

ASSIGNMENT 1

WEB SCRAPING

In [4]:

```
# 1.PHYTHON PROGRAM TO DISPLAY ALL THE HEADER TAGS FROM WIKIPEDIA.ORG

#install libraries

!pip install bs4
!pip install requests

from bs4 import BeautifulSoup
import requests

#send get requests to the webpage server to get the source code of the page

page = requests.get('https://www.wikipedia.org/')

page

#page content
soup = BeautifulSoup(page.content)

soup

# creating a list of all common heading tags
heading_tags = ['h1','h2','h3']
for tags in soup.find_all(heading_tags):
    print(tags.name + ' -> ' + tags.text.strip())

Requirement already satisfied: bs4 in c:\users\dell\anaconda3\lib\site-packages (0.0.1)
Requirement already satisfied: beautifulsoup4 in c:\users\dell\anaconda3\lib\site-packages (from bs4) (4.11.1)
Requirement already satisfied: soupsieve>1.2 in c:\users\dell\anaconda3\lib\site-packages (from beautifulsoup4->bs4) (2.3.1)
Requirement already satisfied: requests in c:\users\dell\anaconda3\lib\site-packages (2.28.1)
Requirement already satisfied: certifi<=2017.4.17 in c:\users\dell\anaconda3\lib\site-packages (from requests) (2022.9.14)
Requirement already satisfied: idna<4,>=2.5 in c:\users\dell\anaconda3\lib\site-packages (from requests) (3.3)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in c:\users\dell\anaconda3\lib\site-packages (from requests) (1.26.11)
Requirement already satisfied: charset-normalizer<3,>=2 in c:\users\dell\anaconda3\lib\site-packages (from requests) (2.0.4)
h1 -> Wikipedia

The Free Encyclopedia
h2 -> 1 000 000+

articles
h2 -> 100 000+

articles
h2 -> 10 000+

articles
h2 -> 1 000+
```

In [19]:

```
# 2.PROGRAM TO DISPLAY IMBD'S TOP 100 RATED MOVIES

# install libraries
!pip install bs4
!pip install requests

from bs4 import BeautifulSoup
import requests

#send get requests to the webpage server to get the source code of the page

page = requests.get('https://www.imdb.com/list/ls091520106/')
page

# page content

soup=BeautifulSoup(page.content)
soup

# creating a list of top rated imbd's 100 movies

movie_title = soup.find('h3',class_='lister-item-header')
movie_title

movie_title.text

movie_year = soup.find('span',class_='lister-item-year text-muted unbold')
movie_year

movie_year.text

rating = soup.find('span',class_='ipl-rating-star__rating')
rating

rating.text

movie=[]

for i in soup.find_all('h3',class_='lister-item-header'):
    movie.append(i.text)

movie

movie_year=[]

for i in soup.find_all('span',class_='lister-item-year text-muted unbold'):
    movie_year.append(i.text)

movie_year

rating=[]

for i in soup.find_all('div',class_='ipl-rating-star small'):
    rating.append(i.text)

rating

print(len(movie),len(movie_year),len(rating))          #printing length

# making dataframe

import pandas as pd

df=pd.DataFrame({'movie':movie, 'movie_year':movie_year, 'rating':rating})
df

Requirement already satisfied: bs4 in c:\users\dell\anaconda3\lib\site-packages (0.0.1)
Requirement already satisfied: beautifulsoup4 in c:\users\dell\anaconda3\lib\site-packages (from bs4) (4.11.1)
Requirement already satisfied: soupsieve>1.2 in c:\users\dell\anaconda3\lib\site-packages (from beautifulsoup4->bs4) (2.3.1)
Requirement already satisfied: requests in c:\users\dell\anaconda3\lib\site-packages (2.28.1)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in c:\users\dell\anaconda3\lib\site-packages (from requests) (1.26.11)
Requirement already satisfied: charset-normalizer<3,>=2 in c:\users\dell\anaconda3\lib\site-packages (from requests) (2.0.4)
Requirement already satisfied: idna<4,>=2.5 in c:\users\dell\anaconda3\lib\site-packages (from requests) (3.3)
Requirement already satisfied: certifi<=2017.4.17 in c:\users\dell\anaconda3\lib\site-packages (from requests) (2022.9.14)
100 100 100
```

Out[19]:

	movie	movie_year	rating
0	1n1.nThe Shawshank Redemption1n	(1994)	1n1n1n1n1n1n9.3n
1	1n2.nThe Godfather1n	(1972)	1n1n1n1n1n1n9.2n
2	1n3.nThe Godfather Part II1n	(1974)	1n1n1n1n1n1n9n
3	1n4.nThe Dark Knight1n	(2008)	1n1n1n1n1n1n9n
4	1n5.n12 Angry Men1n	(1957)	1n1n1n1n1n1n9n
...
95	1n96.nNorth by Northwest1n	(1959)	1n1n1n1n1n1n8.3n
96	1n97.nA Clockwork Orange1n	(1971)	1n1n1n1n1n1n8.3n
97	1n98.nSnatch1n	(2000)	1n1n1n1n1n1n8.2n
98	1n99.nLe fabuleux destin d'Amélie Poulain1n	(2001)	1n1n1n1n1n1n8.3n
99	1n100.nThe Kid1n	(1921)	1n1n1n1n1n1n8.3n

100 rows × 3 columns

In [20]:

```
# 3.PROGRAM TO DISPLAY IMBD'S TOP 100 RATED INDIAN MOVIES

# install libraries
!pip install bs4
!pip install requests

from bs4 import BeautifulSoup
import requests

#send get requests to the webpage server to get the source code of the page

page = requests.get('https://www.imdb.com/list/ls009997493/')
page

# page content

soup=BeautifulSoup(page.content)
soup

# creating a list of top rated imbd's 100 movies

movie_title = soup.find('h3',class_='lister-item-header')
movie_title

movie_title.text

movie_year = soup.find('span',class_='lister-item-year text-muted unbold')
movie_year

movie_year.text

rating = soup.find('span',class_='ipl-rating-star__rating')
rating

rating.text

movie=[]

for i in soup.find_all('h3',class_='lister-item-header'):
    movie.append(i.text)

movie

movie_year=[]

for i in soup.find_all('span',class_='lister-item-year text-muted unbold'):
    movie_year.append(i.text)

movie_year

rating=[]

for i in soup.find_all('div',class_='ipl-rating-star small'):
    rating.append(i.text)

rating

print(len(movie),len(movie_year),len(rating))          #printing length

# making dataframe

import pandas as pd

df=pd.DataFrame({'movie':movie, 'movie_year':movie_year, 'rating':rating})
df

Requirement already satisfied: bs4 in c:\users\dell\anaconda3\lib\site-packages (0.0.1)
Requirement already satisfied: beautifulsoup4 in c:\users\dell\anaconda3\lib\site-packages (from bs4) (4.11.1)
Requirement already satisfied: soupsieve>1.2 in c:\users\dell\anaconda3\lib\site-packages (from beautifulsoup4->bs4) (2.3.1)
Requirement already satisfied: requests in c:\users\dell\anaconda3\lib\site-packages (2.28.1)
Requirement already satisfied: idna<4,>=2.5 in c:\users\dell\anaconda3\lib\site-packages (from requests) (3.3)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in c:\users\dell\anaconda3\lib\site-packages (from requests) (1.26.11)
Requirement already satisfied: certifi<=2017.4.17 in c:\users\dell\anaconda3\lib\site-packages (from requests) (2022.9.14)
Requirement already satisfied: charset-normalizer<3,>=2 in c:\users\dell\anaconda3\lib\site-packages (from requests) (2.0.4)
100 100 100
```

Out[20]:

	movie	movie_year	rating
0	1n1.nRang De Basanth1n	(2006)	1n1n1n1n1n1n8.1n
1	1n2.n3 Idiots1n	(2009)	1n1n1n1n1n1n8.4n
2	1n3.nTaare Zameen Par1n	(2007)	1n1n1n1n1n1n8.4n
3	1n4.nDil Chahta Hai1n	(2001)	1n1n1n1n1n1n8.1n
4	1n5.nSwades: We, the People1n	(2004)	1n1n1n1n1n1n8.2n
...
95	1n96.nWake Up Sid1n	(2009)	1n1n1n1n1n1n7.6n
96	1n97.nRangeela1n	(1995)	1n1n1n1n1n1n7.4n
97	1n98.nShatranj Ke Khilari1n	(1977)	1n1n1n1n1n1n7.5n
98	1n99.nPyaar Ka PUNCHnama1n	(2011)	1n1n1n1n1n1n7.6n
99	1n100.nEk Hasina Thi1n	(2004)	1n1n1n1n1n1n7.5n

100 rows × 3 columns

In [21]:

```
# 4.PROGRAM TO DISPLAY LIST OF RESPECTED FORMER PRESIDENT OF INDIA

#install libraries
!pip install bs4
!pip install requests

from bs4 import BeautifulSoup
import requests

#send get requests to the webpage server to get the source code of the page

page = requests.get('https://presidentofindia.nic.in/former-presidents.htm')
page

#page content

soup=BeautifulSoup(page.content)
soup

#creating a list of respected former president of india

name = soup.find('div',class_='presidentListing')
name

name.text

for i in soup.find_all('div',class_='presidentListing'):
    name.append(i.text)

name

Requirement already satisfied: bs4 in c:\users\dell\anaconda3\lib\site-packages (0.0.1)
Requirement already satisfied: beautifulsoup4 in c:\users\dell\anaconda3\lib\site-packages (from bs4) (4.11.1)
Requirement already satisfied: soupsieve>1.2 in c:\users\dell\anaconda3\lib\site-packages (from beautifulsoup4->bs4) (2.3.1)
Requirement already satisfied: requests in c:\users\dell\anaconda3\lib\site-packages (2.28.1)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in c:\users\dell\anaconda3\lib\site-packages (from requests) (1.26.11)
Requirement already satisfied: certifi<=2017.4.17 in c:\users\dell\anaconda3\lib\site-packages (from requests) (2022.9.14)
Requirement already satisfied: charset-normalizer<3,>=2 in c:\users\dell\anaconda3\lib\site-packages (from requests) (2.0.4)

<div class="presidentListing">
<h3>Shri Ram Nath Kovind (birth - 1945)</h3>
<p><span class="terms">Term of Office: 25 July, 2017 to 25 July, 2022 </span> </p>
<p><a href="https://ramnathkovind.nic.in" target="_blank">https://ramnathkovind.nic.in</a></p>

Shri Ram Nath Kovind (birth - 1945)
Term of Office: 25 July, 2017 to 25 July, 2022
https://ramnathkovind.nic.in

Shri Pranab Mukherjee (1935-2020)
Term of Office: 25 July, 2012 to 25 July, 2017
http://pranabmukherjee.nic.in

Smt Pratibha Devisingh Patil (birth - 1934)
Term of Office: 25 July, 2007 to 25 July, 2012
http://pratibhapatil.nic.in

DR. A.P.J. Abdul Kalam (1931-2015)
Term of Office: 25 July, 2002 to 25 July, 2007
http://abdulkalam.nic.in

Shri K. R. Narayanan (1920 - 2005)
Term of Office: 25 July, 1997 to 25 July, 2002

Dr Shankar Dayal Sharma (1918-1999)
Term of Office: 25 July, 1992 to 25 July, 1997

Shri R Venkataraman (1910-2009)
Term of Office: 25 July, 1987 to 25 July, 1992

Giani Zail Singh (1916-1994)
Term of Office: 25 July, 1982 to 25 July, 1987

Shri Neelam Sanjiva Reddy (1913-1996)
Term of Office: 25 July, 1977 to 25 July, 1982

Dr. Fakhruddin Ali Ahmed (1905-1977)
Term of Office: 24 August, 1974 to 11 February, 1977

Shri Varahagiri Venkata Giri (1894-1980)
Term of Office: 3 May, 1969 to 20 July, 1969 and 24 August, 1969 to 24 August, 1974

Dr. Zakir Husain (1897-1969)
Term of Office: 13 May, 1967 to 3 May, 1969

Dr. Sarvepalli Radhakrishnan (1888-1975)
Term of Office: 13 May, 1962 to 13 May, 1967

Dr. Rajendra Prasad (1884-1963)
Term of Office: 26 January, 1950 to 13 May, 1962
</div>
```

Out[21]:

In []: