

```

# D:\Dropbox\Doug Eagan\Weight\soccer players scraper.py
01| import pandas as pd
02| import re
03| import requests
04| from bs4 import BeautifulSoup
05|
06| # Get basic players information for all players
07| base_url = "https://sofifa.com/players?offset="
08| columns = ['ID', 'Name', 'Age', 'Photo', 'Nationality', 'Flag', 'Overall',
'Potential', 'Club', 'Club Logo', 'Value', 'Wage', 'Special']
09| data = pd.DataFrame(columns = columns)
10|
11| for offset in range(0, 290):
12|     url = base_url + str(offset * 61)
13|     source_code = requests.get(url)
14|     plain_text = source_code.text
15|     soup = BeautifulSoup(plain_text, 'html.parser')
16|     table_body = soup.find('tbody')
17|     for row in table_body.findAll('tr'):
18|         td = row.findAll('td')
19|         picture = td[0].find('img').get('data-src')
20|         pid = td[0].find('img').get('id')
21|         nationality = td[1].find('a').get('title')
22|         flag_img = td[1].find('img').get('data-src')
23|         name = td[1].findAll('a')[1].text
24|         age = td[2].find('div').text.strip()
25|         overall = td[3].text.strip()
26|         potential = td[4].text.strip()
27|         club = td[5].find('a').text
28|         club_logo = td[5].find('img').get('data-src')
29|         value = td[6].text.strip()
30|         wage = td[7].text.strip()
31|         special = td[8].text.strip()
32|         player_data = pd.DataFrame([[pid, name, age, picture, nationality,
flag_img, overall, potential, club, club_logo, value, wage, special]])
33|         player_data.columns = columns
34|         data = data.append(player_data, ignore_index=True)
35| data = data.drop_duplicates()
36|
37| # Get detailed player information from player page
38| detailed_columns = ['Preferred Foot', 'International Reputation', 'Weak Foot',
'Skill Moves', 'Work Rate', 'Body Type', 'Real Face', 'Position', 'Jersey Number',
'Joined', 'Loaned From', 'Contract Valid Until', 'Birthdate', 'Height', 'Weight',
'LS', 'ST', 'RS', 'LW', 'LF', 'CF', 'RF', 'RW', 'LAM', 'CAM', 'RAM', 'LM', 'LCM',
'CM', 'RCM', 'RM', 'LWB', 'LDM', 'CDM', 'RDM', 'RWB', 'LB', 'LCB', 'CB', 'RCB', 'RB',
'Crossing', 'Finishing', 'HeadingAccuracy', 'ShortPassing', 'Volleys', 'Dribbling',
'Curve', 'FKAccuracy', 'LongPassing', 'BallControl', 'Acceleration', 'SprintSpeed',
'Agility', 'Reactions', 'Balance', 'ShotPower', 'Jumping', 'Stamina', 'Strength',
'LongShots', 'Aggression', 'Interceptions', 'Positioning', 'Vision', 'Penalties',
'Composure', 'Marking', 'StandingTackle', 'SlidingTackle', 'GKDividing', 'GKHandling',
'GKkicking', 'GKPositioning', 'GKReflexes', 'ID']
39| detailed_data = pd.DataFrame(index = range(0, data.count()[0]), columns =
detailed_columns)
40| detailed_data.ID = data.ID.values
41|
42| player_data_url = 'https://sofifa.com/player/'
43| for id in data.ID:
44|     url = player_data_url + str(id)
45|     source_code = requests.get(url)
46|     plain_text = source_code.text
47|     soup = BeautifulSoup(plain_text, 'html.parser')
48|     skill_map = {}
49|     columns = soup.find('div', {'class': 'teams'}).find('div', {'class':
'columns'}).findAll('div', {'class': 'column col-4'})
50|     for column in columns:
51|         skills = column.findAll('li')
52|         for skill in skills:

```

```

53|         if(skill.find('label') != None):
54|             label = skill.find('label').text
55|             value = skill.text.replace(label, '').strip()
56|             skill_map[label] = value
57|     meta_data = soup.find('div', {'class': 'meta'}).text.split(' ')
58|     length = len(meta_data)
59|     birthyear = meta_data[length - 3]
60|     birthday = meta_data[length - 4]
61|     birthmonth = meta_data[length - 5]
62|     weight = meta_data[length - 1]
63|     height = meta_data[length - 2].split('\')[0] + '\' + meta_data[length -
64| 2].split('\')[1].split('\')[0]
65|     skill_map["Birthdate"] = birthmonth+"/"+birthday+"/"+birthyear
66|     skill_map["Height"] = height
67|     skill_map["Weight"] = weight
68|     if('Position' in skill_map.keys()):
69|         if skill_map['Position'] in ('', 'RES', 'SUB'):
70|             skill_map['Position'] = soup.find('article').find('div', {'class':
71| 'meta'}).find('span').text
72|             if(skill_map['Position'] != 'GK'):
73|                 card_rows = soup.find('aside').find('div', {'class': 'card
74| mb-2'}).find('div', {'class': 'card-body'}).findAll('div', {'class': 'columns'})
75|                 for c_row in card_rows:
76|                     attributes = c_row.findAll('div', {'class': re.compile('column
77| col-sm-2 text-center')})
78|                     for attribute in attributes:
79|                         if(attribute.find('div')):
80|                             name = ''.join(re.findall('[a-zA-Z]', attribute.text))
81|                             value = attribute.text.replace(name, '').strip()
82|                             skill_map[str(name)] = value
83|     sections = soup.find('article').findAll('div', {'class': 'mb-2'})[1:3]
84|     first = sections[0].findAll('div', {'class': 'column col-4'})
85|     second = sections[1].findAll('div', {'class': 'column col-4'})[:-1]
86|     sections = first + second
87|     for section in sections:
88|         items = section.find('ul').findAll('li')
89|         for item in items:
90|             value = int(re.findall(r'\d+', item.text)[0])
91|             name = ''.join(re.findall('[a-zA-Z]*', item.text))
92|             skill_map[str(name)] = value
93|     for key, value in skill_map.items():
94|         detailed_data.loc[detailed_data.ID == id, key] = value
95|
96| full_data = pd.merge(data, detailed_data, how = 'inner', on = 'ID')
97| full_data.to_csv('data.csv', encoding='utf-8-sig')

```