

My name - Vaishali Ramdas S41VC

Batch No - D52311

Sir = Subrata Mishra

Date = 30-1-2024

PAGE NO.	
DATE	/ /

Assignment - 2

Web-scraping

Q1

Write a Python program to scrape data for "data analyst" job position in "Bangalore" location. You have to scrape the job-title, job-location, company name, experience, - required - you have to scrape first 10 jobs data.

This task will be done in following steps:

1. First get the webpage <https://www.sine.com/>
2. Enter "data Analyst" in 'Job title, skills' field and enter "Bangalore" in "enter the location" field.
3. Then click the search button.
4. Then scrape the data for the first 10 jobs results you get.
5. Finally create a DataFrame of the scraped data.

⇒

To achieve this task, you can use the requests library to send HTTP requests to the website and BeautifulSoup for parsing the HTML Content, additionally you may use the 'pandas' library to create a DataFrame for the scraped data. Here's a sample Python program to get you started.

python

copy code

```
import requests
from bs4 import BeautifulSoup
import pandas as pd

# Step 1: Get the webpage
url = "https://www.shine.com/"
response = requests.get(url)

# Check if the request was successful
if response.status_code == 200:
    print("Successfully fetched the webpage.")
else:
    print(f"Failed to fetch the webpage. Status code: {response.status_code}")
    exit()

# Step 2: Enter search criteria and click the search button
search_criteria = {
    'q': 'Data Analyst',
    'l': 'Bangalore'
}

# Search - URL = https://www.shine.com/job-search/dm-analyst-jobs-in-bangalore
search_url = "https://www.shine.com/job-search/dm-analyst-jobs-in-bangalore"
search_response = requests.get(search_url)

# If search-response status code == 200:
if search_response.status_code == 200:
    print("Successfully performed the search.")
else:
    print(f"Failed to perform the search. Status code: {search_response.status_code}")
    exit()

# Step 3: Scrape the data for the first 10 jobs
soup = BeautifulSoup(search_response.text, 'html.parser')
```

```
job_listings = soup.find('ul', class_ = 'search-listing')
data = []
```

```
for job in job_listings[1:10]:
```

```
    title = job.find('h3', class_ = 'job-title').text.strip()
```

```
    location = job.find('span', class_ = 'job-location').text.strip()
```

```
    company_name = job.find('span', class_ = 'job-comp-name')  
    (text.strip())
```

```
    experience_required = job.find('span', class_ = 'job-exp')  
    (text.strip())
```

```
    data.append({
```

```
        'Job Title': title,
```

```
        'Job Location': location,
```

```
        'Company Name': company_name,
```

```
        'Experience Required': experience_required}
```

3)

```
# step 4: create a Data Frame
```

```
df = pd.DataFrame(data)
```

```
# Display the Data Frame
```

```
print(df)
```

Note

The code above assumes the structure of the HTML elements on the Shine.com website as of my last knowledge update in January 2022. Website structure may change, so you may need to inspect the HTML of the website and adjust the code accordingly if there are any changes.

Q2

write a python program to scrape data for "data scientist" job position in "Bangalore" location, you have to scrape the job-title, job-location, company name. you have to scrape first 10 jobs data. This task will be done in following steps:

- 1 First get the web page <https://www.shine.com>
- 2 Enter Data "Scientist" in "job title", skills field and "Bangalore" in "enter the location" field.
- 3 Then click the search button.

Then scrape the data for the first 10 jobs result you get.

Finally create a data frame.

→ To achieve achieve this task, you can follow a similar approach as in the previous example. I'll use the 'requests' library for sending HTTP requests, 'Beautiful soup' for parsing the HTML content, and 'pandas' to create a Data Frame. Here's sample python program:

Python

import requests

from bs4 import BeautifulSoup

import pandas as pd

A copy code

```
# Step 1: Get the web page
url = "https://www.shine.com"
response = requests.get(url)
# Check if the request was successful
if response.status_code == 200:
    print("Successfully fetched the webpage")
else:
```

print ("F" printed to fetch the webpages. status code
 if response.status_code == 200):
 exit ()

step 2 & Enter search criteria and click the
 'Search' Button

search-data = {

'q': 'Data scientist',

'j': 'Bangalore',

3

search-url = "https://www.shine.com/job-search/
 data-scientist-jobs-in-Bangalore"

search-response = requests.get(search-url, params=search
 data)

if search-response.status_code == 200:

print ("Successfully performed the search.")

else:

print ("F" failed to perform the search, status code: (

if search-response.status_code == 3):

exit ()

step 3: Scrape the data for the first 10 jobs

soup = BeautifulSoup(search-response.text, 'html.parser')

job-listing = soup.find('div', class_ = 'search-listing')

data = []

for job in job-listing[:10]:

title = job.find('h3', class_ = 'job-title').text.strip()

location = job.find('span', class_ = 'job-location').text.strip()

company-name = job.find('span', class_ = 'job-name').text.strip()

data.append({

'Job Title': title,

'Job Location': location,

'Company Name': company-name}

3)

Step 4 is create a Data Frame

$df = pd.DataFrame(datal)$

Display the Data Frame

`print(df)`

As mentioned earlier, the structure of the HTML elements may change over time, so make sure to inspect the HTML of the website and adjust the code accordingly if there are any changes.

Q3

In this question you have to scrape data using the filters available on the webpage you have to use the location and salary filter.

You have to scrape data for "Data Scientist" designation for first 10 job results.

You have to scrape the job - title, job - location, company name, experience required.

The location filter to be used is "Delhi/NCR" the salary filter to be used is '3-6' lakhs

The task will be page <https://www.shine.com/> First get the web page <https://www.shine.com/>

Enter "Data scientist" in skill designation and company filter.

Then click the search button.

Then apply the location filter and salary filter by checking the respective boxes.

Then scrape the data for the first 10 job results you get.

Finally create a DataFrame of the scraped data.

→ To achieve this task you can the request library to send HTTP request to the website and BeautifulSoup for parsing the HTML content. Additionally, you may use the pandas library to create a data frame for the scraped data. Here's a sample Python program

Python

copy code

```

import requests
from bs4 import BeautifulSoup
import pandas as pd

# Step 1: Get the webpage
url = "https://www.shine.com/"
response = requests.get(url)

# check if the request was successful
if response.status_code == 200:
    print("Successfully fetched the webpage")
else:
    print(f"Failed to fetch the webpage. status code {response.status_code}")

# Step 2: Enter search criteria and click the search button
search_data = {
    'q': 'Data Scientist'
}

Search_url = "https://www.shine.com/job-search"
Search_response = requests.get(Search_url, params=search_data)

if Search_response.status_code == 200:
    print("Successfully performed the search")
else:
    print(f"Failed to perform the search. status code {Search_response.status_code}")

exit()

```

Step 2 : Enter search criteria and click the search button
search-data = {
 "q": "Data Scientist",
 "location": "Delhi/NCR",
 "minSalary": 300000,
 "maxSalary": 600000,
 "experience": "Entry level",
 "industry": "Technology",
 "function": "Software Development",
 "employmentType": "Full-time",
 "sortOrder": "Relevance",
 "page": 1, "perPage": 10}

'q': 'Data Scientist'

2

search-url = "https://www.shine.com/job-search/jobs?query=Data-Scientist-jobs"

search-response = requests.get(search-url, params=search-data)

if search-response.status_code == 200:

print("Successfully performed the search.")

else:

print("Failed to perform the search. Status code: ", search-response.status_code)

exit()

Step 3 : Apply location and salary filters.

filters-data = {
 "location": "Delhi/NCR",
 "minSalary": 300000, "maxSalary": 600000, "experience": "Entry level",
 "industry": "Technology", "function": "Software Development",
 "employmentType": "Full-time", "sortOrder": "Relevance", "page": 1, "perPage": 10}

"location": "Delhi/NCR",

"salary": "3-6-lakhs"

3

filter-url = "https://www.shine.com/job-search/filter-jobs?job=Data-Scientist-jobs"

filter-response = requests.post(filter-url, params=search-data)

data = filter-response.json()

if filter-response.status_code == 200:

print("Successfully applied filters")

else:

print("Failed to apply filters. Status code: ", filter-response.status_code)

exit()

Step 4 : Scrape the data for the first 10 jobs

soup = BeautifulSoup(filter-response.text, "lxml")

job-listing = soup.find_all("div", class_="search-listing")

data = []

```
for job in job_listing[3:10]:  
    title = job.find('h3', class_='job-title').text.strip()  
    location = job.find('span', class_='job-location').text.strip()  
    company_name = job.find('span', class_='job-comp-name').text.strip()  
    experience_required = job.find('span', class_='job-exp').text.strip()  
  
    data.append({  
        'Job Title': title,  
        'Job Location': location,  
        'Company Name': company_name,  
        'Experience Required': experience_required  
    })  
  
# Step 5: Create a Data Frame  
df = pd.DataFrame(data)  
  
# Display the Data Frame  
print(df)
```

As always, keep in mind that website structures may change. So you may need to inspect the HTML of the website and adjust the code accordingly if there are any changes.

Q 4 Scrape duty of first 100 Sunglasses listing on flipkart.com
you have to scrape four attributes.

Brand (7) product description (8) price.

The attributes which you have to scrape is ticked marked in the below image.

To scrape the data you have to go through following

- 1 Go to FLIPKART webpage by url: <https://www.flipkart.com>
- 2 Enter "Sunglasses" in the search field where "Search" box products brand and more? is written and click the search icon.

- 3 After that you will reach to the page having a lot of sunglasses. From this page you can scrap required data as usual.
- 4 After scraping data from the first page, go to the 'next' button at the bottom of the page then click on it.

- 5 Now scrape data from this page as usual.
Repeat this until you get a total of 100 sunglasses.

→ To achieve this task, you can use requests' library to send HTTP requests to the website, BeautifulSoup for parsing the HTML content, and pandas to create a data frame for the scraped data. Additionally, you might need to handle pagination by navigating through the "next button". Here's a simple Python program.

Pythoncopy codeimport requests:from bs4 import BeautifulSoupimport pandas as pd

Function to scrape sunglasses data from a given flipkart page.

def scrape_flipkart_sunglasses(url, count):response = requests.get(url)if response.status_code == 200:print("P" "B" "Scraping data from page " "(count" "))else:print("P" "Failed to fetch the web page. Status code:" " "(response.status_code)"")return Nonesoup = BeautifulSoup(response.text, "html.parser")sunglasses_listings = soup.find_all("div", class_="IAtub")data = [{} for _ in range(len(sunglasses_listings))]for sunglasses in sunglasses_listings:brand = sunglasses.find("div", class_="zWkNRV").text.strip()description = sunglasses.find("div", class_="IRwTq").text.strip()price = sunglasses.find("div", class_="Zojeq3").text.strip()data[i] = {}"Brand": brand,"Product Description": description,"Price": price}(3)return data

main scraping logic

base_url = "https://www.flipkart.com"search_url = " /search?q=sunglasses"sunglasses_data = []

Loop through multiple pages until we have info for 100 sunglasses.

page_count = 1

while len(Sunglasses_data) < 100:

current_url = base_url + search_url + str(page_count)

current_data = scrape_flipkart_Sunglasses(current_url)

page_count

if not current_data:

break

Sunglasses_data.append(current_data)

page_count += 1

Create a DataFrame

df = pd.DataFrame(Sunglasses_data[0:100])

Display to DataFrame

print(df)

This code navigates through the pages by incrementing the 'page_count' and considering the duty unit we have information for 100 sunglasses. As always keep in mind that the website structure may change, so you may need to inspect the HTML of the website and adjust the code accordingly if there are any changes.

Q5. Scrape 100 review data from flipkart.com for iPhone 11 Pro. You have to gather the links:
[https://www.flipkart.com/iphone-11-black-64gb/product-reviews/itm4e5041bd101pd?](https://www.flipkart.com/iphone-11-black-64gb/product-reviews/itm4e5041bd101pd?pid=MOBWF6BXGJCEYVY&lid=LSTMOBWF6BXGJCEYV2XSHRj&marketplace)

Place as shown in the above page. You have to scrape the first market attributes. Those are

- (1) Rating (2) Review Summary (3) R/W Review
- (4) You have to scrape this data for first 100 review.

⇒ To scrape 100 reviews data from flipkart for the iPhone 11. You can use the 'request' library to send HTTP request to the specified URL. BeautifulSoup for parsing the HTML content, and pandas to create a data frame for the scraped data. Here's simple Python program:

python

```
import requests
from bs4 import BeautifulSoup
import pandas as pd

# Function to scrape reviews data from a given
# flipkart product review
def scrape_flipkart_reviews(url, count):
    response = requests.get(url)
    if response.status_code == 200:
        print(f"Scraping review data from page {count}...")
    else:
        print("Failed to fetch the web page. Status code:", response.status_code)
```

`if response.status_code == 200:`

`return None`

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

`review_listings = soup.find_all('div', class_='27m-vg')`

`data = []`

`for review in review_listings`

`rating = review.find('div', class_='3lw21k').text.strip()`

`summary = review.find('p', class_='2-nB2T').text.strip()`

`full_review = review.find('div', class_='1-tky').text.strip()`

`data.append({})`

`'Rating': rating}`

`'Review Summary': summary}`

`'Full Review': full_review}`

`)`

`return data`

`# main scraping logic`

`base_url = "https://www.flipkart.com"`

`product_review_url = "https://www.flipkart.com/iphone-11-black-64-gb/p/itm5011ba101pd&pid=MOBPF46F3XGJCFEN"`

`Y81ID = LSTmoBFung6BxGJCFvny`

`12xSHRJB marker place"`

`reviews_data = []`

`# loop through multiple pages until we have data for 100 reviews`

`page_count = 1`

`while len(reviews_data) < 100:`

`current_url = base_url + product_review_url +`

`&page={}&stx=page`

`current_data = scrape_flipkart_reviews(current_url)`

`page_count`

if not current = dutch

brows

review-duty.extend((current-duty))

page-count += 1

Create a Duty Point

dp = pd.Duty Point reviews_dut [81009]

Display the Duty Point

point (dp)

This code navigates through the pages by incrementing the page-count and concatenates the dutch until we have information for 100 reviews. Please note that the website structure may change so you may need to inspect the HTML of the website and adjust the code accordingly if there are any changes.

Q6 Scrape data for first 100 sneakers you find when you visit flipkart.com and search for sneakers in the search field.

You have to scrape 3 attributes of each sneaker:

Find (1) Product Description (2) Price.

⇒ To scrape data for the first 100 sneakers from flipkart you can the requests library to send HTTP requests to the website BeautifulSoup for parsing the HTML content, and pandas to create a data frame for the scraped data. Here's sample python program:

python

```
import requests
from bs4 import BeautifulSoup
import pandas as pd
```

copy code

function to scrape sneakers data from a given flipkart search page

```
def scrape_flipkart_sneakers(url, count):
    response = requests.get(url)
    if response.status_code == 200:
        print("Scraping sneakers data from page")
    else:
        print("Failed to fetch the webpage, code", response.status_code)
```

Soup = BeautifulSoup(response.text, 'html.parser')

sneakers_listings = soup.find_all('div', class_="col-12 bd")

duty = []

for sneakers in sneakers - listings:

brand = sneaker.find('div', class_ = 'zWkVRV').text.strip()

description = sneaker.find('div', class_ = 'JRPwTC').text.strip()

price = sneaker.find('div', class_ = 'ZoJcp3').text.strip()

duty.append({

'brand': brand

'product description': description,

'price': price}

B).

return duty

main scraping logic

base - url = "https://www.flipkart.com"

search - url = "search?q=sneakers"

sneakers - duty = []

loop through multiple pages until we have data for 100 sneakers:

page - count = 1

while len(sneakers - duty) < 100:

current - url = base - url + search - url + "?page=" + str(page - count)

current - duty = scrape.Flipkart.Sneakers(url + page - count)

if not current - duty:

break

sneakers - duty = extend(current - duty)

page - count += 1

created a duty frame

df = pd.DataFrame(sneakers - duty[100:100])

display the duty frame

print(df)

Q7 Go to webpage <https://www.amazon.in/> Enter "Laptop" in the search field and then click the search icon. Then set CPU type filter to Intel Core i7 as shown in the below image. After setting the filters, scroll down to 10 laptops. Now you have to scrape 3 details for each laptop.

- ① Title
- ② Ratings
- ③ Price

⇒ To scrape data for the first 10 laptops with the specified filters from Amazon.in you can use the requests library to send HTTP request to the website. Beautiful for parsing the HTML content and Pandas to create data frame for the scraped data. Here's sample python program.

Python

 copy code

```
import requests
from bs4 import BeautifulSoup
import pandas as pd
# Function to scrape Laptop data from a given
# Amazon search page with filters.
def Scrape_amazon_laptops(url):
    response = requests.get(url)
    if response.status_code == 200:
        print("Successfully fetched the webpage")
    else:
        print("Failed to fetch the webpage status code: ", response.status_code)
    return response
```

Soup = Beautiful soup of response text (html parsed)
 Laptop = listing s = soup.find_all('div', class_ = 's-l-item')
 (content - margin?)

data = []

for laptop in laptop_listings[81078]:

title = laptop.find('span', class_ = 'text-normal')

text.strip()

ratings = laptop.find('span', class_ = 'icon-alt')

ratings_s = ratings.text if ratings else 'not available'

price = laptop.find('span', class_ = 'a-color-price')

text.strip()

data.append({})

'Title': title,

'Ratings': ratings_s,

'Price': price}

3)

return data

main scraping logic

base_url = "https://www.amazon.in"

search_url = "152K=Laptop"

cup_filter = "filter=3A1375424031%2Cn-feature

-thirteen-browse-bin%3A12598163031"

amazon_url = base_url + search_url + cup_filter

laptops_data = scrape_amazon_laptops(amazon_url)

create a Data frame

df = pd.DataFrame(laptops_data)

display the Data Frame

print(df)

This code performs a search for laptops with the specified filters and extracts data for the first 10 laptops, including the title, rating and price. Keep in mind that the website structure may change, so you may need to inspect the HTML of the website and adjust the code accordingly if there are any changes.

Q8:

write a python program to scrape data for top 1000 quotes of all time

The above task will be done in following steps

- 1 First get the webpage <https://www.a2zquotes.com/>
- 2 Click on Top Quotes
- 3 The scrap (a) quote (b) Author (c) Type of quotes

→

To scrape data for the top 1000 quotes from the A2Z Quotes website, you can use the 'requests' library to send HTTP requests to the website, 'BeautifulSoup' for parsing the HTML content, and 'pandas' to create a data frame for the scraped data. Please note that web scraping should be done in compliance with the website's terms of service, and it's always a good practice to check if the website allows automated access.

Here's a simple python program to get you started:

python

A copy code

```
import requests
```

```
from bs4 import BeautifulSoup
```

```
import pandas as pd
```

```
# function to scrape data for the top 100 quotes
```

```
def scrape_quotes(url=top_quotes_url):
```

```
    response = requests.get(url)
```

```
    if response.status_code == 200:
```

```
        print("Successfully fetched the webpage.")
```

```
    else:
```

```
        print(f"Failed to fetch the webpage. Status code: {response.status_code}")
```

```
    return None
```

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

```
quotes_listings = soup.find_all('div', class_='wrap-block')
```

```
data = []
```

```
for quote in quotes_listings:
```

```
    quote_text = quote.find('a', class_='title').text.strip()
```

```
    author = quote.find('div', class_='author').text.strip()
```

```
    quote_type = quote.find('div', class_='kw-box').text.strip()
```

```
    data.append({
```

```
        'Quote': quote_text,
```

```
        'Author': author,
```

```
        'Type of Quote': quote_type
```

```
})
```

```
return data
```

```
# Main scraping logic
```

```
base_url = "https://www.123quotest.com/"
```

```
top_quotes_url = "https://www.123quotest.com/top-quotes"
```

$$\text{quotes_url} = \text{base_url} + \text{top_quotes}$$
$$\text{quotes_data} = \text{scrape}(\text{quotes_url})$$
$$(129\text{ quotes} - \text{url})$$

Create a Data Frame

~~df = pd.DataFrame(quotes_data)~~

Display the Data Frame

~~print(df)~~

This code performs the steps necessary to scrape data for the top 100 quotes from the [Top Quotes](#) website, including the quote text, author, and type of quote. Remember to check the website's terms of service and adjust the code if there are any changes in the website structure.

Q9

Write a python program to display list of respected former prime minister of India, i.e. (Name, Born-Dead)

(Term of office, Remarks) from <https://www.jagranjosh.com/>

This task will be done in following steps:

- 1 First get the web page <https://www.jagranjosh.com/>
- 2 Then you have to click on the GK option
- 3 Then click on the list of all prime ministers of India
- 4 Then scrap the mentioned data and make the data Python friendly

⇒

unfortunately, I cannot interact with websites or perform web scraping in real-time. However, I can provide you with a python code snippet that you can use as a starting point for web scraping the data you need. please note that web scraping should be done in compliance with the website's terms of service.

python

 Copy code.

import requests

from bs4 import BeautifulSoup

import pandas as pd

Function to scrape data for former prime minister of India

def Scrape_prime_ministers_data(url):

 response = requests.get(url)

 if response.status_code == 200:

 print("Successfully fetched the webpage")

 else:

Print: ("P" paired to help the web page. Selects for response. Status code "200")

return None

Soup: BeautifulSoup(response, "html.parser")

Assuming the structure of the HTML on the page, you may need to inspect the actual HTML

prime_ministers_table = soup.find('table', class_="table-bordered-table")

data = []

for row in prime_ministers_table.find_all('tr')[1:]

skipping the header row

columns = row.find_all('td')

Name = columns[0].text.strip()

Born-Dead = columns[1].text.strip()

Term-of-office = columns[2].text.strip()

Remarks = columns[3].text.strip()

data.append({

'Name': name

'Born-Dead': born-dead,

'Term-of-office': term-of-office,

'Remarks': remarks}

})

return data

main scraping logic

base_url = "https://www.jagranjosh.com/1"

gk_option_url = "gks / list-of-prime-ministers-

-India-144029026-1"

Jagranjosh_url = base_url + gk_option_url.

prime-ministers-data = Scrape-prime-ministers
data (jigzangash - 1x1)

Create a Data Frame

df = pd.DataFrame(prime_ministers_data)

Display the Data Frame

print(df)

This code assumes a certain structure of the HTML on the page. You may need to inspect the actual HTML of the website to ensure that the class names and structure match. Additionally, be aware that web scraping can be sensitive to changes in the website structure. Adjust the code accordingly if there are any changes to the website.

Q 10

write a python program to display of 50 most expensive cars in the world (i.e car name and price) from <https://www.motor1.com>

This task will be done in following steps:

- 1 First get the web page <https://www.motor1.com>
- 2 Then you have to type in the search bar so most expensive cars
- 3 Then click on 50 expensive cars in the world
- 4 Then Scrape the mentioned data and make the data frame

⇒

To achieve this task, you can use the 'request' library to send HTTP requests to the website, BeautifulSoup for parsing the HTML content, and pandas to create a data frame for the scraped data. Here's a simple Python program.

Python

```
import requests
from bs4 import BeautifulSoup
import pandas as pd
# Function to scrape data for the 50 most expensive cars from motor1
def scrape_expensive_cars(url):
    response = requests.get(url)
    if response.status_code == 200:
        print("Successfully fetched the webpage")
    else:
        print(f"Failed to fetch the webpage status code {response.status_code}")
```



copy code.

return func

Soup = BeautifulSoup(response.text, 'html.parser')

Assuming the structure of the HTML on the page (you may need to inspect the actual HTML)

cars_listings = soup.find_all('div', class_ = 'article-widget--list-item')

data = []

for car in cars_listings[0:50]:

car_name = car.find('a', class_ = 'article-widget--list-item-title').text.strip()

car_price = car.find('div', class_ = 'article-widget--list-item-details').text.strip()

data.append({

'car name': car_name,
'price': car_price

})

main scraping logic

base_url = "https://www.motor1.com"

search_query = "50 most expensive cars"

search_url = ("F" "search / ?q={search_query.replace(' ', '_')}")

motor1_url = base_url + search_url

expensive_cars_data = scrape_expensive_cars(data, motor1_url)

create a Data Frame

df = pd.DataFrame(expensive_cars_data)

display the Data Frame

print(df)