

BASF Working Student Assignment

Shankar Anbalagan

shankanb09@gmail.com

Source code to my solution: <https://github.com/ShankarAnbalagan/basf-coding-challenge>

Website 1

Scrapy Spider

```
website1_spider.py x
website1 > website1_spider.py > Website1Spider > parse_project
1 import scrapy
2 from urllib.parse import urlparse, parse_qs
3 import time
4
5 import website1_scrape
6 from utils import save_data
7
8 class Website1Spider(scrapy.Spider):
9     name = "website1"
10
11     def start_requests(self):
12         urls = list((f"https://www.uvp-verbund.de/freitextsuche?rstart=0&currentSelectorPage={page_number}" for page_number in range(1, 6)))
13         for url in urls:
14             yield scrapy.Request(url=url, callback=self.parse)
15             time.sleep(5)
16
17     def parse(self, response):
18         parsed_url = urlparse(response.url)
19         base_url = f"{parsed_url.scheme}://{parsed_url.netloc}"
20         project_links = website1_scrape.get_project_links(base_url, response.body)
21         for index, link in enumerate(project_links):
22             yield scrapy.Request(url=link, callback=self.parse_project, meta={'index': index})
23             time.sleep(5)
24
25     def parse_project(self, response):
26         project_index = response.meta['index']
27
28         project_url = response.url
29         parsed_url = urlparse(response.url)
30         base_url = f"{parsed_url.scheme}://{parsed_url.netloc}"
31         page_num = parse_qs(parsed_url.query)['currentSelectorPage'][0]
32
33         project_data, project_html = website1_scrape.extract_data(base_url, response.body)
34         save_data.save('website1\\output', project_data, project_html, page_num, project_index)
```

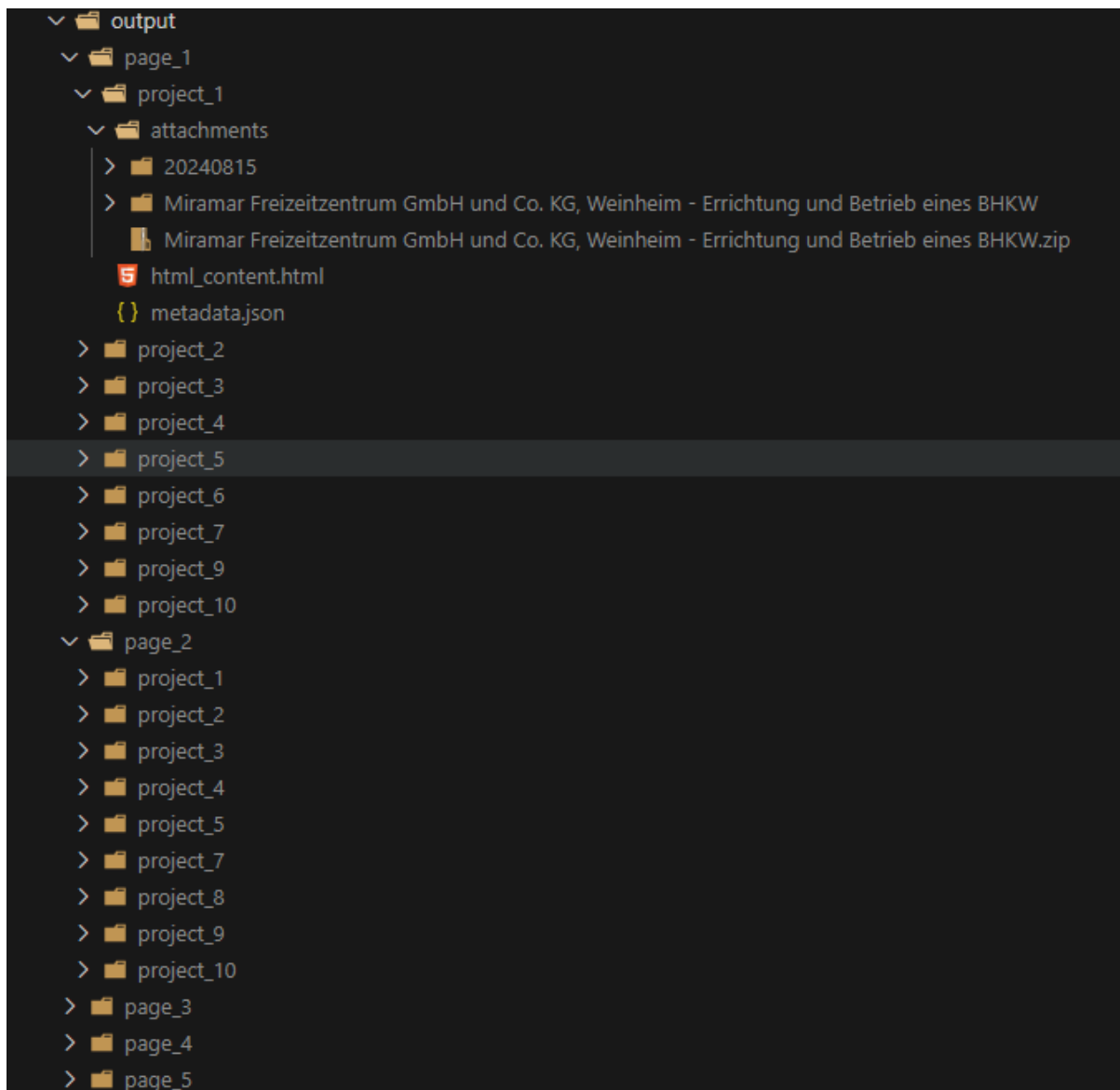
BeautifulSoup

```
website1_scraper.py M X
website1 > website1_scraper.py > extract_data
1 from bs4 import BeautifulSoup
2
3 def get_project_links(base_url, html):
4     soup = BeautifulSoup(html, features="lxml")
5
6     results_div = soup.find('div', id='results')
7     data_divs = results_div.find_all('div', class_='data')
8
9     links = []
10    for data_div in data_divs:
11        links.append(f'{base_url}{data_div.find("div").find("a")["href"]}')
12    print(links)
13    return links
14
15 def extract_data(base_url, html):
16     data = {}
17     soup = BeautifulSoup(html, features="lxml")
18
19     data['title'] = soup.title.string
20
21     description_title_element = soup.find('h3', string='Allgemeine Vorhabenbeschreibung')
22     description = description_title_element.find_next_sibling('p').string
23     data['description'] = description
24
25     date_title_element = soup.find('h4', string='Datum der Entscheidung')
26     data['date'] = date_title_element.find_next_sibling('p').string if not date_title_element else ''
27
28     download_link = soup.find('a', title='Alle Dokumente als ZIP-Datei herunterladen')['href']
29     data['download_link'] = download_link
30
31     print(data)
32     return data, soup.prettify()
```

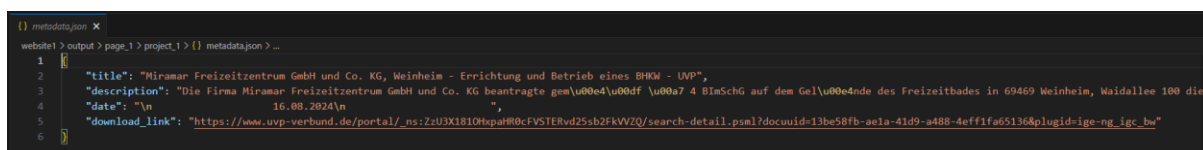
Utility to save metadata and download attachments

```
save_data.py X
website1 > utils > save_data.py > ...
1 import os
2 import json
3 import urllib.request
4 import zipfile
5
6 def save(output_dir, data, html, page_num, project_index):
7     if not os.path.exists(output_dir):
8         os.makedirs(output_dir)
9
10    page_dir = os.path.join(output_dir, f'page_{page_num}')
11    project_dir = os.path.join(page_dir, f'project_{project_index + 1}')
12
13    if not os.path.exists(page_dir):
14        os.makedirs(page_dir)
15
16    if not os.path.exists(project_dir):
17        os.makedirs(project_dir)
18
19    with open(os.path.join(project_dir, 'metadata.json'), 'w') as file:
20        json.dump(data, file, indent=4)
21
22    with open(os.path.join(project_dir, 'html_content.html'), 'w') as file:
23        file.write(html)
24
25    if data['download_link']:
26        attachment_dir = os.path.join(project_dir, 'attachments')
27        if not os.path.exists(attachment_dir):
28            os.makedirs(attachment_dir)
29
30        download_file_path = ''
31        with urllib.request.urlopen(data['download_link']) as response:
32            content_disposition = response.headers.get('Content-Disposition')
33            if content_disposition:
34                filename = content_disposition.split('filename=')[1].strip('"')
35                download_file_path = os.path.join(attachment_dir, filename)
36
37                if len(download_file_path) > 256:
38                    extra_characters = len(download_file_path) - 256
39                    filename = filename.strip('.zip')[:-extra_characters] + '.zip'
40                    download_file_path = os.path.join(attachment_dir, filename)
41
42                with open(download_file_path, 'wb') as file:
43                    file.write(response.read())
44
45        if '.zip' in download_file_path:
46            os.makedirs(download_file_path.strip('.zip'))
47            with zipfile.ZipFile(download_file_path, 'r') as zip_file:
48                zip_file.extractall(attachment_dir)
```

Output folder structure



Metadata.json



Observations/Challenges:

After initial inspection of the html behind the web page, I decided to implement the task using scrapy and BeautifulSoup. This website blocks very frequent requests which prompted me to add a delay between requests. Unfortunately, this slows down the complete extraction process. When downloading attachments, some failed initially because the name of the zip file would be greater than 256 characters (character limit for path in windows os), so I had to trim the name of zip file in case the path character count was more than 256.

Website 2

Main.py

```
main.py x
website2 > main.py > main
1 from SeleniumDriver import SeleniumDriver
2 from utils import save_data
3
4 def main():
5     url = r'https://edl.doe.gov.my/discovery?query=+where+category%3d+%2527EIA+Report%2527&category=EIA%20Report&main=Digital'
6     selenium_driver = SeleniumDriver()
7     selenium_driver.open(url)
8
9     for page in range(1,6):
10         selenium_driver.navigate_to_page_number(page)
11         projects = selenium_driver.get_projects()
12
13         for index, project in enumerate(projects):
14             selenium_driver.navigate_to_project(index)
15             data, pretty_html = selenium_driver.extract_data()
16             save_data.save('website2/output', data, pretty_html, page, index)
17             selenium_driver.open(url)
18
19     selenium_driver.close()
20
21
22 if __name__ == "__main__":
23     main()
```

SeleniumDriver

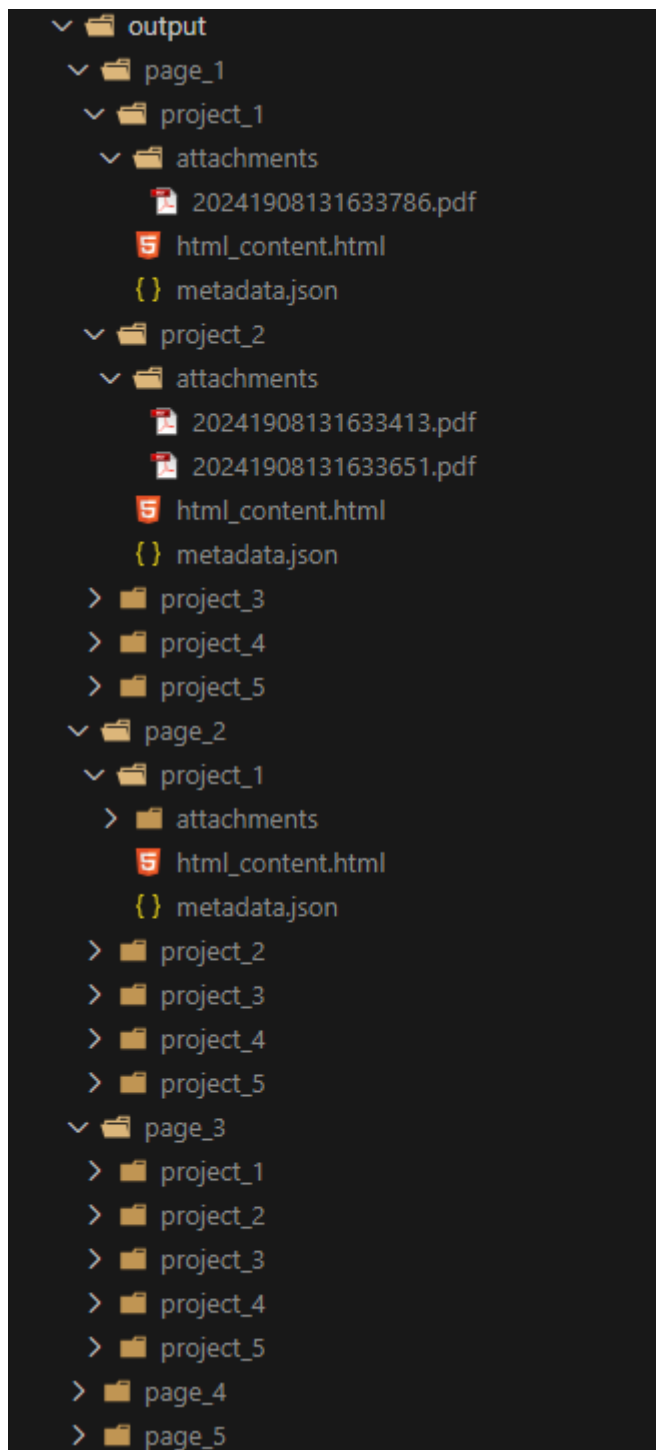
```
SeleniumDriver.py M X
website2 > SeleniumDriver.py > ...
1 from selenium import webdriver
2 from selenium.webdriver.common.by import By
3 from selenium.webdriver.support.wait import WebDriverWait
4 from selenium.webdriver.support import expected_conditions
5 from bs4 import BeautifulSoup
6 from urllib.parse import urlparse, parse_qs
7
8 class SeleniumDriver:
9     def __init__(self):
10         chrome_options = webdriver.ChromeOptions()
11         chrome_options.add_experimental_option("detach", True)
12         service = webdriver.ChromeService(executable_path=r'website2\chromedriver.exe')
13         self._driver = webdriver.Chrome(service=service, options=chrome_options)
14
15     def close(self):
16         self._driver.close()
17
18     def open(self, url):
19         self._driver.get(url)
20
21     def navigate_to_page_number(self, page_num):
22         id = f'content_body_rptPaging_lbPaging_{page_num - 1}'
23         WebDriverWait(self._driver, 30).until(expected_conditions.presence_of_element_located((By.ID, id)))
24         next_page_a = self._driver.find_element(By.ID, id)
25         self._driver.execute_script('arguments[0].click();', next_page_a)
26
27     def navigate_to_project(self, project_index):
28         print(f'content_body_rptResult_lbtnOverview2_{project_index}')
29         element = self._driver.find_element(By.ID, f'content_body_rptResult_lbtnOverview2_{project_index}')
30         self._driver.execute_script('arguments[0].click();', element)
31
32     def get_projects(self):
33         project_elements = self._driver.find_elements(By.CLASS_NAME, 'lbltitle_')
34         return project_elements
35
36     def get_page_html(self):
37         return self._driver.page_source
38
39     def _get_data_by_id(self, soup, id):
40         return soup.find('textarea', id=id).string
41
42     def extract_data(self):
43         WebDriverWait(self._driver, 30).until(expected_conditions.presence_of_element_located((By.ID, 'content_body_txtTitle')))
44
45         data = {}
46         soup = BeautifulSoup(self._driver.page_source, features="lxml")
47
48         data['title'] = self._get_data_by_id(soup, 'content_body_txtTitle')
49         data['category'] = self._get_data_by_id(soup, 'content_body_txtsubcategory')
50         data['release_date'] = self._get_data_by_id(soup, 'content_body_txtreleased_eia')
51         data['consultant'] = self._get_data_by_id(soup, 'content_body_txtconsultant_eia')
52         data['project_developer'] = self._get_data_by_id(soup, 'content_body_txtconsultant_eia')
53         data['location'] = self._get_data_by_id(soup, 'content_body_txtlocation_eia')
54         data['post_date'] = soup.find('input', id='content_body_txtdatepost_eia')['value']
55
56         # https://edl.doe.gov.my/images/items/42681/attachment/20241908131633413.pdf
57         page_url = self._driver.current_url
58         parsed_url = urlparse(page_url)
59         download_base_url = f"{parsed_url.scheme}://{parsed_url.netloc}/images/items/"
60         bib_number = parse_qs(parsed_url.query)['bibnumber'][0]
61
62         download_links = []
63         download_a_tags = soup.find_all('a', class_='lbtn')
64         for tag in download_a_tags:
65             download_links.append(f'{download_base_url}{bib_number}/attachment/{tag.string}')
66         data['download_links'] = download_links
67
68         print(data)
69         return data, soup.prettify()
70
```

Utility to save metadata and download attachments

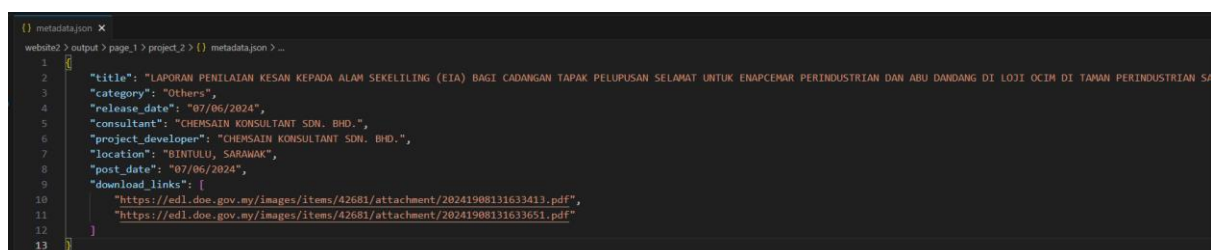
```
save_data.py x
website2 > utils > save_data.py > save

1 import os
2 import json
3 import urllib.request
4 import zipfile
5
6 def save(output_dir, data, html, page_num, project_index):
7     if not os.path.exists(output_dir):
8         os.makedirs(output_dir)
9
10    page_dir = os.path.join(output_dir, f'page_{page_num}')
11    project_dir = os.path.join(page_dir, f'project_{project_index + 1}')
12
13    if not os.path.exists(page_dir):
14        os.makedirs(page_dir)
15
16    if not os.path.exists(project_dir):
17        os.makedirs(project_dir)
18
19    with open(os.path.join(project_dir, 'metadata.json'), 'w') as file:
20        json.dump(data, file, indent=4)
21
22    with open(os.path.join(project_dir, 'html_content.html'), 'w') as file:
23        file.write(html)
24
25    if(data['download_links']):
26        attachment_dir = os.path.join(project_dir, 'attachments')
27        if not os.path.exists(attachment_dir):
28            os.makedirs(attachment_dir)
29
30        for download_link in data['download_links']:
31            download_file_path = ''
32            with urllib.request.urlopen(download_link) as response:
33                filename = download_link.split('/')[-1]
34                download_file_path = os.path.join(attachment_dir, filename)
35
36            with open(download_file_path, 'wb') as file:
37                file.write(response.read())
```

Output folder structure



Metadata.json



Observations/Challenges:

On inspection on the html code behind website 2, I realized it is a ASP.NET web page. This made navigation and actions difficult with scrapy, so I decided to use Selenium for this. Downloading attachments was a challenge too as the download link was not directly available.