

S-109A Introduction to Data Science

Harvard University Summer 2018

Instructors: Pavlos Protopapas, Kevin Rader

Prepared by Group #28

Mark Dinneen Mohammad Karim Ramandeep Harjai

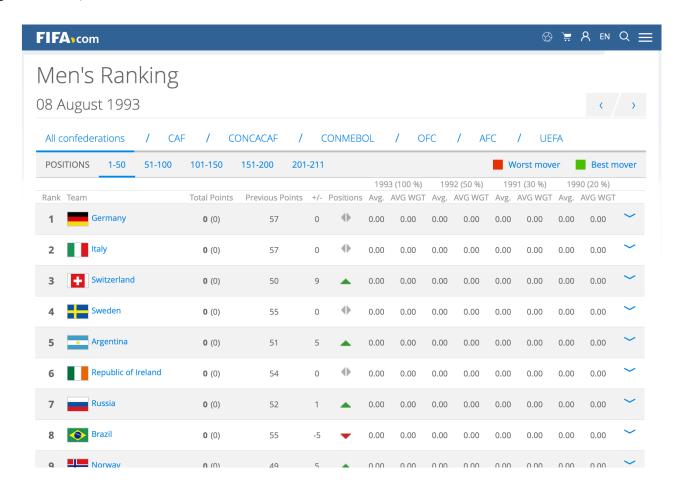
Supplemental Notebook for Web Scraping

This notebook documents all the data which has been collected via web scraping technique

```
In [8]: # import the necessary libraries
        %matplotlib inline
        import numpy as np
        import scipy as sp
        import matplotlib as mpl
        import matplotlib.cm as cm
        import matplotlib.pyplot as plt
        import pandas as pd
        from pandas import Series
        import random
        # libraries for scraping data from web
        import re
        from bs4 import BeautifulSoup
        from sys import argv
        import requests
        from urllib.request import urlopen
        from urllib.error import HTTPError
```

Historical FIFA Ranking Data

The latest men's FIFA ranking data is available at: https://www.fifa.com/fifa-world-ranking/ranking-table/men/index.html), which is published on 07–June–2018. The oldest men's FIFA ranking data is available at: https://www.fifa.com/fifa-world-ranking/ranking-table/men/rank=2/index.html (https://www.fifa.com/fifa-world-ranking/ranking-table/men/rank=2/index.html), which was published on 08–August–1993. We plan to parse this data from 288+ pages, which will provide us the FIFA ranking for various men's teams, for time period between 08–August–1993, and 07–June–2018.



FIFA ranking page for 08-August-1993

```
In [6]: def get_month_number(month_name: str) -> str:
    month_number = 0
    if month_name.lower() == "january":
        month_number = "01"
    elif month_name.lower() == "february":
        month_number = "02"
```

```
elif month name.lower() == "march":
        month number = "03"
    elif month name.lower() == "april":
        month number = "04"
    elif month name.lower() == "may":
        month_number = "05"
    elif month name.lower() == "june":
        month number = "06"
    elif month name.lower() == "july":
        month number = "07"
    elif month name.lower() == "august":
        month number = "08"
    elif month name.lower() == "september":
        month number = "09"
    elif month name.lower() == "october":
        month number = "10"
    elif month name.lower() == "november":
        month number = "11"
    elif month name.lower() == "december":
        month number = "12"
    return month number
def get year(rank date: str) -> str:
    rank date split = rank date.split(" ")
    return rank date split[2]
def format date(rank date: str) -> str:
    rank date split = rank date.split(" ")
    rank date formatted = rank date split[2] + "-" + \
                            get month number(rank date split[1]) + "-"
+ \
                            rank date split[0]
    return rank date formatted
def parse rank data() -> []:
    ranking = []
    page url = "https://www.fifa.com/fifa-world-ranking/ranking-table/
men/rank={}/index.html"
   page_ids = range(2, 288)
    for page id in page ids:
        progress = round( (page id / max(page ids))*100 )
        print("Parsing FIFA ranking data: {}%\r".format(progress), end
="", flush=True)
        rank page = requests.get(page url.format(page id))
        page soup = BeautifulSoup(rank_page.text, 'html.parser')
        rank date = page soup.findAll("ul", {"class":["slider-list","i
tems-1"]})[0].text
        rank soup = page soup.findAll("tr", {"class":"anchor"})
        for rank tag in rank soup:
            date = format date(rank date)
            year = get year(rank date)
            rank = rank tag.find("td", {"class":"tbl-rank"}).text
```

 $df_ranking = (57793, 4)$

Out[6]:

	date	year	team	rank
29624	2007-03-14	2007	Romania	14
2466	1994-12-20	1994	Malawi	82
13262	2000-04-12	2000	Samoa	183
14642	2000-11-15	2000	St Kitts and Nevis	146
53579	2016-11-24	2016	France	7
2533	1994-12-20	1994	Guatemala	149
22041	2004-01-14	2004	Croatia	20
32770	2008-06-04	2008	Tunisia	50
41035	2011-11-23	2011	Hungary	37
15411	2001-03-14	2001	Syria	103

Historical FIFA World Cup Winners

The historical FIFA world-cup finals data is available at:

https://en.wikipedia.org/wiki/List_of_FIFA World Cup_finals

(https://en.wikipedia.org/wiki/List_of_FIFA_World_Cup_finals), which is last published (updated) on 25-July-2018. We'll scrape this data to build a dataset of historical FIFA world cup winners.

		List of fir	nals matches, their ve	enues and locations, the finalist	ts, and final scores		
Year ¢	Winners ♦	Final score ^[2] ♦	Runners-up +	Venue	Location +	Attendance +	References
1930	Uruguay 🔙	4–2	Argentina	Estadio Centenario	Montevideo, Uruguay	80,000	[7][8]
1934	Italy **	2-1 [†] [n 3]	Czechoslovakia	Stadio Nazionale PNF	Rome, Italy	50,000	[9][10]
1938	Italy ■ ■	4–2	Hungary	Stade Olympique de Colombes	Paris, France	45,000	[11][12]
1950 ^[n 4]	Uruguay 🔙	2–1 [n 5]	Brazil	Estádio do Maracanã	Rio de Janeiro, Brazil	199,854 ^[13]	[14][15]
1954	West Germany	3–2	T Hungary	Wankdorf Stadium	Bern, Switzerland	60,000	[16][17]
1958	Brazil 📀	5–2	Sweden	Råsunda Stadium	Solna, Sweden	51,800	[18][19]
1962	Brazil 🔷	3–1	Czechoslovakia	Estadio Nacional	Santiago, Chile	69,000	[20][21]
1966	England +	4-2† [n 6]	West Germany	Wembley Stadium	London, England	93,000	[22][23]
1970	Brazil 📀	4–1	■ Italy	Estadio Azteca	Mexico City, Mexico	107,412	[24][25]
1974	West Germany	2–1	Netherlands	Olympiastadion	Munich, West Germany	75,200	[26][27]
1978	Argentina ==	rgentina 3-1† Netherlands Estadio Monumental Buenos A		Buenos Aires, Argentina	71,483	[28][29]	
1982	Italy 🛮 🖠	3–1	West Germany	Santiago Bernabéu	Madrid, Spain	90,000	[30][31]
1986	Argentina ===	3–2	West Germany	Estadio Azteca	Mexico City, Mexico	114,600	[32][33]
1990	West Germany	1–0	Argentina	Stadio Olimpico	Rome, Italy	73,603	[34][35]
1994	Brazil 📀	0-0‡ [n 8]	■ Italy	Rose Bowl	Pasadena, United States	94,194	[36][37]
1998	France	3–0	♦ Brazil	Stade de France	Saint-Denis, France	80,000	[38][39]
2002	Brazil 📀	2–0	Germany	International Stadium	Yokohama, Japan	69,029	[40][41]
2006	Italy 🛮 🖠	1—1‡ [n 9]	France	Olympiastadion	Berlin, Germany	69,000	[42][43]
2010	Spain	1-0 [†] [n 10]	Netherlands	Soccer City	Johannesburg, South Africa	84,490	[44][45]
2014	Germany ===	1-0 [†] [n 11]	Argentina	Estádio do Maracanã Rio de Janeiro, Brazil		74,738	[46][47]
2018	France	4–2	Croatia	Luzhniki Stadium	Moscow, Russia	78,011	
				Upcoming finals			
Year	Team 1	ν	Team 2	Venue	Location	Attendance	Reference
2022				Lusail Iconic Stadium	Lusail, Qatar		
2026				TBD	North America		

Historical FIFA World Cup Final Match Results

```
row soup = table soup.findAll("tr")
    results = []
    for row in range(1,len(row soup)-4):
        result = []
        result.append(re.findall(r"[0-9]+", row soup[row].find("th").t
ext)[0])
        result.append(row soup[row].findAll("td")[0].find("a").text)
        result.append(row soup[row].findAll("td")[2].find("a").text)
        result.append(re.findall(r"[0-9]+", row soup[row].findAll("td"
)[1].find("a").text)[0])
        result.append(re.findall(r"[0-9]+", row soup[row].findAll("td"
)[1].find("a").text)[1])
        result.append(row soup[row].findAll("td")[0].find("a").text)
        result.append(row_soup[row].findAll("td")[3].findAll("a")[0].t
ext)
        result.append(row soup[row].findAll("td")[4].findAll("a")[0].t
ext)
        result.append(row soup[row].findAll("td")[4].findAll("a")[1].t
ext)
        result.append(re.findall(r"[0-9,]+", row soup[row].findAll("td
")[5].text)[0].replace(',',''))
        results.append(result)
    return results
# un-comment the lines below run data scraping
# colnames = ['year', 'team1', 'team2', 'team1 score', 'team2 score',
#
              'winner','venue','city','country','attendance']
# df fifa finals = pd.DataFrame.from records(parse past winners(), col
umns=colnames)
# df fifa finals["year"] = df fifa finals["year"].astype(int)
# df fifa finals["team1"] = df fifa finals["team1"].astype(str)
# df fifa finals["team2"] = df fifa finals["team2"].astype(str)
# df fifa finals["team1 score"] = df fifa finals["team1 score"].astype
(int)
# df fifa finals["team2 score"] = df fifa finals["team2 score"].astype
(int)
# df fifa finals["winner"] = df fifa finals["winner"].astype(str)
# df fifa finals["venue"] = df fifa finals["venue"].astype(str)
# df fifa finals["city"] = df fifa finals["city"].astype(str)
# df fifa finals["country"] = df fifa finals["country"].astype(str)
# df fifa finals["attendance"] = df fifa finals["attendance"].astype(i
nt)
# df fifa finals = df fifa finals.replace("West Germany", "Germany")
# df_fifa_finals.to_pickle("datasets/fifa_finals_data.pkl") # store da
taframe to local disk
# read back web scraped dataset
df fifa finals = pd.read pickle("datasets/fifa finals data.pkl")
print("\ndf fifa finals = ", df fifa finals.shape, "\n")
df fifa finals
```

Out[103]:

	year	team1	team2	team1_score	team2_score	winner	ven
0	1930	Uruguay	Argentina	4	2	Uruguay	Estadio Centenario
1	1934	Italy	Czechoslovakia	2	1	Italy	Stadio Nazionale PN
2	1938	Italy	Hungary	4	2	Italy	Stade Olympique de Colombes
3	1950	Uruguay	Brazil	2	1	Uruguay	Estádio do Maracanã
4	1954	Germany	Hungary	3	2	Germany	Wankdorf Stadium
5	1958	Brazil	Sweden	5	2	Brazil	Råsunda Stadium
6	1962	Brazil	Czechoslovakia	3	1	Brazil	Estadio Nacional
7	1966	England	Germany	4	2	England	Wembley Stadium
8	1970	Brazil	Italy	4	1	Brazil	Estadio Aztec
9	1974	Germany	Netherlands	2	1	Germany	Olympiastadi
10	1978	Argentina	Netherlands	3	1	Argentina	Estadio Monumental
11	1982	Italy	Germany	3	1	Italy	Santiago Bernabéu
12	1986	Argentina	Germany	3	2	Argentina	Estadio Azteo
13	1990	Germany	Argentina	1	0	Germany	Stadio Olimpico
14	1994	Brazil	Italy	0 0 Brazil		Rose Bowl	
15	1998	France	Brazil 3 0 France		France	Stade de France	
16	2002	Brazil	Germany	2	0	Brazil	International Stadium
17	2006	Italy	France	1	1	Italy	Olympiastadi

18	2010	Spain	Netherlands	1	0	Spain	Soccer City
19	2014	Germany	Argentina	1	0	Germany	Estádio do Maracanã
20	2018	France	Croatia	4	2	France	Luzhniki Stadium

Out[119]:

year
5
4
4
2
2
2
1
1

FIFA World Cup — All Time Team Rankings

FIFA all time team rankings, and associated team statistics is available at https://www.fifa.com/fifa-tournaments/statistics-and-records/worldcup/teams/index.html). We'll scrape this data to build a dataset of FIFA team rankings.

FIFA _v com							⊗ `	≓ Α EI	v Q ≡
FIFA World Cup TM Teams Statistics Victories All-time rankings Participations Matches Cards KIA OFFICIAL MATCH BALL CARRIER LINE MATCH BALL CARRIER									
FIFA World Cup™ - All-tii	me rankings								
RANK TEAM	PTS	MP	W	D	L	GS	GA	AV. PTS	APPS.
1 BRAZIL	227	104	70	17	17	221	102	2.2	20
2 GERMANY	218	106	66	20	20	224	121	2.1	18
3 ITALY	156	83	45	21	17	128	77	1.9	18
4 ARGENTINA	140	77	42	14	21	131	84	1.8	16
5 SPAIN	99	59	29	12	18	92	66	1.7	14
6 ENGLAND	98	62	26	20	16	79	56	1.6	14

FIFA World Cup - All Time Team Rankings

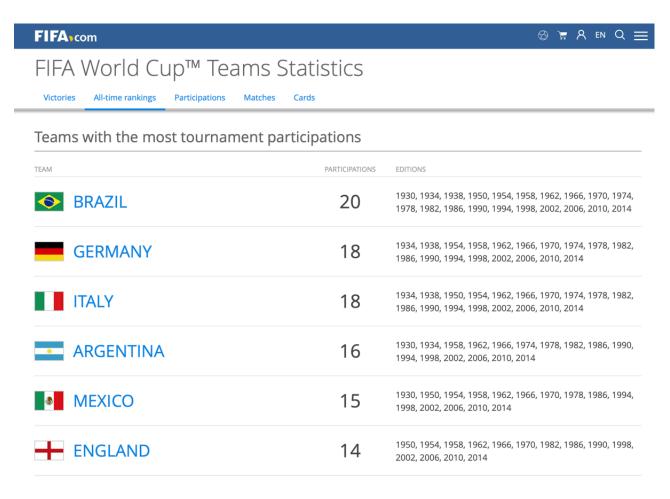
```
ext"}).text)
        result.append(row.findAll("td")[3].find("span", {"class":"text
"}).text)
        result.append(row.findAll("td")[4].find("span", {"class":"text
"}).text)
        result.append(row.findAll("td")[5].find("span", {"class":"text
"}).text)
        result.append(row.findAll("td")[6].find("span", {"class":"text
"}).text)
        result.append(row.findAll("td")[7].find("span", {"class":"text
"}).text)
        result.append(row.findAll("td")[8].find("span", {"class":"text
"}).text)
        result.append(row.findAll("td")[9].find("span", {"class":"text
"}).text)
        result.append(row.findAll("td")[10].text)
        result.append(row.findAll("td")[11].find("span", {"class":"tex
t"}).text)
        results.append(result)
    return results
# un-comment the lines below run data scraping
# colnames = ['rank','team','points','matches','win','draw','lost',
              'goal for', 'goal against', 'points avg', 'appearances']
# df fifa wc ranking = pd.DataFrame.from records(parse team rankings()
, columns=colnames)
# df fifa wc ranking["rank"] = df fifa wc ranking["rank"].astype(int)
# df fifa_wc_ranking["points"] = df_fifa_wc_ranking["points"].astype(i
nt)
# df_fifa_wc_ranking["matches"] = df_fifa_wc_ranking["matches"].astype
(int)
# df fifa wc ranking["win"] = df fifa wc ranking["win"].astype(int)
# df fifa wc ranking["draw"] = df fifa wc ranking["draw"].astype(int)
# df fifa wc ranking["lost"] = df fifa wc ranking["lost"].astype(int)
# df fifa wc ranking["goal for"] = df fifa wc ranking["goal for"].asty
pe(int)
# df fifa wc ranking["goal against"] = df fifa wc ranking["goal agains
t"].astype(int)
# df fifa wc ranking["points avg"] = df fifa wc ranking["points avg"].
astype(float)
# df fifa wc ranking["appearances"] = df fifa wc ranking["appearances"
].astype(int)
# df fifa wc ranking.to pickle("datasets/fifa wc ranking.pkl") # store
dataframe to local disk
# read back web scraped dataset
df fifa wc ranking = pd.read pickle("datasets/fifa wc ranking.pkl")
print("\ndf fifa wc ranking = ", df fifa wc ranking.shape, "\n")
df fifa wc ranking.head(10)
```

Out[141]:

	rank	team	points	matches	win	draw	lost	goal_for	goal_against	points_a
0	1	Brazil	227	104	70	17	17	221	102	2.2
1	2	Germany	218	106	66	20	20	224	121	2.1
2	3	Italy	156	83	45	21	17	128	77	1.9
3	4	Argentina	140	77	42	14	21	131	84	1.8
4	5	Spain	99	59	29	12	18	92	66	1.7
5	6	England	98	62	26	20	16	79	56	1.6
6	7	France	96	59	28	12	19	106	71	1.6
7	8	Netherlands	93	50	27	12	11	86	48	1.9
8	9	Uruguay	72	51	20	12	19	80	71	1.4
9	10	Sweden	61	46	16	13	17	74	69	1.3

FIFA World Cup — Participations

FIFA World Cup - team participations, and associated team statistics is available at https://www.fifa.com/fifa-tournaments/statistics-and-records/worldcup/teams/index.html (https://www.fifa.com/fifa-tournaments/statistics-and-records/worldcup/teams/index.html). We'll scrape this data to build a dataset of FIFA team participations in world cup tournament.



FIFA World Cup — Team Participations

```
def parse team participations() -> []:
In [147]:
              url = "https://www.fifa.com/fifa-tournaments/statistics-and-record
          s/worldcup/teams/index.html"
              page = requests.get(url)
              page soup = BeautifulSoup(page.text, 'html.parser')
              table_soup = page_soup.findAll("table", {"class":["table","tbl-all
          timeranking"]})[2].find("tbody")
              row soup = table soup.findAll("tr")
              results = []
              for row in row soup:
                  result = []
                  result.append(row.findAll("td")[0].find("span", {"class":"t-nT
          ext"}).text)
                  result.append(row.findAll("td")[2].find("span", {"class":"text
          "}).text)
                  result.append(row.findAll("td")[3].find("span", {"class":"text
          "}).text)
                  results.append(result)
              return results
          # un-comment the lines below run data scraping
          # colnames = ['team', 'participations', 'years']
          # df fifa participations = pd.DataFrame.from records(parse team partic
          ipations(), columns=colnames)
          # df fifa participations["participations"] = df_fifa_participations["p
          articipations"].astype(int)
          # df_fifa_participations.to_pickle("datasets/fifa_participations.pkl")
          # store dataframe to local disk
          # read back web scraped dataset
          df fifa participations = pd.read pickle("datasets/fifa participations.
          pkl")
          print("\ndf fifa participations = ", df fifa participations.shape, "\n
          ")
          df fifa participations.head(10)
```

Out[147]:

	team	participations	years
0	Brazil	20	1930, 1934, 1938, 1950, 1954, 1958, 1962, 1966
1	Germany	18	1934, 1938, 1954, 1958, 1962, 1966, 1970, 1974
2	Italy	18	1934, 1938, 1950, 1954, 1962, 1966, 1970, 1974
3	Argentina	16	1930, 1934, 1958, 1962, 1966, 1974, 1978, 1982
4	Mexico	15	1930, 1950, 1954, 1958, 1962, 1966, 1970, 1978
5	England	14	1950, 1954, 1958, 1962, 1966, 1970, 1982, 1986
6	France	14	1930, 1934, 1938, 1954, 1958, 1966, 1978, 1982
7	Spain	14	1934, 1950, 1962, 1966, 1978, 1982, 1986, 1990
8	Belgium	12	1930, 1934, 1938, 1954, 1970, 1982, 1986, 1990
9	Uruguay	12	1930, 1950, 1954, 1962, 1966, 1970, 1974, 1986

```
In [120]: from IPython.core.display import HTML
    def css_styling(): styles = open("cs109.css", "r").read(); return HTML
        (styles)
        css_styling()
```

Out[120]: