

My Name Vaishali Rambus Saini

Batch No : DS2311

Sir :- Shwetank Mishry

Date :- 9-2-2023

Assignment 3 Web scraping

Q 1

wrote a python program which searches the product under a particular product from www.amazon.in. The product to be searched will be taken as input from user. For e.g. if user input is guitar, Then search for guitars.

⇒

Certainly ! Here's Python program using the Beautiful Soup library to scrape Amazon.in and search for products based on user input.

Python

A copy code

```
import requests
from bs4 import BeautifulSoup
def search_amazon(Product):
    URL = f"https://www.amazon.in/s?k={product}"
    headers = {
        "User-Agent": "Mozilla/5.0 Windows NT 10.0"
    }
    response = requests.get(URL, headers=headers)
    response.raise_for_status()
```

try

```
response = requests.get(URL, headers=headers)
response.raise_for_status()
```

Soup = BeautifulSoup \in response.text, ('html.parser')
products = soup.find_all('div', {'class': 'result-item'})

if len(products) > 0:

~~Print ("F" showing results for Product 3: 1 b")~~

for item in products > 0

Print F "showing results for E product 3':In")

~~for item in products:~~

product-name-item. Find ('span', 'class': 'c')
ct-text-normal & .text-start

~~product - price = item. Find ("span" ∈ class "a-price")~~

~~unit price = product - name 3 / unit price~~

Σ product - price³ (n): diagonal elements
also = sum of all off-diagonal elements

~~else~~ Print ("Two products found.")

except requests.exceptions.HTTPError as e:

~~print(F"HTTP error occurred: {e}")~~

Except Exception clause:

print ("An error occurred see 'e'"')

Search-query = Input (Enter the product you want to search for on Amazon IPo")

Search - amazon Search - query).

This program will prompt the user to input a product, then it will scrape Amazon.in for that product, and point out the names and prices of the products found. Make sure to install the requests and BeautifulSoup libraries if you haven't already.

A) Copy Code

pip install requests

pip install BeautifulSoup4

This code can be further extended to include more detailed information or additional functionality as needed.

Q2

In the above question, you scrape the following details of each product listed in first 3 page of your search.

results and save it in a data frame and csv. In case if any product has less than 3 pages in search results then Scrape all the products available under that product name. Details to be scraped are:

"Brand name", "Name of the product", "Price", "Return Exchange", "Expected Delivery", "Availability" and "Product URL", in case, if any of the details are missing for any of the products then replace it by "N/A".

⇒

To accomplish this task, we'll need to extend the previous program to scrape the required details for each product listed on the first three pages of the search results on Amazon.in. We'll use the pandas library to organise the data frame and save it to a CSV file. Here's the updated code,

Python

 Copy code

import requests

from BS4 import BeautifulSoup

import random as rd

def scrape_product_details(url):

headers = {

"User-Agent": "Mozilla/5.0 Windows NT

10.0; Win64; x64) AppleWebKit/537.36 (KHTML,

like Gecko) Chrome / 58.0.3029.110 Safiri /

537.3"}

products_details = []

try

response = requests.get(url, headers=headers)

response.raise_for_status()

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

product_name = soup.find('div', {'id': 'productTitle'})

.text.strip()

products_price = soup.find('span', {'id': 'productTitle'})

.text.strip()

price = soup.find('span', {'id': 'priceblock_ourprice'})

.text.strip()

return_exchange = soup.find('div', {'id': 'productSupport'})

.find('div', {'id': 'returnPolicy'})

.text.strip()

else:

expected_delivery = soup.find('div', {'id': 'dmPromise'})

.text.strip()

availability = soup.find('div', {'id': 'dmDeliveryPromise'})

.text.strip()

availability = soup.find('div', {'id': 'availability'})

.text.strip()

availability = soup.find('div', {'id': 'availability'})

.text.strip()

Product - details - append (Exception handling)

'brand' name'	:	brand - name
'name of the product'	:	product - name
'price'	:	price
'Return / Exchange'	:	return - exchange
'Expected delivery'	:	expected - delivery
'Availability'	:	availability
'product URL'	:	URL

3)

```
except requests.exceptions.HTTPError as e:
    print("An error occurred : ", e)
except Exception as e:
    print("An error occurred : ", e)
```

```
return product - details
def search_amazon(product):
    base - url = "https://www.amazon.in/s?k="
    headers = {"User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/58.0.3029.110 Safari/537.36"}
    product - details - cell = []
    for i:
```

```
        for page - num in range(1, 4):
            url = f'{base - url}{product}{page}{page - page - num}'
            response = requests.get(url, headers=headers)
            response.raise_for_status()
            soup = BeautifulSoup(response.text, 'html.parser')
            products = soup.find_all('div', {'class': 's-result-item'})
            item3 = products[3].find('div', {'class': 's-item-container'})
```

```
for item in products:
```

```
    product_link = item.find('a' in class_ == 'a-link').get('href')
```

```
    if product_link:
```

```
        product_url = "https://www.amazon.in" + product_link[11:-3])
```

```
        product_details = scrape_product_details(product_url)
```

```
        product_details = cell.extend(product_details)
```

```
        except requests.exceptions.HTTPError as e:
```

```
            print(f"HTTP error occurred: {e}")
```

```
        return product_details[cell]
```

```
if __name__ == '__main__':
```

```
    search_query = input("Enter the product you want  
to search for on Amazon.in: ")
```

```
    product_details = search_amazon(search_query)
```

```
    if product_details:
```

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

```
        df.to_csv(f'{search_query}.csv', index=False)
```

```
        print("Product details saved to CSV file.")
```

```
    else:
```

```
        print("No products found.")
```

This code will search for the given product on amazon in scrape the details for each product listed on the first three pages of the search result (overall available page is less than three); and save the details to a CSV file named "`search_query-details.csv`". If any detail is missing for product, it will be replaced by "2".

Q3

write a python program to access the search box and search button on <image.google.com> and Scrape 10 images each for keyword 'Purists', 'car' and machine learning, humor, 'california coffee'.



To accomplish this task we can use the selenium library in python which allows us to automate interactions with a web browser. First make sure you have the selenium library installed pip install selenium) and you also need to download a compatible web driver for the browser you want to automate (e.g chrome driver for google chrome). Here's the python program.

Python

Copy code

```
from selenium import webdriver
from selenium.webdriver.common.keys import Keys
import time
import os
import requests

# Function to save image
def save_image(image_url, directory):
    try:
        img_data = requests.get(image_url).content
        with open(os.path.join(directory, os.path.basename(image_url)), 'wb') as handle:
            handle.write(img_data)
        print(f"Image saved: {os.path.basename(image_url)}")
    except Exception as e:
        print(f"Error saving image: {e}")
```

~~H Function to Scrape Image~~

```

def scrap_image(keyword, num, images):
    driver = webdriver.Chrome(executable_path="path
        -to-chromedriver")
    driver.get("https://www.google.com/search")
    search_box = driver.find_element(Byxpath('
        input[@title='Search'])')
    search_box.send_keys(keyword + keys.RETURN)
    # wait for results to load
    time.sleep(2)
    # scroll down to load more images
    last_height = driver.execute_script("return docum
        d.height; scrollHeight return document.documentElement.scrollHeight")
    while True:
        driver.find_elements_by_xpath('
            [reclass="rg - i o (l u wd")'])
        driver.execute_script("window.scrollTo(0,
            document.body.scrollHeight")
        time.sleep(2)
        new_height = driver.execute_script("return docum
            d.height; scrollHeight")
    if new_height == last_height:
        break
    last_height = new_height

```

H. Get image URLs

```

image_elements = driver.find_elements(Byxpath(
    '
        image[@class="rg - i o (l u wd")'])
image_urls = elem.get_attribute("src") for
in image_elements if elem.get_attribute("src")]
# create directory to save image
os.makedirs(keyword, exist_ok=True)

```

save image

for i, url in enumerate(image_urls, [run
= images]):

Sure-image(url, keyword)

if i == num_images - 1:

break.

dollar.quit()

if --frame-- == '--main--':

Keywords = ['fruits', 'cars', 'machine learning',
'Guitar', 'Cakes',]

num_images = 10

for keyword in Keywords:

Scrape-image(keyword, num_images)

This script will open Google images in a chrome browser, search for the specified keywords and Scrape the URLs of the images. Then, it will save the first 10 images for each keyword in separate directories, named after the keywords. Make sure to replace "path-to-chromedriver" with the actual path and your chrome driver executable.

Q4 Write a Python program to search for a smart phone (e.g. oneplus, nokia, pixel 4A, etc). on www.flipkart.com and scrape following details from the search results displayed on 1st page. Details to be scraped: Brand, name, "Smart phone Name", "Colour", "RAM", "Storage" (Ram), primary camera, price

→ Python

A copy code.

import requests

from bs4 import BeautifulSoup

def scrape_smart_phones(product):

url = "https://www.flipkart.com/search?q=(product)&otracker=search%20smartphone%20-%20FLIPKART&as-show-on%20as-off"

headers = {

"User-Agent": "Mozilla/5.0 windows NT 10.0"

"win64: x64) AppleWebKit/537.36 (KHTML, like Gecko)"

Chrome/58.0.3029.110 Safari/537.3"

3

try:

response = requests.get(url, headers=headers)

response.raise_for_status()

soup = BeautifulSoup(response.text, "html.parser")

products = soup.find_all("div", {"class": "a1wbe3"})

for product in products:

brand_name = product.find("div", {"class": "row13"}).text.strip()

Smart Phone Name = product.find("div", {"class": "row13"}).text.strip()

specs = product.find_all("li", {"class": "row13"})

color = specs[0].text.strip()

```

xum = spec[2].text.strip()
primary = cameras[3].text.strip()
print(f"Brand name : {brand_name}")
print(f"Smart phone name : {smart_phone_name}")
print(f"Color : {color}")
print(f"RAM : {RAM}")
print(f"Storage (ROM) : {ROM}")
print(f"Primary camera : {primary_camera}")
print(f"RAM : {RAM}")
except requests.exceptions.HTTPError as e:
    print(f"An error occurred: {e}")
except Exception as e:
    print(f"An error occurred: {e}")

```

ip = --name= --main --

search_query = input("Enter the smart phone you want to search for on Pilipketat")

scope = smart_phone(search_query)

This program takes user input for the smart phone to search on Pilipketat, scrapes the details for each product listed on the first page of the search results, and prints out the details for each smart phone found. You can run this program with the smart phone name (e.g. Oneplus Nord), and it will display the brand name, smart phone name, color, RAM, ROM, and primary camera for each product listed on the first page of the search results on Pilipketat.

Q5 Write a program to scrap geospatial coordinates (latitude, longitude) of a city searched on google maps.

→ To scrape geospatial coordinates (latitude and longitude) of a city searched on google maps, we can use Google Maps Geocoding API. This API allows us to retrieve geolocation data based on search queries.

Here's Python program to accomplish this using the 'requests' library to make API requests:

python

copy code

```
import requests
def get_coordinates(city_name):
    api_key = "YOUR-Google-MAPS-API-KEY"
    url = f"https://maps.googleapis.com/api/geocode/json?address={city_name}&key={api_key}"
    try:
        response = requests.get(url)
        response.raise_for_status()
        data = response.json()
        if data['status'] == 'OK':
            location = data['results'][0]['geometry']['location']
            latitude = location['lat']
            longitude = location['lng']
            print(f"Geospatial coordinates for {city_name}:")
            print(f"Latitude: {latitude}")
            print(f"Longitude: {longitude}")
        else:
            print(f"Error: unable to get coordinates.")
```

```
except requests.exceptions.HTTPError as e:  
    print(F"HTTP error occurred: {e}")  
except Exception as e:  
    print(F"An error occurred: {e}")
```

if __name__ == "__main__":

city = input("Enter the city name to get its coordinates: ")
get_coordinates(city)

Before running the program, make sure to replace your Google Maps API Key with your actual Google Maps API key, you can obtain an API key by following the instruction here: <https://developers.google.com/maps/api-key-starter>.

Once once you have your API key, you can run the program, input the city name you want to search for on Google maps, and it will print out the geographical coordinates (latitude and longitude) of that city.

Q6

write a program to scrap all the available details of Best gaming laptops from digit.in.

⇒ To scrap all the available details of the gaming laptop from digit.in we can inspect the website to understand its structure and then use python along with the Beautiful soup library to extract the necessary information. Here's python program to accomplish this task.

python

Copy code

import requests

from Bs4 import BeautifulSoup

def soups -> gaming_laptops ():

url = "https://www.digit.in/top-products/Best-gaming-laptops-4.0.html#1"

headers = {

"User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like gecko) Chrome/58.0.3029.110 Safari/537.36"

}

try

response = requests.get(url, headers=headers)

response.raise_for_status()

soup = BeautifulSoup(response.text, "html.parser")

laptops = soup.find_all("div", {"class": "laptop"})

for laptop in laptops:

laptop_name = laptop.find("h3").text.strip()

specs = laptop.find("ul")

for spec in specs.children:

Processor = Spec[0].text.strip()

ram = Specs[1].text.strip()

os = Specs[2].text.strip()

display = Specs[3].text.strip()

storage = Specs[4].text.strip()

price = Specs[5].text.strip()

print("Laptops names:", laptop_names)

print("Processor, Processor")

print("Ram:", ram)

print("Operating system:", os)

print("Display:", display)

print("Storage:", storage)

print("Price:", price)

print("By")

except requests.exceptions.ConnectionError as e:

print("A connection error occurred (%s)" % e)

except Exception as e:

print("An error occurred (%s)" % e)

if __name__ == "__main__":

scrape_gaming_laptops()

This program will scrape the details of the best gaming laptop from digit.in and print out the laptop name, processor, RAM, Operating system, display, storage & price for each laptop. You can further process this info as per your requirement, such as saving it a file or storing it in a database.

Q7 Write a Python program to scrape the details for all billionaires from www.forbes.com. Details to be scrapped:

"Rank", "Name", "Net worth", "Age", "Citizenship", "Source", "Industry", ..

→ Sure, here's a Python program that the Beautiful Soup library to scrape the details of billionaires from Forbes website.

Python



Copy code

```
import requests
from bs4 import BeautifulSoup
def scrape_billionaires():
    url = "https://www.forbes.com/billionaires/"
    response = requests.get(url)
    soup = BeautifulSoup(response.content, "html.parser")
    billionaires_data = []
    for row in soup.find_all("div", class_="table-row"):
        rank = row.find("div", class_="rank").text.strip()
        name = row.find("div", class_="person-name").text.strip()
        net_worth = row.find("div", class_="net-worth").text.strip()
        age = row.find("div", class_="age").text.strip()
        citizenship = row.find("div", class_="country-of-citizenship")
        source = row.find("div", class_="source").text.strip()
        industry = row.find("div", class_="category").text.strip()
        billionaire = {
            "Rank": rank,
            "Name": name,
            "Net worth": net_worth,
            "Age": age,
            "Citizenship": citizenship,
            "Source": source,
            "Industry": industry
        }
        billionaires_data.append(billionaire)
    return billionaires_data
```

"Citizenship" :- citizenship,

"Source" :- source,

"Industry" :- industry

3.

billionaires - dual-capped (Billionaire)

return billionaires - dutch

If -- name -- == "-- main --" &

Billionaires = scope - Billionaires ()

for billionaires in billionaires:

print (billionaire)

This code will extract the required details (Rank, Name, net worth, Age, Citizenship, Source, and Industry) for each billionaire listed on Forbes website and print them out. Make sure to have the Beautiful Soup library installed ('pip install BeautifulSoup4') before running this code.

Q8 Write a program to extract at least 500 comment upvote and time when comment was posted from any YouTube video.

⇒ To extract comments, comment upvotes, and the time when each comment was posted from a YouTube video, you can use the YouTube Data API. First, you need to set up a project in the Google Cloud Console and enable the YouTube Data API for that project. Then, you can use the google-api-python-client library to interact with the API. Below is a Python program that demonstrates how to do this.

Python

Copy code

```
from googleapiclient.discovery import Build
from datetime import datetime
# Set up the YouTube Data API service
api_key = "YOUR-API-KEY" # Replace "YOUR-API-KEY" with your actual API key
youtube = build("youtube", "V3", developerKey=api_key)

def get_video_comments(video_id=max_results=500):
    comments = []
    next_page_token = None
    while len(comments) < max_results:
        response = youtube.commentThreads().list(
            part="snippet",
            videoId=video_id,
            maxResults=max_results,
            pageToken=next_page_token)
        for item in response["items"]:
            comment = item["snippet"]["topLevelComment"]
            author = comment["snippet"]["authorDisplayName"]
            text = comment["snippet"]["textDisplay"]
            timestamp = comment["snippet"]["publishedAt"]
            upvotes = comment["snippet"]["likeCount"]
            comments.append((author, text, timestamp, upvotes))
        if "nextPageToken" in response:
            next_page_token = response["nextPageToken"]

    return comments
```

~~maxResults = min(maxResults, len(Comments), 100)~~

PageToken = next - pageToken
).execute()

for item in response["items"]:

Comment = item["snippet"] # Placeholder Comment
"SNIPPET"

CommentText = Comment["textDisplay"]

CommentUpvotes = Comment.get("likeCount", 0)

CommentTime = Comment["publishedAt"]

CommentTime = datetime.datetime.strptime(CommentTime, "%Y-%m-%dT%H:%M:%S")

Comments.append({

"text": CommentText,

"upvotes": CommentUpvotes,

"time": CommentTime

})

if "nextPageToken" in responses:

nextPageToken = response["nextPageToken"]

else:

break

return Comments

If --name-- == "main":

videoId = "YOUR-VIDEO-ID".

Replace "YOUR-VIDEO-ID" with the ID
of the YouTube video

(comments = get - Video - comments (video - id))

for comment in comments:

```
print ("Comment", comment["text"])
print ("Upvotes", comment ["upvotes"])
print ("times", comment ["time"])
print (x).
```

make sure to replace "YOUR-API-KEY" with your actual YouTube Data API key and "YOUR-VIDEO-ID", with the ID of the YouTube video you want to extract comments from. This code will print out the comments, the number of upvotes each comment has received, and the time when each comment was posted.

Q9

Write a python program to Scrape data for all available Hostels from <https://www.hostelworld.com/in/London>. You have to Scrape hostel name, distance from city centre, rating, total reviews, overall reviews, facilities from price, dorms, from price, facilities and property description.

⇒ To Scrape data for all available hostels in London from Hostelworld you can use the Beautiful Soup library in python. Here's python program to achieve this.

Python

A copy code

```
import requests
```

```
from bs4 import BeautifulSoup
```

```
def Scrape_hostels_in_london():
```

```
url = "https://www.hostelworld.com/find/hostel?city=London&chosenCountry=England"
```

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

```
soup = BeautifulSoup(response.content, "html.parser")
```

```
hostels = []
```

```
hostel_list = soup.find_all("div", class_="property")
```

```
for hostel in hostel_list:
```

```
name = hostel.find("h2", class_="title").text.strip()
```

```
distance = hostel.find("span", class_="description")
```

```
scoring = hostel.find("div", class_="score orange")
```

```
total_reviews = hostel.find("div", class_="key")
```

```
overall_reviews = hostel.find("div", class_="key")
```

```
privileges_price = hostel.find("div", class_="price")
```

```
total_price = hostel.find("div", class_="price")
```

Text .strip()

`facilities = [item.text.strip() for item in hostel.find_all("li", class_="Facility-badge")]`

`description = hostel.find("div", class_="keycard").find_next_sibling("div").text.strip()`

`hostel_data = {`

`"name": Name,`

`"distance from city centre": distance,`

`"Rating": Rating,`

`"total_reviews": total_reviews,`

`"overall_reviews": overall_reviews,`

`"privates_from_price": -privates - price,`

`"dorms_from_price": dorms - price,`

`"Facilities": facilities,`

`"Description": description,`

3

`hostels.append(hostel_data)`

`return hostels`

`if __name__ == "__main__":`

`london_hostels = scrape_hostels_in_london()`

`for hostel in london_hostels:`

`print(hostel)`

`print()`

This code will scrape data for all available hostels in London from Hostelworld, including hostel name, distance from the city centre, rating, total reviews, overall reviews, prices for private rooms and dorms, facilities, and property description. Make sure to have the BeautifulSoup library installed ("pip install BeautifulSoup4") before running this code.