

Top 5 Leagues goals by nationality

November 19, 2024

0.1 Data Parsing: Top 5 Football Leagues Historical Goalscorers by Nation

Ever wondered the amount of goals that players from your country have scored in Europe Top 5 Football Leagues? Well, thanks to powerful libraries such as BeautifulSoup and Pandas, this is possible.

All the data is collected from <https://www.worldfootball.net/goalgetter/>.

*The data begins in the season of 1963-1964 because this was the year where Bundesliga was founded. Therefore, it would be unfair to consider previous years.

Import the necessary libraries

```
[12]: # import libraries for data manipulation
import numpy as np
import pandas as pd

# import libraries for parsing
import requests
from bs4 import BeautifulSoup

# to suppress warnings
import warnings
warnings.filterwarnings('ignore')
```

Helper function to generate text for the columns having two dates as parameters.

```
[13]: def generate_seasons_years(from_date, to_date):
    seasons_text = []
    for year in range(from_date, to_date):
        seasons_text.append(str(year) + "-" + str(year + 1))
    return seasons_text
```

Recursive function that fills the goals into the dictionary per season.

```
[14]: def fill_data(global_dict, country, season, goals,
    ↪ empty_seasons_dictionary=None):
    if (country in global_dict.keys()):
        global_dict[country][season] += goals
    return
```

```

    else: # In case this is the first time that a country appears, first it is
    ↪ initialized and then filled through recursion.
        global_dict[country] = empty_seasons_dictionary.copy()
        fill_data(global_dict, country, season, goals)

```

Function that parses the webpages. Then it extracts the relevant keywords to populate the dictionary

```

[15]: def parse_and_fill(global_dict, url, season, empty_seasons_dictionary):
    # Fetch the webpage content
    response = requests.get(url)

    # Parse the HTML using BeautifulSoup
    soup = BeautifulSoup(response.content, 'html.parser')

    # Locate the table containing the data (goal scorers, etc.)
    table = soup.find('table', class_='standard_tabelle') # Look for the
    ↪ specific class used in the table (in case there's multiple)

    # Extract the data
    rows = table.find_all('tr')
    for row in rows[1:]: # Skip the header row
        cols = row.find_all('td')
        cols = [col.text.strip() for col in cols] # Clean the text
        # Save important data
        country = cols[3]
        goals = int(cols[5].rsplit(" ")[0])
        # Use the data to populate the dictionary
        fill_data(global_dict, country, season, goals, empty_seasons_dictionary)

```

Function that produces the right url. Some webpages have a non-intuitive webpage, therefore some if statements are introduced

```

[16]: def get_urls(league, season):
    # Base website url
    base_url = 'https://www.worldfootball.net/goalgetter/'
    urls = []

    # Checks for specific cases
    if (league == "esp-primera-division" and season == "2016-2017"):
        urls.append(base_url + league + "-" + season + "_2/")
    elif (league == "esp-primera-division" and season == "1986-1987"):
        spain_leagues = ["esp-primera-division-1986-1987-playoff-1-6",
    ↪ "esp-primera-division-1986-1987-playoff-13-18",
    ↪ "esp-primera-division-1986-1987-playoff-7-12",
    ↪ "esp-primera-division-1986-1987-vorrunde"]

```

```

        # In this season, the league was divided in sub-leagues. This for loop
        ↪ makes sure that all of them are included.
        for spain_league in spain_leagues:
            urls.append(base_url + spain_league + "/")
    else:
        # The most common case.
        urls.append(base_url + league + "-" + season + "/")
    return urls

```

Function that initializes a dictionary row with 0s as values.

```

[17]: def create_empty_seasons_dictionary(seasons):
    # Empty dictionary is defined
    seasons_dictionary = {}
    # Populate the dictionary
    for season in seasons:
        seasons_dictionary[season] = 0
    return seasons_dictionary

```

Main Function that iterates over each season and each league and populates the dictionary using helper functions.

```

[18]: def extract_values_top_5_leagues(from_date, to_date):
    # Sets the relevant parameters for the iterations
    goals_per_nation_and_year = {}
    seasons = generate_seasons_years(int(from_date), int(to_date))
    empty_seasons_dictionary = create_empty_seasons_dictionary(seasons)
    leagues = ["eng-premier-league", "fra-ligue-1", "bundesliga",
    ↪ "ita-serie-a", "esp-primera-division"]

    # Main loop that iterates through every leaguer per each season.
    for season in seasons:
        for league in leagues:
            urls = get_urls(league, season)
            for url in urls: # for the case where there are multiple urls in
            ↪ one league in a single season
                parse_and_fill(goals_per_nation_and_year, url, season,
                ↪ empty_seasons_dictionary)

    #Returns a sorted dictionary based on the name of the keys.
    return dict(sorted(goals_per_nation_and_year.items()))

```

Calls the main function

```

[19]: # Main dictionary produced by the program stored in a variable
final_dictionary = extract_values_top_5_leagues(1963,2024)

```

```
[43]: # Check if an empty key exists and delete it if so
      if "" in final_dictionary:
          del final_dictionary[""]

      # Show the first 5 seasons for the first 20 countries (alphabetically sorted)
      for element in list(final_dictionary.items())[:20]:
          print(element[0], list(element[1].items())[:5])
```

```
Albania [('1963-1964', 0), ('1964-1965', 0), ('1965-1966', 0), ('1966-1967', 0),
('1967-1968', 0)]
Algeria [('1963-1964', 91), ('1964-1965', 49), ('1965-1966', 58), ('1966-1967',
27), ('1967-1968', 34)]
Angola [('1963-1964', 2), ('1964-1965', 13), ('1965-1966', 10), ('1966-1967',
6), ('1967-1968', 7)]
Antigua & Barbuda [('1963-1964', 0), ('1964-1965', 0), ('1965-1966', 0),
('1966-1967', 0), ('1967-1968', 0)]
Argentina [('1963-1964', 129), ('1964-1965', 92), ('1965-1966', 104),
('1966-1967', 76), ('1967-1968', 58)]
Armenia [('1963-1964', 0), ('1964-1965', 0), ('1965-1966', 0), ('1966-1967', 0),
('1967-1968', 0)]
Australia [('1963-1964', 0), ('1964-1965', 0), ('1965-1966', 0), ('1966-1967',
0), ('1967-1968', 0)]
Austria [('1963-1964', 18), ('1964-1965', 12), ('1965-1966', 17), ('1966-1967',
8), ('1967-1968', 11)]
Azerbaijan [('1963-1964', 0), ('1964-1965', 0), ('1965-1966', 0), ('1966-1967',
0), ('1967-1968', 0)]
Barbados [('1963-1964', 0), ('1964-1965', 0), ('1965-1966', 0), ('1966-1967',
0), ('1967-1968', 0)]
Belarus [('1963-1964', 0), ('1964-1965', 0), ('1965-1966', 0), ('1966-1967', 0),
('1967-1968', 0)]
Belgium [('1963-1964', 7), ('1964-1965', 7), ('1965-1966', 6), ('1966-1967', 8),
('1967-1968', 9)]
Benin [('1963-1964', 0), ('1964-1965', 0), ('1965-1966', 0), ('1966-1967', 0),
('1967-1968', 0)]
Bermuda [('1963-1964', 0), ('1964-1965', 0), ('1965-1966', 0), ('1966-1967', 0),
('1967-1968', 0)]
Bolivia [('1963-1964', 0), ('1964-1965', 0), ('1965-1966', 0), ('1966-1967', 0),
('1967-1968', 0)]
Bosnia-Herzegovina [('1963-1964', 0), ('1964-1965', 0), ('1965-1966', 0),
('1966-1967', 0), ('1967-1968', 0)]
Brazil [('1963-1964', 139), ('1964-1965', 103), ('1965-1966', 107),
('1966-1967', 93), ('1967-1968', 78)]
Bulgaria [('1963-1964', 0), ('1964-1965', 0), ('1965-1966', 0), ('1966-1967',
0), ('1967-1968', 0)]
Burkina Faso [('1963-1964', 0), ('1964-1965', 0), ('1965-1966', 0),
('1966-1967', 0), ('1967-1968', 0)]
Burundi [('1963-1964', 0), ('1964-1965', 0), ('1965-1966', 0), ('1966-1967', 0),
```

```
('1967-1968', 0)]
```

0.1.1 Creation of CSV/Excel file.

Initialize a list where the dictionary will be transformed.

```
[44]: list_for_csv = []
      # Name for the outer keys stored in the header
      headers = ["Countries"]
      for country, inner_dict in final_dictionary.items():
          for key in inner_dict.keys():
              # Names of the inner keys (seasons) stored in the header
              headers.append(key)

          # Just add it once
          break
```

Now populate the list with the correct format

```
[45]: # Loop that iterates over the inner dictionary items
      for country, inner_dict in final_dictionary.items():
          # Creates a row with the country as its first value
          country_goals = [country]
          for value in inner_dict.values():
              # Appends the goals per season in the right order
              country_goals.append(int(value))
          list_for_csv.append(country_goals)
```

Makes the necessary arrangements to convert it into a dataframe

```
[46]: # Saves the python list as a numpy array
      list_as_numpy_array = np.array(list_for_csv)
      # Creates the dataframe
      df = pd.DataFrame(list_as_numpy_array, columns=headers)
      # Forces numerical value
      df.iloc[:, 1:] = df.iloc[:, 1:].apply(pd.to_numeric)
      # Creates a column that accumulates all the goals per country
      df['sum'] = df.iloc[:, 1:].sum(axis=1)
      # Sorts the dataframe by cumulative total sum.
      df = df.sort_values(by="sum", ascending = False)
      # Resets index to assure proper display
      df.reset_index(drop = True, inplace=True)

      # Saves the file as CSV or Excel

      ### df.to_csv('top_5_leagues_countries_cumulative.csv', index=False)
      ### df.to_excel('top_5_leagues_countries_cumulative.xlsx', index=False)
```

```
[47]: df
```

```

[47]:      Countries 1963-1964 1964-1965 1965-1966 1966-1967 1967-1968 1968-1969 \
0      Germany      830      760      953      845      913      776
1      England     1132     1061     953      895     982      890
2      France       721      701     967      801     771     585
3      Spain        495      495     499      578     576     491
4      Italy        373      457     466      468     383     415
..      ...
145     Tanzania      0      0      0      0      0      0
146  North Korea      0      0      0      0      0      0
147     Cambodia      0      0      0      0      0      0
148   Kazakhstan      0      0      0      0      0      0
149      Iraq         0      0      0      0      0      0

      1969-1970 1970-1971 1971-1972 ... 2015-2016 2016-2017 2017-2018 \
0      881      869      940 ...      360      359      389
1      836      736      783 ...      297      291      275
2      699      759      788 ...      552      603      513
3      516      474      596 ...      573      647      584
4      396      423      427 ...      439      472      431
..      ...
145      0      0      0 ...      0      0      0
146      0      0      0 ...      0      1      0
147      0      0      0 ...      0      0      0
148      0      0      0 ...      0      0      0
149      0      0      0 ...      1      0      0

      2018-2019 2019-2020 2020-2021 2021-2022 2022-2023 2023-2024      sum
0      382      363      351      365      419      440  39159
1      284      366      372      358      387      459  37994
2      614      493      614      700      681      568  37710
3      608      600      613      618      503      555  34418
4      417      424      374      392      310      321  27812
..      ...
145      0      1      0      0      0      0      1
146      0      0      0      0      0      0      1
147      0      0      0      0      0      0      1
148      0      0      0      0      0      0      1
149      0      0      0      0      0      0      1

[150 rows x 63 columns]

```