

Name : Aryan Vigyat

Registration Number : 20BCE1452

Question 1:

Given a seed/root URL, e.g., "Vit.ac.in", Design a simple crawler to return all pages (URLs) that contains a keyword "research" from this site. (25 pages)

```
In [ ]: import re
from bs4 import BeautifulSoup as soup
import requests
seed_url="http://www.vit.ac.in"
searchWord="research"
maxPages=25
response=requests.get(seed_url)
print("Response Status Code :",response.status_code)
rootPage=soup(response.content,'html.parser')
atags=rootPage.find_all('a')
result=list()
pages=0
for atag in atags:
    if(pages==maxPages):
        break
    link=atag['href']
    if re.search(searchWord,link,re.IGNORECASE):
        result.append(link)
    pages+=1

print("The Links in the SEED URL Page which consists of the word research are as follows")
for res in result:
    print("\t",res)
```

```
Response Status Code : 200
The Links in the SEED URL Page which consists of the word research are as follows
    https://vit.ac.in/admissions/research
    https://vit.ac.in/research
    https://vit.ac.in/research
    https://vit.ac.in/research/academic
    https://vit.ac.in/research/sponsored-research
    https://vit.ac.in/research/centers-list
    https://vit.ac.in/schools-centres-list-research-guides-2022
    3d-printing-play-major-role-mitigating-spread-covid-19-say-researchers-vi
t
    3d-printing-play-major-role-mitigating-spread-covid-19-say-researchers-vi
t
    https://vit.ac.in/research
```

Question 2:

Find documents that contain the word “admissions” and the word “international” within the URL “Vit.ac.in” using Python. (25 pages)

```
In [ ]: import requests
from bs4 import BeautifulSoup as soup
import re
seed_URL = "http://www.vit.ac.in"
searchWords = ['admissions', 'international']
maxPages=25
response=requests.get(seed_URL)
print("Response Status Code : ", response.status_code)
rootPage=soup(response.content, 'html.parser')
atags=rootPage.find_all('a')
validLinks=list()
for atag in atags:
    link=atag['href']
    if link.startswith("http") :
        if link not in validLinks :
            validLinks.append(link)
print("Total Number of Documents is {}".format(len(validLinks)))
final=list()
foundPages=0
failed=list()
pages=0
for link in validLinks :
    if(pages==maxPages):
        break
    try :
        page = requests.get(link).text
    except requests.ConnectionError :
        try :
            page = requests.get(link, verify=False).text
        except :
            failed.append(link)
            continue
    if (re.search(searchWords[0], page, re.IGNORECASE)) and (re.search(searchWords[1], page, re.IGNORECASE)):
        final.append(link)
        foundPages+=1
    pages+=1
print("The Documents that Contain the Words Admission and International is ")
for i in final:
    print("\t",i)
```

Response Status Code : 200

Total Number of Documents is 166

C:\Users\ayuar\AppData\Local\Programs\Python\Python310\lib\site-packages\urllib3\connectionpool.py:1045: InsecureRequestWarning: Unverified HTTPS request is being made to host 'chennai.vit.ac.in'. Adding certificate verification is strongly advised. See: https://urllib3.readthedocs.io/en/1.26.x/advanced-usage.html#ssl-warnings
warnings.warn(

The Documents that Contain the Words Admission and International is

```
https://vitap.ac.in/
https://vitbhopal.ac.in/
https://vit.ac.in
https://vit.ac.in/about-vit
https://vit.ac.in/about/vision-mission
https://vit.ac.in/vit-milestones
https://vit.ac.in/about/leadership
https://vit.ac.in/governance
https://vit.ac.in/about/administrative-offices
https://vit.ac.in/about/infrastructure
https://vit.ac.in/about/sustainability
https://vit.ac.in/true-green
https://vit.ac.in/about/community-outreach
https://vit.ac.in/about/communityradio
https://vit.ac.in/all-news-archieved
https://vit.ac.in/all-events
https://vit.ac.in/national-institutional-ranking-framework-nirf
https://vit.ac.in/mhrdugcaicte
https://vit.ac.in/about/news-letter
https://vit.ac.in/academics/home
https://vit.ac.in/programmes-offered-1
https://vit.ac.in/programmes-offered-2021-22
https://vit.ac.in/programmes-offered-2020-21
https://vit.ac.in/schools
```

Question 3:

Find documents that contain the word “Programme” but not the word “programming” within the URL “Vit.ac.in” using Python. (5 pages)

```
In [ ]: import re
from bs4 import BeautifulSoup as soup
import requests
seedURL="http://vit.ac.in"
searchWords=['Programme','Programming']
response=requests.get(seedURL)
maxPages=5
print("Response Status Code :",response.status_code)
rootPage=soup(response.content,'html.parser')
atags=rootPage.find_all('a')
validLinks=list()
for atag in atags:
    link=atag['href']
    if link.startswith("http") :
        if link not in validLinks :
            validLinks.append(link)
print("Total Number of Documents Linked to the Root URL is {}".format(len(validLinks)))
final=list()
foundPages=0
failed=list()
pages=0
for link in validLinks :
    if(pages==maxPages):
        break
```

```

try :
    page = requests.get(link).text
except requests.ConnectionError :
    try :
        page = requests.get(link, verify=False).text
    except :
        failed.append(link)
    continue
if (re.search(searchWords[0], page, re.IGNORECASE)) and (not re.search(searchW
    final.append(link)
    foundPages+=1
    pages+=1
print("The Documents that Contain the Word Programme and not Programming is ")
for i in final:
    print("\t",i)

```

Response Status Code : 200

Total Number of Documents Linked to the Root URL is 166

C:\Users\ayuar\AppData\Local\Programs\Python\Python310\lib\site-packages\urllib3\c
onnectionpool.py:1045: InsecureRequestWarning: Unverified HTTPS request is being m
ade to host 'chennai.vit.ac.in'. Adding certificate verification is strongly advis
ed. See: https://urllib3.readthedocs.io/en/1.26.x/advanced-usage.html#ssl-warnings
warnings.warn(

The Documents that Contain the Word Programme and not Programming is

```

https://vitap.ac.in/
https://vitbhopal.ac.in/
https://vit.ac.in
https://vit.ac.in/about-vit
https://vit.ac.in/about/vision-mission

```

Question 4:

**Write a web crawler program which takes as input a url
(Educational website) and a search key word and maximum
number of pages (15-20 Pages) to be searched and returns as
output all the web pages it searched till it found the search word
on a web page or return failure.**

```

In [ ]: import re
from bs4 import BeautifulSoup as soup
import requests
seedURL4 = input("Enter the Input URL:")
searchWord = input("Enter the Search Word: ")
maxPages = int(input("Enter the Max Pages:"))
response = requests.get(seedURL4)
print("Status of the response : ", response.status_code)
rootPage=soup(response.content,'html.parser')
atags=rootPage.find_all('a')
validLinks=list()
for atag in atags:
    try:
        link=atag['href']
        if link.startswith("http") :
            if link not in validLinks :
                validLinks.append(link)

```

```

except:
    pass
print("Total Number of Documents is {}".format(len(validLinks)))
final=list()
foundPages=0
failed=list()
pages=0
for link in validLinks :
    if(pages==maxPages):
        break
    try :
        page = requests.get(link).text
    except requests.ConnectionError :
        try :
            page = requests.get(link, verify=False).text
        except :
            failed.append(link)
            continue
    if (re.search(searchWord, page, re.IGNORECASE)):
        final.append(link)
        foundPages+=1
    pages+=1
if(foundPages==0):
    print("Failure")
else:
    print("The Documents that Contain the Word {} is {}".format(searchWord))
    for i in final:
        print("\t",i)

```

Enter the Input URL:<https://pes.edu/>
Enter the Search Word: campus
Enter the Max Pages:19
Status of the response : 200
Total Number of Documents is 66
The Documents that Contain the Word campus is
<https://pessat.com/>
<https://pessat.com>
<https://pes.edu/btech/>
<https://pes.edu/programs/>
<https://pes.edu/scholarships/>
<https://pes.edu/placements/>
<https://news.pes.edu/>
<https://news.pes.edu/11033>
<https://news.pes.edu/11031>
<https://news.pes.edu/11032>
<https://events.pes.edu/>
https://events.pes.edu/10657#new_tab
<https://www.youtube.com/watch?v=5PllbH06p50>
<https://pes.edu/bca>
<https://pes.edu/tag/new/>
<https://pes.edu/tag/undergraduate/>
<https://pes.edu/bscmedical>
<https://pes.edu/bsceconomics>
<https://pes.edu/law/>

Question 5:

Write a Python program to read the given website and extract the phone numbers and emails and contact addresses from Chennai, Amaravathi, Bhopal vit website. (5 Marks)

```
In [ ]: import re
from bs4 import BeautifulSoup as soup
import requests
seedURL5=["https://vitap.ac.in/","https://vitbhopal.ac.in/","https://chennai.vit.ac.in/"]
f = open("myfile.txt", "w")
for url in seedURL5:
    response = requests.get(url,verify=False)
    print("Status of the response : ", response.status_code)
    rootPage=soup(response.content,'html.parser')
    phpattern = re.compile(r'[7-9][0-9]{9}')
    emailpattern=re.compile(r'\b[A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Z|a-z]{2,}\b')
    phone_numbers=re.findall(phpattern,str(rootPage))
    email_ids=re.findall(emailpattern,str(rootPage))
    f.write("-----{} Email IDS and Phone Number-----\n".format(url))
    f.write(( ' '.join(email_ids)))
    f.write(( ' '.join(phone_numbers)))
    f.write("\n")
f.close()
f2 = open("myfile.txt", "r")
print(f2.read())
```

```
C:\Users\ayuar\AppData\Local\Programs\Python\Python310\lib\site-packages\urllib3\connectionpool.py:1045: InsecureRequestWarning: Unverified HTTPS request is being made to host 'vitap.ac.in'. Adding certificate verification is strongly advised. See: https://urllib3.readthedocs.io/en/1.26.x/advanced-usage.html#ssl-warnings
  warnings.warn(
```

```
Status of the response : 200
```

```
C:\Users\ayuar\AppData\Local\Programs\Python\Python310\lib\site-packages\urllib3\connectionpool.py:1045: InsecureRequestWarning: Unverified HTTPS request is being made to host 'vitbhopal.ac.in'. Adding certificate verification is strongly advised. See: https://urllib3.readthedocs.io/en/1.26.x/advanced-usage.html#ssl-warnings
  warnings.warn(
```

```
Status of the response : 200
```

```
C:\Users\ayuar\AppData\Local\Programs\Python\Python310\lib\site-packages\urllib3\connectionpool.py:1045: InsecureRequestWarning: Unverified HTTPS request is being made to host 'chennai.vit.ac.in'. Adding certificate verification is strongly advised. See: https://urllib3.readthedocs.io/en/1.26.x/advanced-usage.html#ssl-warnings
  warnings.warn(
```

Status of the response : 200

-----<https://vitap.ac.in/> Email IDs and Phone Number-----

7868934148 7868934148 7868934148 7797590012 9520281387

-----<https://vitbhopal.ac.in/> Email IDS and Phone Number-----

-----<https://chennai.vit.ac.in/> Email IDS and Phone Number-----
cw.cc@vit.ac.in wlh.cc@vit.ac.in transport.cc@vit.ac.in cw.cc@vit.ac.in wlh.cc@vit.ac.in transport.cc@vit.ac.in admin.chennai@vit.ac.in 8272122876 8272215041 8272239968 8272270359 8272298574 8659071608 8588684987 8679772278 8634261762 8634270077 9038225429 7121194580 8634261762 8588708437 8634261762 8588708437 8201114257 820114257 8588708437 8588708437 8634261762 859971350 9315459946 8679772278 7358782569 8272122876 8634261762 9038225429 8634270077 7121194580 7358782569 8634261762 8272215041 8588708437 8272239968 8634261762 8272270359 8588708437 8272298574 8201114257 8201114257 8588708437 8659071608 8588708437 8659971350 9315459946 8588684987