iset_item_mgr(loc, value)

File C:\ProgramData\anaconda3\Lib\site-packages\pandas\core\internals\mana
gers.py:1419, in BlockManager.insert(self, loc, item, value)
1417 else:
1418     self._insert_update_mgr_locs(loc)
-> 1419     self._insert_update_blklocs_and_blkgnos(loc)
1421 self.axes[0] = new_axis
1422 self.blocks += (block,)

File C:\ProgramData\anaconda3\Lib\site-packages\pandas\core\internals\mana
gers.py:1456, in BlockManager._insert_update_blklocs_and_blkgnos(self, loc)
1453 if loc == self.blklocs.shape[0]:
1454     # np.append is a lot faster, let's use it if we can.
1455     self._blklocs = np.append(self._blklocs, 0)
-> 1456     self._blkgnos = np.append(self._blkgnos, len(self.blocks))
1457 elif loc == 0:
1458     # np.append is a lot faster, let's use it if we can.
1459     self._blklocs = np.append(self._blklocs[:-1], 0)[::-1]

File <__array_function__ internals>:200, in append(*args, **kwargs)

File C:\ProgramData\anaconda3\Lib\site-packages\numpy\lib\function_base.p
y:5499, in append(arr, values, axis)
5497     values = ravel(values)
5498     axis = arr.ndim-1
-> 5499 return concatenate((arr, values), axis=axis)

File <__array_function__ internals>:200, in concatenate(*args, **kwargs)

```

KeyboardInterrupt:

```
In [ ]: 1 group_matches[group_matches['Result']!=group_matches['Random_Result']]
```

```
In [95]:
1
2
3 # Calculate total home goals for each team
4 home_goals = group_matches.groupby('Home_Team')['home_goals_elo'].sum
5
6 # Calculate total away goals for each team
7 away_goals = group_matches.groupby('Away_Team')['away_goals_elo'].sum
8
9 # Combine home and away goals to get the total goals for each team
10 total_goals = home_goals.add(away_goals, fill_value=0)
11
12 # Convert the Series to DataFrame for better presentation
13 total_goals_df = total_goals.reset_index()
14 total_goals_df.columns = ['Team', 'Total_Goals']
15 total_goals_df.sort_values('Total_Goals', ascending=False, inplace=True)
16 total_goals_df
```

	Team	Total_Goals
3	Spain	4.0
0	England	2.0
1	France	1.0
2	Netherlands	1.0

7 Top Goal scorer

	Pos.	Player	Date of birth	age	Caps	Goals	G/C	Club	Country
0	FW	Niclas Füllkrug	1993-02-09	31	15	11	0.733333	Germany Borussia Dortmund	Germany
1	FW	Romelu Lukaku	1993-05-13	31	114	83	0.728070	Italy Roma	Belgium
2	FW	Harry Kane (captain)	1993-07-28	30	90	63	0.700000	Germany Bayern Munich	England
3	FW	Gonçalo Ramos	2001-06-20	22	12	8	0.666667	France Paris Saint-Germain	Portugal
4	FW	Tomáš Chory	1995-01-26	29	3	2	0.666667	Czech Republic Viktoria Plzeň	Czech Republic
5	FW	Aleksandar Mitrović	1994-09-16	29	90	57	0.633333	Saudi Arabia Al Hilal	Serbia

In [97]: 1 euro_players.columns

Out[97]: Index(['Pos.', 'Player', 'Date of birth ', 'age', 'Caps', 'Goals', 'G/C',
'Club', 'Country', 'Manager', 'Group'],
dtype='object')

In [98]: 1 sorted_df = euro_players.groupby('Country').apply(lambda x: x.sort_val
2 sorted_df.reset_index(drop=True, inplace=True)
3 sorted_df['Rank'] = sorted_df.groupby('Country')['G/C'].rank(method='c
4 final_players=sorted_df[sorted_df['Rank'] <=5][['Player', 'Country', 'G
5 final_players

Out[98]:

	Player	Country	G/C	Rank
0	Jasir Asani	Albania	0.333333	1.0
1	Ernest Muçi	Albania	0.300000	2.0
2	Armando Broja	Albania	0.250000	3.0
3	Rey Manaj	Albania	0.212121	4.0
4	Mirlind Daku	Albania	0.200000	5.0
27	Christoph Baumgartner	Austria	0.378378	1.0
28	Maximilian Entrup	Austria	0.333333	2.0
29	Marko Arnautović	Austria	0.321429	3.0
30	Michael Gregoritsch	Austria	0.277778	4.0
31	Marcel Sabitzer	Austria	0.217949	5.0
56	Romelu Lukaku	Belgium	0.728070	1.0

In [99]:

```
1 final_players_sort = pd.merge(final_players, total_goals_df, left_on=  
2 final_players_sort  
3 # Calculate proportional goals based on G/C ratio  
4 # Normalize the G/C ratio within each country to distribute goals  
5 # final_players_sort['G/C Normalized'] = final_players_sort.groupby('C'  
6 # final_players_sort['Player Goals'] = (final_players_sort['G/C Normal  
7  
8 # final_players_sort.sort_values('Player Goals',ascending=False,inplace=  
9
```

Out[99]:

	Player	Country	G/C	Rank	Team	Total_Goals
0	Harry Kane (captain)	England	0.700000	1.0	England	2.0
1	Ivan Toney	England	0.500000	2.0	England	2.0
2	Bukayo Saka	England	0.343750	3.0	England	2.0
3	Cole Palmer	England	0.333333	4.0	England	2.0
4	Ollie Watkins	England	0.250000	5.0	England	2.0
5	Kylian Mbappé (captain)	France	0.597403	1.0	France	1.0
6	Warren Zaïre-Emery	France	0.500000	2.0	France	1.0
7	Olivier Giroud	France	0.435115	3.0	France	1.0
8	Antoine Griezmann	France	0.346457	4.0	France	1.0
9	Randal Kolo Muani	France	0.200000	5.0	France	1.0
10	Memphis Depay	Netherlands	0.488889	1.0	Netherlands	1.0
11	Cody Gakpo	Netherlands	0.391304	2.0	Netherlands	1.0
12	Georginio Wijnaldum	Netherlands	0.307692	3.0	Netherlands	1.0
13	Wout Weghorst	Netherlands	0.290323	4.0	Netherlands	1.0
14	Steven Bergwijn	Netherlands	0.250000	5.0	Netherlands	1.0
15	Joselu	Spain	0.500000	1.0	Spain	4.0
16	Álex Baena	Spain	0.500000	1.0	Spain	4.0
17	Álvaro Morata (captain)	Spain	0.478873	2.0	Spain	4.0
18	Ferran Torres	Spain	0.450000	3.0	Spain	4.0
19	Lamine Yamal	Spain	0.333333	4.0	Spain	4.0
20	Mikel Oyarzabal	Spain	0.250000	5.0	Spain	4.0

In [100]:

```
1 final_players_sort['G/C_Normalized'] = final_players_sort.groupby('Country')['G/C'].transform('mean')
2 final_players_sort
```

Out[100]:

	Player	Country	G/C	Rank	Team	Total_Goals	G/C_Normalized
0	Harry Kane (captain)	England	0.700000	1.0	England	2.0	0.329089
1	Ivan Toney	England	0.500000	2.0	England	2.0	0.235064
2	Bukayo Saka	England	0.343750	3.0	England	2.0	0.161606
3	Cole Palmer	England	0.333333	4.0	England	2.0	0.156709
4	Ollie Watkins	England	0.250000	5.0	England	2.0	0.117532
5	Kylian Mbappé (captain)	France	0.597403	1.0	France	1.0	0.287355
6	Warren Zaïre-Emery	France	0.500000	2.0	France	1.0	0.240503
7	Olivier Giroud	France	0.435115	3.0	France	1.0	0.209293
8	Antoine Griezmann	France	0.346457	4.0	France	1.0	0.166648
9	Randal Kolo Muani	France	0.200000	5.0	France	1.0	0.096201
10	Memphis Depay	Netherlands	0.488889	1.0	Netherlands	1.0	0.282888
11	Cody Gakpo	Netherlands	0.391304	2.0	Netherlands	1.0	0.226422
12	Georginio Wijnaldum	Netherlands	0.307692	3.0	Netherlands	1.0	0.178041
13	Wout Weghorst	Netherlands	0.290323	4.0	Netherlands	1.0	0.167991
14	Steven Bergwijn	Netherlands	0.250000	5.0	Netherlands	1.0	0.144659
15	Joselu	Spain	0.500000	1.0	Spain	4.0	0.199028
16	Álex Baena	Spain	0.500000	1.0	Spain	4.0	0.199028
17	Álvaro Morata (captain)	Spain	0.478873	2.0	Spain	4.0	0.190619
18	Ferran Torres	Spain	0.450000	3.0	Spain	4.0	0.179125
19	Lamine Yamal	Spain	0.333333	4.0	Spain	4.0	0.132685
20	Mikel Oyarzabal	Spain	0.250000	5.0	Spain	4.0	0.099514

In [101]:

```
1 final_players_sort['Player_Goals'] = np.floor(final_players_sort['G/C_Normalized'] * 10)
2
3 unique_total_goals = final_players_sort.drop_duplicates(subset=['Country']).sum('Player_Goals')
4
5 # Calculate remaining goals after initial assignment
6 assigned_goals = final_players_sort['Player_Goals'].sum()
7 remaining_goals = unique_total_goals - assigned_goals
8
9 remaining_goals
```

Out[101]: 8.0

```
In [102]: 1 # Distribute remaining goals starting from the players with the highest
2 sorted_indices = final_players_sort.sort_values(by=['G/C_Normalized'],
3
4 # Ensure we do not exceed the number of available players
5 num_players = len(sorted_indices)
6 to_distribute = min(remaining_goals, num_players)
7
8 # Loop to increment goals for the top ranked players based on remaining
9 for i in range(int(to_distribute)):
10     final_players_sort.loc[sorted_indices[i], 'Player_Goals'] += 1
11
12 final_players_sort
```

Out[102]:

	Player	Country	G/C	Rank	Team	Total_Goals	G/C_Normalized	Player.
0	Harry Kane (captain)	England	0.700000	1.0	England	2.0	0.329089	
1	Ivan Toney	England	0.500000	2.0	England	2.0	0.235064	
2	Bukayo Saka	England	0.343750	3.0	England	2.0	0.161606	
3	Cole Palmer	England	0.333333	4.0	England	2.0	0.156709	
4	Ollie Watkins	England	0.250000	5.0	England	2.0	0.117532	
5	Kylian Mbappé (captain)	France	0.597403	1.0	France	1.0	0.287355	
6	Warren Zaire-Emery	France	0.500000	2.0	France	1.0	0.240503	
7	Olivier Giroud	France	0.435115	3.0	France	1.0	0.209293	
8	Antoine Griezmann	France	0.346457	4.0	France	1.0	0.166648	
9	Randal Kolo Muani	France	0.200000	5.0	France	1.0	0.096201	
10	Memphis Depay	Netherlands	0.488889	1.0	Netherlands	1.0	0.282888	
11	Cody Gakpo	Netherlands	0.391304	2.0	Netherlands	1.0	0.226422	
12	Georginio Wijnaldum	Netherlands	0.307692	3.0	Netherlands	1.0	0.178041	
13	Wout Weghorst	Netherlands	0.290323	4.0	Netherlands	1.0	0.167991	
14	Steven Bergwijn	Netherlands	0.250000	5.0	Netherlands	1.0	0.144659	
15	Joselu	Spain	0.500000	1.0	Spain	4.0	0.199028	
16	Álex Baena	Spain	0.500000	1.0	Spain	4.0	0.199028	
17	Álvaro Morata (captain)	Spain	0.478873	2.0	Spain	4.0	0.190619	
18	Ferran Torres	Spain	0.450000	3.0	Spain	4.0	0.179125	
19	Lamine Yamal	Spain	0.333333	4.0	Spain	4.0	0.132685	
20	Mikel Oyarzabal	Spain	0.250000	5.0	Spain	4.0	0.099514	

In [103]:

1 final_players_sort['Player_Goals'].sum()

Out[103]: 8


```
In [104]: 1 final_players_sort.sort_values('Player_Goals',ascending=False,inplace=True)
          2 final_players_sort[:10]
```

Out[104]:

	Player	Country	G/C	Rank	Team	Total_Goals	G/C_Normalized	Player
0	Harry Kane (captain)	England	0.700000	1.0	England	2.0	0.329089	
6	Warren Zaïre-Emery	France	0.500000	2.0	France	1.0	0.240503	
15	Joselu	Spain	0.500000	1.0	Spain	4.0	0.199028	
11	Cody Gakpo	Netherlands	0.391304	2.0	Netherlands	1.0	0.226422	
1	Ivan Toney	England	0.500000	2.0	England	2.0	0.235064	
7	Olivier Giroud	France	0.435115	3.0	France	1.0	0.209293	
10	Memphis Depay	Netherlands	0.488889	1.0	Netherlands	1.0	0.282888	
5	Kylian Mbappé (captain)	France	0.597403	1.0	France	1.0	0.287355	
8	Antoine Griezmann	France	0.346457	4.0	France	1.0	0.166648	
9	Randal Kolo Muani	France	0.200000	5.0	France	1.0	0.096201	

```
In [105]: 1 "Niclas Füllkrug"
```

Out[105]: 'Niclas Füllkrug'

As we can see, the teams that are more likely to win its respecptive group are Portugal, Belgium and England. And the team which have more difficult to pass the round is Albania.

8 Conclusion

Predicting football match outcomes is inherently uncertain, but statistical models like the Poisson distribution offer a structured approach to quantify this uncertainty. Our predictions for the 2024 UEFA Euro provide a glimpse into possible outcomes based on historical performance data. As with all models, the key is to continually refine our approach as new data becomes available.