Python Web Scraper Application

Phase 2 Source

Autumn Capasso, Michael Galyen, Kacey Gambill,

Eric Viera, Anthony Washington

Department of Computer Science, University of Maryland Global Campus

CMSC-495

Professor Janak Rajani

July 26, 2022

# Table of Contents

## Phase 2 Plan

The plan for Phase 2 was to add the following functionalities:
- Performance tests for all functions.
- Migrate Unittest tests to Pytest.
- UI styling.
- Functions to scrape FreeCodeCamp for resume and interview advice (text and links).

## Additions to the Project

In this phase of the project, the team added functions to app.py to scrape FreeCodeCamp.org for resume and interview advice, started adding more styling to the user interface, and added more performance tests to the test suite. The acceptance tests in the use guide were updated to match the current UI of the application.

Figure 1 image of the user interface

Figure 2 image of the updated links function in WebScraper.py

```python
30    # This block scrapes links
31    def links(url: str):
32        text_dict = {
33                "site": [
34                    ]
35                }
36        r = requests.get(url)
37        linksoup = BeautifulSoup(r.content, 'html.parser')
38    ##    print(linksoup.a.prettify())
39        for link_index, link in enumerate(linksoup.find_all('a')):
40            # set limit at 100 so it doesn't take forever
41            if link_index == 100:
42                break
43            link_href = link.get('href')
44            if str(link_href).startswith('http'):
45    #            print(link_href)
46                try:
47                    text_dict_info = text(link_href)
48                    text_dict['site'].append(text_dict_info)
49                except:
50                    print('did not scrape')
51        return text_dict['site']
52
```

Figure 3 image of the newWebScraper.py function, text

```python
53     # This block scrapes the text
54 def text(url: str) -> dict:
55         r = requests.get(url)
56         textsoup = BeautifulSoup(r.content, 'html.parser')
57         title_text = textsoup.find_all("title")
58         text_body = textsoup.find_all('p')
59
60
61 #     print(textsoup.title.prettify())
62     try:
63         text_dict_info = {
64                 "url": url,
65                 "title": title_text,
66                 "p": text_body
67                 }
68
69     except:
70         print('no info found')
71     return text_dict_info
72
```

Figure 4 image of benchmark test of the app.py index page function

```python
23     # Test flask app index page function performance
24 def test_index_page_performance(client):
25     start = time.time()
26     response = client.get('/')
27     if response.status_code == 200:
28         stop = time.time()
29     else:
30         print(response)
31     total_time = stop - start
32     print(f"Total time to execute index_page function was {total_time}")
33
```

Figure 5 image of benchmark test of the app.py jobs page function

```python
41      # Test flask app index page function performance
42    ⊟def test_jobs_page_performance(client):
43          start = time.time()
44          response = client.get('/jobs')
45          if response.status_code == 200:
46              stop = time.time()
47          else:
48              print(response)
49          total_time = stop - start
50          print(f"Total time to execute jobs_page function was {total_time}")
```

Figure 6 image of benchmark test of the app.py resume page function

```python
59      # Test flask app index page function performance
60    ⊟def test_resume_page_performance(client):
61          start = time.time()
62          response = client.get('/resume')
63          if response.status_code == 200:
64              stop = time.time()
65          else:
66              print(response)
67          total_time = stop - start
68          print(f"Total time to execute resume_page function was {total_time}")
```

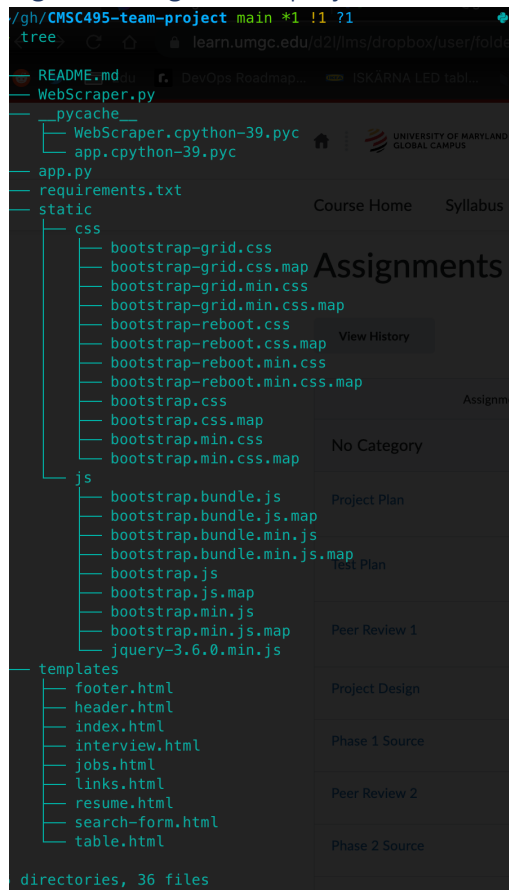Figure 7 image of benchmark test of the app.py interview page function

```python
77      # Test flask app index page function performance
78    ⊟def test_interview_page_performance(client):
79          start = time.time()
80          response = client.get('/interview')
81          if response.status_code == 200:
82              stop = time.time()
83          else:
84              print(response)
85          total_time = stop - start
86          print(f"Total time to execute interview_page function was {total_time}")
```

## Current Project Structure

The structure generally didn't change in a significant way from Phase 1 to Phase 2.

Figure 8 image of the project structure



## Challenges

Migrating tests from Unittest to Pytest and adding styles to the UI were both proposed for this week. More research and development time was needed in terms of the font end work so that will be represented better in phase 3. Migrating tests from Unittest to Pytest was an unexpected challenge so that will be pushed to phase 3 as well.

While functions were added to scrape FreeCodeCamp for links and text in the categories, resume tips and interview tips, it was discovered that these links of FreeCodeCamp are likely being generated by JavaScript so more time will be needed to figure out how to get the intended links. The text being returned by these scrapes is also a little messy with HTML tags mixed in with the text so that will also be dealt with in the next phase.

## Milestones

June 28 – Finalize project plan and specifications
July 5 – Web scraper is built and ready to test
July 5 – Test plan is finalized

July 12 – Project design is finalized
July 12 – Flask application is finalized
July 19 – Web scraper is integrated with the Flask framework
July 19 – User guide is finalized
July 19 – Alternative design write up (if needed)
July 26 – Submit test results
August 2 – Finalize history write up
August 9 – Finalize documentation

## Milestones Review

According to the milestones listed in our original project plan document, all testing with results was supposed to be finalized by this phase. More time is needed for this so it will get pushed to phase 3. The user guide was finalized so that is up to date. The web scraper has been integrated with the Flask app. The general structure and main functions of the Flask app have been finalized. The project design has been finalized. The project plan and specifications have been finalized.

# Conclusions

While this sprint was productive, some things will carry over to the next sprint. UI styling, test migration, and better parsing will all be pushed to Phase 3.