

Python Web Scraping with Selenium

Peter Njoroge Mwangi

September 19, 2024

1 Introduction

This document explains how to use Python for web scraping using Selenium. We extract football match data (date, home team, score, and away team) from a website and store the data in a CSV file for analysis. The target website is <https://www.adamchoi.co.uk/overs/detailed> Adam Choi's football stats page.

2 Required Libraries

We use the following libraries for this project:

- **Selenium**: For web browser automation.
- **pandas**: For data manipulation and saving the scraped data.
- **webdriver_manager**: For automatically handling ChromeDriver installation.

Install these libraries using:

```
1 pip install selenium pandas webdriver_manager
```

3 Code Breakdown

3.1 Importing Required Modules

We start by importing the necessary modules for Selenium to control the browser and pandas for data processing.

```

1 import pandas as pd
2 from selenium import webdriver
3 from selenium.webdriver.chrome.service import Service
4 from selenium.webdriver.common.by import By
5 from selenium.webdriver.support.ui import Select
6 from webdriver_manager.chrome import ChromeDriverManager

```

3.2 Setting up the WebDriver

Here, we use `ChromeDriverManager` to handle `ChromeDriver` installation. This simplifies the process of setting up a compatible driver for Selenium.

```

1 driver = webdriver.Chrome(service=Service(ChromeDriverManager
    ().install()))

```

3.3 Navigating to the Website

We use Selenium's `get()` function to open the desired website.

```

1 website = 'https://www.adamchoi.co.uk/overs/detailed'
2 driver.get(website)

```

3.4 Interacting with the Website

We locate and click the "All matches" button and select the 2024/2025 football season from a dropdown menu using XPath.

```

1 all_matches_button = driver.find_element(By.XPATH, "//label[
    contains(text(), 'All matches')"])
2 all_matches_button.click()
3
4 try:
5     season_dropdown = driver.find_element(By.XPATH, "//select
6     [@id='season']")
7     select = Select(season_dropdown)
8     select.select_by_visible_text('2024/2025')
9 except Exception as e:
10    print(f"Error selecting season: {e}")

```

3.5 Scraping Match Data

We use the following code to locate table rows (<tr>) and extract relevant details for each match.

```
1 matches = driver.find_elements(By.TAG_NAME, 'tr')
2
3 date = []
4 home_team = []
5 score = []
6 away_team = []
7
8 for match in matches:
9     try:
10         cells = match.find_elements(By.TAG_NAME, 'td')
11         if len(cells) == 4:
12             date.append(cells[0].text)
13             home_team.append(cells[1].text)
14             score.append(cells[2].text)
15             away_team.append(cells[3].text)
16     except Exception as e:
17         print(f"Error processing row: {e}")
```

3.6 Cleaning and Saving the Data

Here, we clean the score column by replacing dashes with colons to avoid Excel misinterpreting them as dates, and save the data to a CSV file.

```
1 df = pd.DataFrame({
2     'Date': date,
3     'home_team': home_team,
4     'score': score,
5     'away_team': away_team
6 })
7
8 df['score'] = df['score'].apply(lambda x: x.replace('-', ':'))
9 df.to_csv('Football_data.csv', index=False)
```

3.7 Closing the Browser

Finally, we close the browser after the scraping process is complete.

```
1 driver.quit()
```

4 Additional Information

- **Error Handling:** We use `try-except` blocks to manage potential issues while interacting with the webpage or processing rows.
- **XPath and CSS Selectors:** XPath is used to target elements. In case the website's structure changes, update the XPaths accordingly.
- **Score Formatting:** By replacing dashes with colons in the score column, we prevent Excel from interpreting certain scores as dates.

This script can be extended to scrape additional data or adapted for use on other websites.

5 Full Code

```
1 import pandas as pd
2 from selenium import webdriver
3 from selenium.webdriver.chrome.service import Service
4 from selenium.webdriver.common.by import By
5 from selenium.webdriver.support.ui import Select
6 from webdriver_manager.chrome import ChromeDriverManager
7
8 # Define the website to scrape
9 website = 'https://www.adamchoi.co.uk/overs/detailed'
10
11 # Use ChromeDriverManager to automatically download and
12 # manage the correct ChromeDriver version
13 driver = webdriver.Chrome(service=Service(ChromeDriverManager
14                              ()).install()))
15
16 # Open the website
17 driver.get(website)
18
19 # Click the "All matches" button to display all match data
20 all_matches_button = driver.find_element(By.XPATH, "//label[
21     contains(text(), 'All matches')]")
22 all_matches_button.click()
23
24 # Select the 2024/2025 season from the dropdown
25 try:
26     # Locate the season dropdown
```

```

24     season_dropdown = driver.find_element(By.XPATH, "//select
    [id='season']")
25
26     # Create a Select object for interacting with the
    dropdown
27     select = Select(season_dropdown)
28
29     # Select the '2024/2025' season by its visible text
30     select.select_by_visible_text('2024/2025')
31
32 except Exception as e:
33     print(f"Error selecting season: {e}")
34
35 # Find all match rows by the 'tr' tag name
36 matches = driver.find_elements(By.TAG_NAME, 'tr')
37
38 # Initialize lists to hold the extracted data
39 date = []
40 home_team = []
41 score = []
42 away_team = []
43
44 # Loop through the table rows to extract match data
45 for match in matches:
46     try:
47         # Find all 'td' elements (columns) in the current row
48         cells = match.find_elements(By.TAG_NAME, 'td')
49
50         # Ensure that the row has exactly 4 columns (date,
51         home team, score, away team)
52         if len(cells) == 4:
53             # Append each piece of data to its corresponding
54             list
55             date.append(cells[0].text)
56             home_team.append(cells[1].text)
57             score.append(cells[2].text)
58             away_team.append(cells[3].text)
59
60             # Print each row for debugging
61             print(f>Date: {cells[0].text}, Home: {cells[1].
62             text}, Score: {cells[2].text}, Away: {cells[3].text}")
63
64     except Exception as e:
65         print(f"Error processing row: {e}")

```

```

64 # Close the browser after extracting the data
65 driver.quit()
66
67 # Create a pandas DataFrame from the extracted data
68 df = pd.DataFrame({
69     'Date': date,
70     'home_team': home_team,
71     'score': score,
72     'away_team': away_team
73 })
74
75 # Replace hyphens in the 'score' column to avoid Excel
    interpreting scores as dates
76 df['score'] = df['score'].apply(lambda x: x.replace('-', ':'))
77
78 # Save the DataFrame to a CSV file
79 df.to_csv('Football_data.csv', index=False)
80
81 # Print the DataFrame
82 print(df)

```