

My Name - Vaishali Salve

Batch No - IS 231

Date: 30/02/2024

Str 8 - Shwetank Mishra

Assignment Name :- Web-Scraping

Assignment 4

Q1 Scrape the details of most viewed videos on YouTube from wikipedia. URL - https://en.wikipedia.org/wikilist_of_most-viewed_Youtube_videos You need to find following details: (A) Rank, (B) name (C) Artist (D) Upload date (E) views

→ To scrape the details of the most viewed videos on the YouTube from wikipedia you can use python with libraries like BeautifulSoup and requests. It's a basic example of how you can achieve this:

python

copy code

```
import requests
```

```
from bs4 import BeautifulSoup
```

```
# URL of the wikipedia page
```

```
url = 'https://en.wikipedia.org/wikilist_of-most-viewed_Youtube_videos'
```

```
It sends a GET request to the URL
```

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

Parse the HTML - Content

Soup = BeautifulSoup(response.text('htm.parser'))

Find the table containing the most viewed videos.

table = soup.find('table', {'class': 'wikitable'})

Initialize lists to store data

rank = list() = []

name_list = []

artist_list = []

upload_date_list = []

views_list = []

Iterate over each row in the table (skipping the header row)

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

Extract data from each column in the row.

columns = row.find_all('td')

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

name = columns[1].text.strip()

artist = columns[2].text.strip()

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

views = columns[4].text.strip()

Append data to respective lists.

rank_list.append(rank)

name_list.append(name)

artist_list.append(artist)

upload_date_list.append(upload_date)

views_list.append(views)

```
A print the scraped data  
for i in range (len (Rank - list)) :  
    print ("F" "Rank : " rank - list [i] [3])  
    print ("F" "Name : " name - list [i] [3])  
    print ("F" "Artist : " artist - list [i] [3])  
    print ("F" "Upload Date : " upload - date - list [i] [3])  
    print ("F" "Views : " views - list [i] [3])  
print ()
```

The this Script sends a GET request to the wikipedia page, parses the HTML content using BeautifulSoup, find the table containing the most viewed videos, extracts the required details (Rank, Name, Artist, Upload Date, views) from each row and then prints the scraped data. You might need to install the BeautifulSoup library if you haven't already ('pip install BeautifulSoup4').

Q2. Scrape the details from Team India's international fixtures from bcci.tv. URL = <http://www.bcci.tv/>. You need to find following details:
 (A) Series (B) Place (C) Date (D) Time
 Note: From bcci.tv home page you have reach to the international fixture page through code.

→ To scrape the details of Team India's (international) fixtures from the BCCI website, you can use Python with libraries like BeautifulSoup and requests.
 Here's how you can do it.

Python

A Copy Code

```
import requests
```

```
from bs4 import BeautifulSoup
```

```
# URL of the BCCI website
```

```
url = "https://www.bcci.tv/"
```

```
# send a GET request to the URL
```

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

```
# parse the HTML content
```

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

```
# find the link to the international fixtures page
fixtures_link = soup.find('a', text = "International
Fixtures")['href']
```

```
# Construct the URL for the international fixtures page
```

```
fixtures_url = url + fixtures_link
```

```
# send a GET request to the international fixtures page
```

Fixtures-response = request.get('fixtures-ur11')

Parse the HTML - Content of the fixtures page
fixtures-soup = BeautifulSoup(fixtures-response.text, 'html.parser')
Find the section containing the fixtures
fixtures-section = fixtures-soup.find('div', class_='js-list')

Initialize list to store data

series-list = []
place-list = []
date-list = []
time-list = []

Iterate over each fixture
for fixture in fixtures-section.find_all('li', 'list-item'):

Extract date from the fixture
series = fixture.find('p', class_='fixtures-additional-info').text.strip()
place = fixture.find('span', class='l-disknowet-with-text').text.strip()
date = fixture.find('span', class='fixture-date').text.strip()
time = fixture.find('span', class='fixture-time').text.strip()

Append data to respective list

```
series-list.append(series)
```

```
place-list.append(place)
```

```
date-list.append(date)
```

```
time-list.append(time)
```

```
# print the scraped data
```

```
for i in range(len(series-list)):
```

```
    print("Series : " + series-list[i][3])
```

```
    print("Place : " + place-list[i][3])
```

```
    print("Date : " + date-list[i][3])
```

```
    print("Time : " + time-list[i][3])
```

```
print()
```

This script first navigates to the BCCI website, finds the link to the international fixtures page, constructs the URL for that page, sends a GET request to it and then parses the HTML content to extract the required details (Series, Place, Date, Time) for each fixture. Finally it prints the scraped data.

Q3

Scrape the details of supervisor G.D. of India from statisticsline.com.

URL = <http://www.statisticsline.com/>

- You have to find following details
- (A) Name
 - (B) State
 - (C) GSDP (18-19) - at current prices
 - (D) GSDP (19-20) - at current prices
 - (E) Share (18-19) - (F) GDP (5 billion)

Note & from statisticsline home page you can reach to company page through code.

⇒

Python

Copy code

```
import requests
```

```
from bs4 import BeautifulSoup
```

```
# URL of the Statisticsline website
```

```
url = http://www.statisticsline.com//?
```

```
# Send a GET request to the URL
```

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

```
# Parse the HTML - Content
```

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

```
# Find the link to the economy page
```

```
economy_link = soup.find('div', class_="dropdown-content").find('a', text="Economy")
```

```
# Construct the URL for the economy page
```

```
economy_url = url + economy_link
```

Send a GET request to the economy page.
 economy-response = requests.get('economy-url')

Parse the HTML content of the economy page.
 economy-soup = BeautifulSoup(economy-response.text, 'html.parser')

Find the link to the Indian States GDP page.
 Indian-states-gdp-link = economy-soup.find('a', class_="dropdown-menu").find('a', text="Indian States")["href"]

Construct the URL for the Indian States GDP page.
 Indian-states-gdp-response = requests.get(Indian-states-gdp-link)

Send a GET request to the Indian states GDP page.

Indian-states-gdp-response = requests.get(Indian-states-gdp-link)

Parse the HTML content of the Indian States GDP page.

Indian-states-gdp-soup = BeautifulSoup(Indian-states-gdp-response.text, 'html.parser')

Find the table containing the state-wise GDP details.

table = Indian-states-tables-gdp-soup.find("table", {"id": "table-id"})

Initialize lists to store data.

link-list = []

state-list = []

gdp-18-19 = current-packages-list = []

gdp-19-20 = current-packages-list = []

share-18-19 = list = []

gdp-billion-list = []

Iterate over each row in the table (skipping the headers row)

for row in table.find_all("tr"):

Extract data from each column in the row

columns = row.find_all("td")

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

state = columns[1].text.strip()

gsdp = 18-19 - current_prices = columns[2].text.strip()

gsdp = 19-20 - current_prices = columns[3].text.strip()

share = 18-19 - columns[4].text.strip()

gdp = billion = columns[5].text.strip()

Append data to respective list

rank_list.append(rank)

state_list.append(state)

gsdp_18-19 - current_prices_list.append(gsdp_18-19)

gsdp_19-20 - current_prices_list.append(gsdp_19-20)

share = 18-19 - list.append(share_18-19)

gdp = billion - list.append(gdp_billion)

Print the scraped data

for i in range(len(rank_list)):

print("P" rank[i] & rank_list[i] & "3")

print("P" state[i] & state_list[i] & "3")

print("P" gsdp_18-19[i] & current_prices[i] & gsdp_18-19[i])

print("P" gsdp_19-20[i] & current_prices[i] & gsdp_19-20[i])

print("P" share[i] & share_18-19[i] & "3")

print("P" gdp[i] & gdp_billion[i] & "3")

print("P" gdp[i] & gdp_billion[i] & "3")

This script first navigates to the statisticstim website find the link to the economy page (constur) the URL that page, sends a GET request to it find the link the Home state, GDP page, constructs the URL for that sends a GET request to it and then parse the HTML content to extract required details (Rank, State, GDP (18-19) - current policies, GDP (19-20) current policies, shares (18-19) - GDP (\$billion) for own state. Finally, it print the scraped data.

Q4 Scrape the details of trending repositories of GitHub.com. URL - <https://github.com/>

⇒ To Scrape the details of trending repositories you can use python with libraries like BeautifulSoup and requests. Here's how you can do it in Python.

 copy code

```

import requests
from bs4 import BeautifulSoup
# Function to get the trending repositories data
def get_trending_repositories():
    # URL of GitHub trending page
    url = "https://github.com/Trending"
    # send a GET request to the URL
    response = requests.get(url)
    # parse the HTML content
    soup = BeautifulSoup(response.text, "html.parser")
    # find the list of trending repositories
    repo_list = soup.find_all("article", class_="Box-row")
    # Initialise lists to store data
    titles = []
    descriptions = []
    contributors_counts = []
    languages = []
    # extract data for each repository
    for repo in repo_list:
        title = repo.h1.text
        description = repo.p.text
        contributors_count = repo.div[1].text
        language = repo.div[2].text
        titles.append(title)
        descriptions.append(description)
        contributors_counts.append(contributors_count)
        languages.append(language)
    return {
        "titles": titles,
        "descriptions": descriptions,
        "contributors_counts": contributors_counts,
        "languages": languages
    }

```

```

# Extract repository title
title = repo.find("h1", class_="h3 lh-condensed").text.strip()

titles.append(title)

# Extract language (pse)
language = repo.find("span", class="f inline-block md-o mr-3").text.strip()

languages.append(language)

# Return the scraped data
return title, descriptions, contribution_counts, languages

```

Get the trending repositories details

title, descriptions, contribution_counts, languages = get_trending_repositories()

```

# Print the scraped data
for i in range(len(titles)):
    print(f"Repository Title: {title[i]}")
    print(f"Repository Description: {descriptions[i]}")
    print(f"Contribution Count: {contribution_counts[i]}")
    print(f"Language Used: {languages[i]}")
    print()

```

This script navigates to the GitHub trending page, parses the HTML content to extract details (Repository Title, Repository Description, Contribution Count, Language Used) for each trending repository, and then prints the scraped data.

Q5

Scrape the details of top 100 songs on billboard.com

URI = <http://www.billboard.com>

You have to find the following details:

- (A) Song name
- (B) Artist name
- (C) Last week rank
- (D) Peak rank
- (E) Weeks on Board

Notes: From the home page you have to click on the chart option then hot 100 page link through code



Python

A copy code

```
import requests
```

```
from bs4 import BeautifulSoup
```

```
# Function to get the top 100 songs data
```

```
def get_top_100_songs():
```

```
# URI of Billboard's Hot 100 page
```

```
URI = https://www.billboard.com/charts/hot-100
```

```
# Send a GET request to the URL
```

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

```
# Parse the HTML content
```

```
Soup = BeautifulSoup(response.text, "html.parser")
```

```
# Find the list of top 100 songs
```

```
songs_list = Soup.find_all("li", class_="chart-list-item")
```

```
# Initialize lists to store data
```

- song_numbers = []

- artists = []

- last_week_ranks = []

peak - ranks = []

books - on board = []

Extract Song name

Song - name = Song . find ("span", class, "chart-element-information--song") . text . strip ()
 song - name . append (Song - name)

Extract artist name

artist - name = Song . find ("span", class, "chart-element-information--artist") . text . strip ()
 artist - name . append (Artist - name)

Extract last week rank

last week - rank = Song . find ("span", class, "chart-element-menu-text-center-color-secondary-text-lust") . text . strip ()

#

Extract Peak rank

peak - rank = Song . find ("span", class, "chart-element-menu-text-center-color-secondary-text-peak") . text . strip ()
 peak - rank . append (peak - rank)

Extract week on Board

weeks - on - board = Song . find ("span", class, "chart-element-menu-text-center-colors-secondary-text-week") . text . strip ()

week - on - board . append (weeks - on board)

Return the scraped data

return song - name, artist - name, last week, ranks, peak - ranks, weeks - on - board

Get the top 100 songs duty!

Song - names → artist - names, last - week
rank, peak - ranks, weeks on board
= get - top - 100 songs ()

Print the scraped dut.

for i in song: (len(song) >= 1)

print (F"Song name: {Song - names [i]}")

print (F"Artist name: {Artist - names [i]}")

print (F"Last week Rank: {last - week Rank [i]}")

print (F"Peak Rank: {peak - rank [i]}")

print (F"Weeks on Board: {weeks - on board [i]}")

print ()

This script navigates to the Billboard Hot 100 page parses the HTML content to extract details (Song name, Artist Name, Last week Rank, Peak Rank, Weeks on Board) for each song in the top 100 list and then prints the scraped data.

Scope the details of Highest selling novels.

- (A) Book name (B) Author name (C) Volume no
(D) Publisher (E) Genre

URL - <https://www.theguardian.com/news/1decadebook/01/2012/aug/09/best-selling-book-all-time-fifty-shades-grey> - Compare

~~Python~~

Copy Code

```
import requests  
from bs4 import BeautifulSoup  
# Function to get the highest selling novels url  
def get_highest_selling_novels():  
    url = "https://www.theguardian.com/news/datablog/2017/jul/09/best-selling-book-all-time-fifty-shades-grey-compare"
```

Send a GET request to the URL

~~response = request.get('URL')~~

~~# parse the HTML content~~

~~Soup = BeautifulSoup(response.text, 'lxml')~~

Find the table (containing) the highest selling novels daily

~~table = Soup. find ("table", class_= "id scratchs-Scratch")~~

Initiative 1 is to store data

book - hymns = []

Author - Name = []

~~Volumes sold = 17~~

Publishes - [7]

genres = [

Extract data for each novel

for row in table: find cell ("By") [1:3] :

(columns = row, find cell ("fd"))

book = Name = columns [1].text.strip()

author_name = columns [2].text.strip()

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

publisher = columns [4].text.strip()

genre = columns [5].text.strip()

Append data to respective lists

book_names.append(book_name)

author_names.append(author_name)

volumes.append(volumes)

publishers.append(publisher)

genres.append(genre)

Return the scraped data

return book_names, author_names, volumes, sold,
publishers, genres

Get the highest selling novels

book_names = author_names, volumes, sold,

publishers, genres = get_highest_selling
movies()

Print the scraped data

for i in range (len(Books), name)) :
 3

```

print(F"Book-name : {Book-names[1][3]}")
print(F"Author-name : {author-names[1][3]}")
print(F"Volums-Sold : {volume-sold[1][3]}")
print(F"Publisher : {publisher[1][3]}")
print(F"Genre : {genres[1][3]})")
print()
    
```

print()

This script navigates to the guardian's webpage containing information about the highest selling novels, parses the HTML content extract details (Book name, Author name, Volumes Sold, Publisher, Genres) for each novel and then prints the scraped data.

Q7) Scrape the details of most watched TV series of all time [imdb.com/list/ls095964455/](https://www.imdb.com/list/ls095964455/) you have to find the following details

- (A) name (B) year span (C) genre (D) runtime
- (E) Rating (F) votes



Python

Copy code

```

import requests
from bs4 import BeautifulSoup

# Function to get the most watched TV series
def get_most_watched_tv_series():
    # URL of the IMDB page containing the most watched TV series
    url = "https://www.imdb.com/list/ls095964455/"

    # Send a GET request to the URL
    response = requests.get(url)

    # Parse the HTML content
    soup = BeautifulSoup(response.text, "html.parser")

    # Find the list of most watched TV series
    series_list = soup.find("div", class_="list-item-content")

    # Initialize list to store data
    names = []
    years_spans = []
    genres = []
    runtimes = []
    ratings = []
    votes = []

    # Loop through the list items
    for item in series_list:
        name = item.h3.a["title"]
        year_span = item.p.span["title"]
        genre = item.p.find("span", class_="genre").text
        runtime = item.p.find("span", class_="runtime").text
        rating = item.p.find("span", class_="rating").text
        votes_count = item.p.find("span", class_="votes").text

        names.append(name)
        years_spans.append(year_span)
        genres.append(genre)
        runtimes.append(runtime)
        ratings.append(rating)
        votes.append(votes_count)

    return {
        "names": names,
        "years_spans": years_spans,
        "genres": genres,
        "runtimes": runtimes,
        "ratings": ratings,
        "votes": votes
    }

```

Extract data for each TV series

for series in series - list

Extract series name

name = series . find ("h3") . class . strip()

series . append (name)

Extract year span

year . span = series . find ("span" = "list-item-year-muted-unbold") . text . strip()

year . span . append (year . span)

Extract genres

genre = series . find ("span", class = "genre") . text . strip()

genres . append (year . genre)

Extract run time

run-time = series . find ("span", class = "Runtimes") . text . strip()

run = times . append (run-time)

Extract rating

rating = series . find ("span", class = "ip-rating") . text . strip()

rating . append (rating)

Extract votes

vote = series . find ("span", class = "name") . "nv" . text . strip()

votes . append (vote)

Return the scraped data

return frames, year, spans, genres, run, times, rating, votes,

Get the most watched TV series details
names, year, spans, genres, runtimes,
rating, votes, = get-most-watched-tv-
series ()

Print the scraped data
for i in range(len(names)):

```
print(F"Name : {name[1:3]}")  
print(F"Year Span : {year_spans[1:3]}")  
print(F"Genre : {genres[1:3]}")  
print(F"Run time : {runtimes[1:3]}")  
print(F"Ratings : {ratings[1:3]}")  
print(F"Votes : {votes[1:3]}")  
print()
```

This script provides to the IMDB page (containing information about the most watched TV series) parses the HTML content to extract details (Names, Year, Span, Genre, Run time, Rating, Votes) for each series, and prints the scraped data.

Q3. Details of Datasets from UCI machine learning repositories. URL = <http://archive.ics.uci.edu/ml>
 You have to find the following details:

- (A) Dataset Name (B) Data Type (C) Task
- (D) Attribute Type (E) No of instance (F) year

→ Python

A copy code

```

import requests
from bs4 import BeautifulSoup

# Function to get the dataset details
def get_dataset_details():
    # URL of the UCI machine learning Repository home
    url = "http://archive.ics.uci.edu/ml"

    # Send a GET request to the URL
    response = requests.get(url)

    # Parse the HTML content
    soup = BeautifulSoup(response.text, 'html.parser')

    # Print the links to the show all Dataset' page
    all_dataset_links = soup.find('a', href='dataset')
    # Construct the URL for the show All Dataset' page
    all_dataset_url = url + all_dataset_links['href']

    # Send a GET request to the 'show' All Dataset' page
    all_dataset_response = requests.get(all_dataset_url)
    # Parse the HTML content of the show All
    # Dataset' page
    all_dataset_soup = BeautifulSoup(all_dataset_response.text, 'html.parser')
  
```

Find the table containing the dataset details
~~table = self.dataset - soup - find("table", class_ = "table")~~
~~table = self.dataset - soup - find("table", class_ = "table")~~

Initialize list to store data

dataset - names = []

data - types = []

tasks = []

attribute - types = []

num - attributes = []

years = []

Extract data for each dataset

for row in table - find_all("tr"):

columns = row - find_all("td")

dataset - name = columns[0] - text . strip()

data - type = columns[1] - text . strip()

task = columns[2] - text . strip()

attribute - type = columns[3] - text . strip()

num . instance = columns[4] - text . strip()

num . attributes = columns[5] - text . strip()

year = columns[6] - text . strip()

Append detail to respective lists

dataset - names . append (dataset - name)

data - types . append (data - type)

tasks . append (task)

attribute - type . append (attribute - type)

num . instances . append (num . instance)

num . attributes . append (num . attribute)

years . append (year)

Return the scraped data
~~dataset~~ dataset - name, data-types, tasks, attribute-types-num, instances, num_attributes, years

Get the dataset details

~~dataset~~ dataset - names, data-types, tasks, attribute-types-num, instances, num_attributes, years
= getDatasetDetails()

Print the dataset details

```
for i in range(1, len(dataset['names'])):  

    print("Dataset Name : " + dataset['name'][i])  

    print("Data Type : " + dataset['data-types'][i])  

    print("Task : " + dataset['tasks'][i])  

    print("Attribute Types : " + dataset['attribute-types'][i])  

    print("No of Instances : " + dataset['num-instances'][i])  

    print("No of Attributes : " + dataset['num-attributes'][i])  

    print("Year : " + dataset['years'][i])  

    print()
```

This script navigates to the UCI machine learning Repository homepage, finds the link to the 'Show All Dataset' page, constructs the URL for the that page, sends a GET request to it and then parses the HTML content to extract details (Dataset Name, Data Task, Attribute Type, No. of instances, No. of Attributes, Year). For each dataset, finally it prints the scraped data.