

Task Web Scrapping

Import Library and Web:

```
1 import requests
2 from bs4 import BeautifulSoup
```

Python

```
1 soup = BeautifulSoup(response.text, 'html.parser')
```

Python

```
1 url = 'https://www.baraasallout.com/test.html'
2 response = requests.get(url)
3 print(response.text)
```

Python

Task 1

- **Extract Text Data:**
 - Extract all headings (<h1>, <h2>).
 - Extract all text content inside <p> and tags.
 - Savethis data into a Extract_Text_Data.CSV file.
 - <https://www.pythontutorial.net/python-basics/python-write-csv-file/>

```
1 headers_1 = soup.find_all('h1')
2 headers_1
```

[20]

Python

```
... [<h1>Web Scrapping Practice</h1>]
```

```
1 headers_2 = soup.find_all('h2')
2 headers_2
```

[21]

Python

```
... [<h2>Available Products</h2>,
<h2>Product Table</h2>,
<h2>Watch This Video</h2>,
<h2>Contact Us</h2>,
<h2>Product Information</h2>,
<h2>Featured Products</h2>]
```

```

1 paragraph = soup.find('p')
2 paragraph

[9] Python

... <p>Welcome to the web scraping task! Use your skills to extract the required data from this page.</p>

1 all_p = soup.find_all('p')
2 all_p

[10] Python

... [<p>Welcome to the web scraping task! Use your skills to extract the required data from this page.</p>,
<p><strong>Sharp Objects</strong></p>,
<p style="color: green;">£47.82</p>,
<p style="color: green;">✓ In stock</p>,
<p><strong>In a Dark, Dark Wood</strong></p>,
<p style="color: green;">£19.63</p>,
<p style="color: green;">✓ In stock</p>,
<p><strong>The Past Never Ends</strong></p>,
<p style="color: green;">£56.50</p>,
<p style="color: green;">✓ In stock</p>,
<p><strong>A Murder in Time</strong></p>,
<p style="color: green;">£16.64</p>,
<p style="color: green;">Out stock</p>,
<p class="name">Wireless Headphones</p>,
<p class="price" style="display: none;">$49.99</p>,

```

```

1 first_ul = soup.find('ul')
2

[25] Python

1 all_ul_li = first_ul.find_all('li')
2 all_ul_li

[26] Python

... [<li class="highlight">Laptop</li>,
<li>Smartphone</li>,
<li>Tablet</li>,
<li>Smartwatch</li>]

1 for li in all_ul_li:
2     print(li.get_text())

[27] Python

... Laptop
Smartphone
Tablet
Smartwatch

```

```

24
25 # Prepare the data
26 data = (
27     [('Heading', heading) for heading in headings] +
28     [('Paragraph', paragraph) for paragraph in all_p] +
29     [('List Item', list_item) for list_item in all_ul_li]
30 )
31
32 # Convert the data into a DataFrame
33 df = pd.DataFrame(data, columns=['Type', 'Content'])
34
35 csv_file = 'Extract_Text_Data.csv'
36 df.to_csv(csv_file, index=False, encoding='utf-8')
37
38 print(f"Data has been saved to {csv_file}.")
39

```

Task 2

- Extract Table Data:
 - Extract data from the table, including:
 - Product Name.
 - Price.
 - StockStatus.
 - Savethis data into a Extract_Table_Data.CSV file.
 - <https://www.pythontutorial.net/python-basics/python-write-csv-fil>

```
1 tables = soup.find_all('table')
2 tables

[32] Python

... [<table>
  <tr>
    <th>Product</th>
    <th>Price</th>
    <th>In Stock</th>
  </tr>
  <tr>
    <td>Laptop</td>
    <td>$1000</td>
    <td>Yes</td>
  </tr>
  <tr>
    <td>Smartphone</td>
    <td>$800</td>
    <td>No</td>
  </tr>
  <tr>
    <td>Tablet</td>
    <td>$500</td>
    <td>Yes</td>
  </tr>
</table>]
```

```
1 all_divs = soup.find_all('div')
2 all_divs

[29] Python

... [<div style="display: flex; justify-content: space-around; margin-top: 20px;">
  <div style="text-align: center; width: 200px; border: 1px solid #ddd; padding: 10px; border-radius: 5px;">
    
    <p><strong>Sharp Objects</strong></p>
    <p style="color: green;">£47.82</p>
    <p style="color: green;">✓ In stock</p>
    <button style="background-color: blue; color: white; border: none; padding: 10px 15px; border-radius: 5px; cursor: pointer;">Add to basket</button>
  </div>
  <div style="text-align: center; width: 200px; border: 1px solid #ddd; padding: 10px; border-radius: 5px;">
    
    <p><strong>In a Dark, Dark Wood</strong></p>
    <p style="color: green;">£19.63</p>
    <p style="color: green;">✓ In stock</p>
    <button style="background-color: blue; color: white; border: none; padding: 10px 15px; border-radius: 5px; cursor: pointer;">Add to basket</button>
  </div>
  <div style="text-align: center; width: 200px; border: 1px solid #ddd; padding: 10px; border-radius: 5px;">
    <img alt="The Past Never Ends" href="http://books.toscrape.com/media/cache/c0/59/c05972805aa7201171b8fc71a5b00292.jpg" style="width: 100%; height: auto; border-radius: 5px;"/>
    <p><strong>The Past Never Ends</strong></p>
    <p style="color: green;">£56.50</p>
    <p style="color: green;">✓ In stock</p>
    <button style="background-color: blue; color: white; border: none; padding: 10px 15px; border-radius: 5px; cursor: pointer;">Add to basket</button>
  </div>
  <div style="text-align: center; width: 200px; border: 1px solid #ddd; padding: 10px; border-radius: 5px;">
    
    <p><strong>A Murder in Time</strong></p>
  </div>
</div>]
```

```

25
26 # Prepare the data for saving to CSV
27 data = {
28     'Table Data': table_data,
29     'BG Yellow Elements': bg_yellow_elements,
30     'DIV Content': all_divs
31 }
32
33 # Adjust for different lengths (pad with None)
34 max_length = max(len(data['Table Data']), len(data['BG Yellow Elements']), len(data['DIV Content']))
35 data['Table Data'] += [None] * (max_length - len(data['Table Data']))
36 data['BG Yellow Elements'] += [None] * (max_length - len(data['BG Yellow Elements']))
37 data['DIV Content'] += [None] * (max_length - len(data['DIV Content']))
38
39 # Create a DataFrame
40 df = pd.DataFrame({
41     'Table Data': ['; '.join(row) if row else None for row in data['Table Data']], # Join table rows
42     'BG Yellow Elements': data['BG Yellow Elements'],
43     'DIV Content': data['DIV Content']
44 })
45
46 csv_file = 'Extract_Table_Data.csv'
47 df.to_csv(csv_file, index=False, encoding='utf-8')
48
49 print(f"Data has been saved to {csv_file}.")

```

Task 3

- **Extract Product Information (Cards Section):**
 - Extract data from the book cards at the bottom of the page, including:
 - BookTitle.
 - Price.
 - StockAvailability.
 - Button text (e.g., "Add to basket").
 - Savethedata into a Product_Information.JSON file.
 - <https://www.geeksforgeeks.org/how-to-convert-python-dic/>

```

1 book_cards = soup.find_all('div', style=lambda value: value and 'text-align: center' in value)
2
3 books_data = []
4 for card in book_cards:
5     title = card.find('p').strong.text
6     price = card.find('p', style=lambda value: value and 'color: green' in value).text
7     stock = card.find_all('p', style=lambda value: value and 'color: green' in value)[-1].text.strip()
8     button_text = card.find('button').text.strip()
9
10    books_data.append({
11        'BookTitle': title,
12        'Price': price,
13        'StockAvailability': stock,
14        'ButtonText': button_text
15    })
16
17 for book in books_data:
18     print(book)

```

```

{'BookTitle': 'Sharp Objects', 'Price': '£47.82', 'StockAvailability': '✓ In stock', 'ButtonText': 'Add to basket'}
{'BookTitle': 'In a Dark, Dark Wood', 'Price': '£19.63', 'StockAvailability': '✓ In stock', 'ButtonText': 'Add to basket'}
{'BookTitle': 'The Past Never Ends', 'Price': '£56.50', 'StockAvailability': '✓ In stock', 'ButtonText': 'Add to basket'}
{'BookTitle': 'A Murder in Time', 'Price': '£16.64', 'StockAvailability': 'Out stock', 'ButtonText': 'Add to basket'}

```

```

29
30     # Append extracted information to the products list
31     products.append({
32         "BookTitle": title,
33         "Price": price,
34         "StockAvailability": stock_availability,
35         "ButtonText": button
36     })
37
38     json_file = 'Product_Information.json'
39     with open(json_file, 'w', encoding='utf-8') as file:
40         json.dump(products, file, indent=4, ensure_ascii=False)
41
42     print(f"Product data has been saved to {json_file}.")
43

```

Task 4

- **Extract Form Details:**
 - Extract all input fields from the form, including:
 - Field name (e.g., username, password).
 - Input type (e.g., text, password, checkbox, etc.).
 - Default values, if any.
 - Save the data into a JSON file.
 - <https://www.geeksforgeeks.org/how-to-convert-python-dictionary-to-json/>

```
1 forms = soup.find_all('form')
2 forms

[51] Python

... [<form>
<label for="username">Username:</label>
<input id="username" name="username" placeholder="Enter your username" type="text"/>
<label for="password">Password:</label>
<input id="password" name="password" placeholder="Enter your password" type="password"/>
<label for="options">Choose an option:</label>
<select id="options" name="options">
<option value="option1">Option 1</option>
<option value="option2">Option 2</option>
<option value="option3">Option 3</option>
</select>
<label>
<input name="terms" type="checkbox"/> I agree to the terms and conditions
</label>
<input type="submit" value="Submit"/>
</form>]
```

```
1 # Extract input fields
2 input_fields = first_form.find_all('input')
3 input_fields

[57] Pyt

... [<input id="username" name="username" placeholder="Enter your username" type="text"/>,
<input id="password" name="password" placeholder="Enter your password" type="password"/>,
<input name="terms" type="checkbox"/>,
<input type="submit" value="Submit"/>]

1 first_form.get_text()

[54] Pyt

... '\nUsername:\n\nPassword:\n\nChoose an option:\n\nOption 1\nOption 2\nOption 3\n\n\n I agree to the terms and conditions\n\n\n'
```

+ Code + Markdown

```
1 fields_data = []
2
3 for field in input_fields:
4     field_name = field.get('name', 'N/A')
5     input_type = field.get('type', 'text')
6     default_value = field.get('value', 'N/A')
7     fields_data.append({
8         'Field Name': field_name,
9         'Input Type': input_type,
10        'Default Value': default_value
11    })
12
13 for field in fields_data:
14     print(field)
```

```
{'Field Name': 'username', 'Input Type': 'text', 'Default Value': 'N/A'}
{'Field Name': 'password', 'Input Type': 'password', 'Default Value': 'N/A'}
{'Field Name': 'terms', 'Input Type': 'checkbox', 'Default Value': 'N/A'}
{'Field Name': 'N/A', 'Input Type': 'submit', 'Default Value': 'Submit'}
```

```
15
16 for form in forms:
17     inputs = form.find_all('input')
18     input_details = []
19     for input_tag in inputs:
20         field_name = input_tag.get('name', 'N/A') # Field name
21         input_type = input_tag.get('type', 'text') # Input type
22         default_value = input_tag.get('value', None) # Default value
23         input_details.append({
24             "FieldName": field_name,
25             "InputType": input_type,
26             "DefaultValue": default_value
27         })
28     form_details.append(input_details)
29
30 form_json_file = 'Form_Details.json'
31 with open(form_json_file, 'w', encoding='utf-8') as file:
32     json.dump(form_details, file, indent=4, ensure_ascii=False)
33
```

Task 5

- Extract Links and Multimedia:
 - Extract the hyperlink (<a> tag) and its href value.
 - Extract the video link from the <iframe> tag.
 - Savethedata into a JSON file.
 - <https://www.geeksforgeeks.org/how-to-convert-python-dictionary-to-json/>

```
1 links = soup.find_all('a')
2 links

[73] Python
... []
```

```
1 first_link.get_text()
2 first_link.attrs # Not Existing attributes

Python

... -----
AttributeError                                Traceback (most recent call last)
Cell In[34], <a href='vscode-notebook-cell:?execution_count=34&line=1'>line 1</a>
----> <a href='vscode-notebook-cell:?execution_count=34&line=1'>1</a> first_link.get_text()
      <a href='vscode-notebook-cell:?execution_count=34&line=2'>2</a> first_link.attrs

AttributeError: 'NoneType' object has no attribute 'get_text'

Python

1 first_link.attrs['href'] # Not Existing attributes

Python

... -----
AttributeError                                Traceback (most recent call last)
Cell In[36], <a href='vscode-notebook-cell:?execution_count=36&line=1'>line 1</a>
----> <a href='vscode-notebook-cell:?execution_count=36&line=1'>1</a> first_link.attrs['href']

AttributeError: 'NoneType' object has no attribute 'attrs'
```

>>>>> No <a> in the site

- Extract the video link from the iframe tag.

```
1 iframe = soup.find('iframe')
2 iframe

[79] Python
... <iframe height="315" src="https://www.youtube.com/watch?v=ujf9RNU8dCU" width="560"></iframe>
```

```
8
9 # Prepare the data
10 multimedia_data = {
11     "Hyperlinks": hyperlinks,
12     "Videos": videos
13 }
14
15 links_json_file = 'Links_and_Multimedia.json'
16 with open(links_json_file, 'w', encoding='utf-8') as file:
17     json.dump(multimedia_data, file, indent=4, ensure_ascii=False)
18
```


Task 6

- Scraping Challenge

- Product Name: Located within ``.

- HiddenPrice: Located within ``, which has `style="display: none;"`.

- Available Colors: Located within ``.

- Product ID: The value stored in the data-id attribute.

- ExampleOutput:

```
[  
  {'id': '101', 'name': 'Wireless Headphones', 'price': '$49.99', 'colors': 'Black, White, Blue'},  
  {'id': '102', 'name': 'Smart Speaker', 'price': '$89.99', 'colors': 'Grey, Black'},  
  {'id': '103', 'name': 'Smart Watch', 'price': '$149.99', 'colors': 'Black, Silver, Gold'}  
]
```

```
1 products = []  
2 for card in product_cards:  
3     product_id = card.get("data-id")  
4     name = card.find("p", class_="name").text  
5     price = card.find("p", class_="price").text.strip() # Extract hidden price  
6     colors = card.find("p", class_="colors").text.replace("Available colors: ", "").strip()  
7     products.append({  
8         "id": product_id,  
9         "name": name,  
10        "price": price,  
11        "colors": colors  
12    })  
13  
14 for product in products:  
15     print(product)
```

```
{'id': '101', 'name': 'Wireless Headphones', 'price': '$49.99', 'colors': 'Black, White, Blue'}  
{'id': '102', 'name': 'Smart Speaker', 'price': '$89.99', 'colors': 'Grey, Black'}  
{'id': '103', 'name': 'Smart Watch', 'price': '$149.99', 'colors': 'Black, Silver, Gold'}
```

Thanks