# HOA 7.2: Webscraping using BeasutifulSoap and Request

**Name:** Jann Moises Nyll B. De los Reyes

**Section:** 22S3 2---

# Data Gathering

## Sources of Data

a vast amount of historical data can be found in files such as:

- MS Word Documents
- Emails
- Spreadsheets
- MS Powerpoint
- PDFs
- HTML
- and plaintext files

Publice and Private Archives

CSV,JSON, and XML files use plaintext, a common format, and are compatible with a wide range of applications

The Web can be mined for data using a web scraping application

The IoT uses sensors create data

Sensors in smartphones, cars, airplanes, street lamps, and home appliance capture raw data

## Open data and Private Data

1. Open Data

    The Open Knowledge Foundation describes Open Data as "any content,information or data that people are free to use, reuse, and redistribute without any legal, technological, or social restriction.

2. Private Data

    Data related to an expectation of privacy and regulated by particular country/government

## Structure and Unstructured Data

1. Structured Data

    Data entered and maintained in fixed field within a file or record. Easily entered,classified,queried, and analyzed Realational databases or spreadsheets

2. Unstructured Data (Lacks organization)

    Raw data photo contents, audio, video, webpages, blogs, books, journals, white papers, PowerPoint presentations, articles, emails, wikis, word processing documents, and text in general

S3

# Example of gathering image using webcam

Note: Run this snippet using local jupyter notebook

```python
In [2]:  import cv2
         #from google.colab.patches import cv2_imshow # make this a comment to run the file
         key = cv2.waitKey(1)
         webcam = cv2.VideoCapture(0)

         while True:
           try:
             check, frame = webcam.read()
             print(check) #prints true as long as the webcam is running
             print(frame) #prints matrix values of each framecd
             cv2.imshow("Capturing", frame)
             key = cv2.waitKey(1)
             if key == ord('s'):

               cv2.imwrite(filename = 'saved_img.jpg', img = frame)
               webcam.release()
               img_new = cv2.imread('saved_img.jpg', cv2.IMREAD_GRAYSCALE)
               img_new = cv2.imshow("Captured Image", img_new)
               cv2.waitKey(1650)
               cv2.destroyAllWindows()
               print("Processing image ...")
               img_ =cv2.imread('saved_img.jpg',cv2.IMREAD_ANYCOLOR)
               print("Converting RGB image to grayscale ...")
               gray = cv2.cvtColor(img_,cv2.COLOR_BGR2GRAY)
               print("Converted RGB image to grayscale ...")
               print("Resizing image to 28 x28 scale ...")
               img_ = cv2.resize(gray,(28,28))
               print("Resized ...")
               img_resized = cv2.imwrite(filename = 'saved_img-final.jpg', img = img_)
               print("Image saved!")

               break
             elif key == ord('q'):
               print("Turning off camera.")
               webcam.release()
               print("Camera off.")
               print("Program ended.")
               cv2.destroyAllWindows()
               break

           except(KeyboardInterrupt):
             print("Turning off camera.")
             webcam.release()
             print("Camera off.")
             print("Program ended.")
             cv2.destroyAllWindows()
             break
```

```
True
[[[2 2 4]
  [2 2 4]
  [2 2 4]
  ...
  [2 2 4]
  [2 2 4]
  [2 2 4]]

 [[2 2 4]
  [2 2 4]
  [2 2 4]
  ...
  [2 2 4]
  [2 2 4]
  [2 2 4]]

 [[2 2 4]
  [2 2 4]
  [2 2 4]
  ...
  [2 2 4]
  [2 2 4]
  [2 2 4]]

 ...

 [[2 2 4]
  [2 2 4]
  [2 2 4]
  ...
  [2 2 4]
  [2 2 4]
  [2 2 4]]

 [[2 2 4]
  [2 2 4]
  [2 2 4]
  ...
  [2 2 4]
  [2 2 4]
  [2 2 4]]

 [[2 2 4]
  [2 2 4]
  [2 2 4]
  ...
  [2 2 4]
  [2 2 4]
  [2 2 4]]]
True
[[[2 2 4]
  [2 2 4]
  [2 2 4]
  ...
  [2 2 4]
  [2 2 4]
  [2 2 4]]

 [[2 2 4]
  [2 2 4]
  [2 2 4]
  ...
  [2 2 4]
  [2 2 4]
  [2 2 4]]

 [[2 2 4]
  [2 2 4]
  [2 2 4]
```

```
 [ 14  11  12]
 [ 14  11  12]]

[[ 51  40  66]
 [ 51  40  66]
 [ 52  41  67]
 ...
 [ 17  16  18]
 [ 17  14  17]
 [ 17  14  17]]]
Processing image ...
Converting RGB image to grayscale ...
Converted RGB image to grayscale ...
Resizing image to 28 x28 scale ...
Resized ...
Image saved!
```

## Example of gathering voice data using microphone

Note: Run the snippet of codes using local jupyter notebook

In [3]:
```
!pip3 install sounddevice
```

```
Requirement already satisfied: sounddevice in c:\users\lenovo\anaconda3\lib\site-packages (0.4.6)
Requirement already satisfied: CFFI>=1.0 in c:\users\lenovo\anaconda3\lib\site-packages (from sounddevice)
(1.16.0)
Requirement already satisfied: pycparser in c:\users\lenovo\anaconda3\lib\site-packages (from CFFI>=1.0->soun
ddevice) (2.21)
```

In [4]:
```
!pip3 install wavio
```

```
Requirement already satisfied: wavio in c:\users\lenovo\anaconda3\lib\site-packages (0.0.8)
Requirement already satisfied: numpy>=1.19.0 in c:\users\lenovo\anaconda3\lib\site-packages (from wavio) (1.2
6.2)
```

In [5]:
```
!pip3 install scipy
```

```
Requirement already satisfied: scipy in c:\users\lenovo\anaconda3\lib\site-packages (1.11.4)
Requirement already satisfied: numpy<1.28.0,>=1.21.6 in c:\users\lenovo\anaconda3\lib\site-packages (from sci
py) (1.26.2)
```

In [18]:
```
#!apt-get install libportaudio2
```

The command above is specific to Debian-based Linux Distribution and cannot directly run in a Jupyter notebook on Windows. Instead we can use `pyaudio` library which uses `portaudio`, in Jupyter notebook on Windows.

In [21]:
```
!pip3 install pyaudio
```

```
Requirement already satisfied: pyaudio in c:\users\lenovo\anaconda3\lib\site-packages (0.2.14)
```

In [73]:
```python
#import required libraries
import sounddevice as sd
from scipy.io.wavfile import write
import wavio as wv

#sampling frequency
freq = 44100

#Recording duration
duration = 5

#start recorder with the given values
# of duration and sample frequency
recording =sd.rec(int(duration * freq),
                  samplerate =freq, channels= 2)
#Record audio for the given number of seconds
sd.wait()

#This will convert the NumPy array to anaudio
# This will the given sampling frequency
```

```
write("recording0.wav",freq,recording)

#Convert the NumPy array to audio file
wv.write("recording1.wav",recording, freq, sampwidth =2)
```

## Web Scraping

**Web Scraping, Web Harvesting, or web data extraction** is data scraping used for extracting data from websites. The web scraping software may directly acces the World Wide Web using the Hypertext Transfer Protocol or a web browser. While web scraping can be done manually by a software user, the term typically refers to automated processes implemented using a bot or a web crawler. It is a form of copying in which specific data gathered and copied from the web, typically into a central local database or spreadsheet, for later retrieval or analysis.

Reference: link text

## Image Scraping using BeautifulSoup and Request

In [22]:
```
!pip install bs4
```

```
Collecting bs4
  Downloading bs4-0.0.2-py2.py3-none-any.whl.metadata (411 bytes)
Requirement already satisfied: beautifulsoup4 in c:\users\lenovo\anaconda3\lib\site-packages (from bs4) (4.1
2.2)
Requirement already satisfied: soupsieve>1.2 in c:\users\lenovo\anaconda3\lib\site-packages (from beautifulso
up4->bs4) (2.5)
Downloading bs4-0.0.2-py2.py3-none-any.whl (1.2 kB)
Installing collected packages: bs4
Successfully installed bs4-0.0.2
```

In [24]:
```
pip install requests
```

```
Requirement already satisfied: requests in c:\users\lenovo\anaconda3\lib\site-packages (2.31.0)
Requirement already satisfied: charset-normalizer<4,>=2 in c:\users\lenovo\anaconda3\lib\site-packages (from
requests) (2.0.4)
Requirement already satisfied: idna<4,>=2.5 in c:\users\lenovo\anaconda3\lib\site-packages (from requests)
(3.4)
Requirement already satisfied: urllib3<3,>=1.21.1 in c:\users\lenovo\anaconda3\lib\site-packages (from reques
ts) (1.26.18)
Requirement already satisfied: certifi>=2017.4.17 in c:\users\lenovo\anaconda3\lib\site-packages (from reques
ts) (2024.2.2)
Note: you may need to restart the kernel to use updated packages.
```

In [27]:
```
import requests
from bs4 import BeautifulSoup

def getdata(url):
    r = requests.get(url)
    return r.text

htmldata =getdata("https://www.google.com/")
soup = BeautifulSoup(htmldata, 'html.parser')
for item in soup.find_all('img'):
    print(item['src'])
```

```
/images/branding/googlelogo/1x/googlelogo_white_background_color_272x92dp.png
```

In [28]:
```
pip install selenium
```

```
Collecting selenium
  Downloading selenium-4.18.1-py3-none-any.whl.metadata (6.9 kB)
Requirement already satisfied: urllib3<3,>=1.26 in c:\users\lenovo\anaconda3\lib\site-packages (from urllib3
[socks]<3,>=1.26->selenium) (1.26.18)
Collecting trio~=0.17 (from selenium)
  Downloading trio-0.25.0-py3-none-any.whl.metadata (8.7 kB)
Collecting trio-websocket~=0.9 (from selenium)
  Downloading trio_websocket-0.11.1-py3-none-any.whl.metadata (4.7 kB)
Requirement already satisfied: certifi>=2021.10.8 in c:\users\lenovo\anaconda3\lib\site-packages (from seleni
um) (2024.2.2)
Collecting typing_extensions>=4.9.0 (from selenium)
  Downloading typing_extensions-4.10.0-py3-none-any.whl.metadata (3.0 kB)
Collecting attrs>=23.2.0 (from trio~=0.17->selenium)
  Downloading attrs-23.2.0-py3-none-any.whl.metadata (9.5 kB)
Requirement already satisfied: sortedcontainers in c:\users\lenovo\anaconda3\lib\site-packages (from trio~=0.
17->selenium) (2.4.0)
Requirement already satisfied: idna in c:\users\lenovo\anaconda3\lib\site-packages (from trio~=0.17->seleniu
m) (3.4)
Collecting outcome (from trio~=0.17->selenium)
  Downloading outcome-1.3.0.post0-py2.py3-none-any.whl.metadata (2.6 kB)
Collecting sniffio>=1.3.0 (from trio~=0.17->selenium)
  Downloading sniffio-1.3.1-py3-none-any.whl.metadata (3.9 kB)
Requirement already satisfied: cffi>=1.14 in c:\users\lenovo\anaconda3\lib\site-packages (from trio~=0.17->se
lenium) (1.16.0)
Requirement already satisfied: exceptiongroup in c:\users\lenovo\anaconda3\lib\site-packages (from trio~=0.17
->selenium) (1.0.4)
Collecting wsproto>=0.14 (from trio-websocket~=0.9->selenium)
  Downloading wsproto-1.2.0-py3-none-any.whl.metadata (5.6 kB)
Requirement already satisfied: PySocks!=1.5.7,<2.0,>=1.5.6 in c:\users\lenovo\anaconda3\lib\site-packages (fr
om urllib3[socks]<3,>=1.26->selenium) (1.7.1)
Requirement already satisfied: pycparser in c:\users\lenovo\anaconda3\lib\site-packages (from cffi>=1.14->tri
o~=0.17->selenium) (2.21)
Collecting h11<1,>=0.9.0 (from wsproto>=0.14->trio-websocket~=0.9->selenium)
  Downloading h11-0.14.0-py3-none-any.whl.metadata (8.2 kB)
Downloading selenium-4.18.1-py3-none-any.whl (10.0 MB)
   ---------------------------------------- 0.0/10.0 MB ? eta -:--:--
   - -------------------------------------- 0.3/10.0 MB 9.6 MB/s eta 0:00:02
   -- ------------------------------------- 0.7/10.0 MB 9.2 MB/s eta 0:00:02
   ---- ----------------------------------- 1.1/10.0 MB 8.7 MB/s eta 0:00:02
   ----- ---------------------------------- 1.5/10.0 MB 8.6 MB/s eta 0:00:01
   ------- -------------------------------- 2.0/10.0 MB 9.0 MB/s eta 0:00:01
   ---------- ----------------------------- 2.5/10.0 MB 9.5 MB/s eta 0:00:01
   ----------- ---------------------------- 2.8/10.0 MB 9.3 MB/s eta 0:00:01
   ------------ --------------------------- 3.1/10.0 MB 8.6 MB/s eta 0:00:01
   -------------- ------------------------- 3.6/10.0 MB 8.9 MB/s eta 0:00:01
   -------------- ------------------------- 3.9/10.0 MB 8.8 MB/s eta 0:00:01
   ---------------- ----------------------- 4.4/10.0 MB 8.8 MB/s eta 0:00:01
   ------------------ --------------------- 4.9/10.0 MB 8.9 MB/s eta 0:00:01
   -------------------- ------------------- 5.3/10.0 MB 8.8 MB/s eta 0:00:01
   ---------------------- ----------------- 5.7/10.0 MB 8.9 MB/s eta 0:00:01
   ---------------------- ----------------- 5.9/10.0 MB 8.6 MB/s eta 0:00:01
   ---------------------- ----------------- 6.1/10.0 MB 8.3 MB/s eta 0:00:01
   ------------------------ --------------- 6.5/10.0 MB 8.2 MB/s eta 0:00:01
   -------------------------- ------------- 6.8/10.0 MB 8.3 MB/s eta 0:00:01
   -------------------------- ------------- 7.1/10.0 MB 8.2 MB/s eta 0:00:01
   ---------------------------- ----------- 7.5/10.0 MB 8.2 MB/s eta 0:00:01
   ----------------------------- ---------- 7.7/10.0 MB 8.0 MB/s eta 0:00:01
   ------------------------------- -------- 8.2/10.0 MB 8.0 MB/s eta 0:00:01
   -------------------------------- ------- 8.5/10.0 MB 8.1 MB/s eta 0:00:01
   ---------------------------------- ----- 8.8/10.0 MB 8.1 MB/s eta 0:00:01
   ------------------------------------ --- 9.2/10.0 MB 8.0 MB/s eta 0:00:01
   -------------------------------------- - 9.5/10.0 MB 8.0 MB/s eta 0:00:01
   ---------------------------------------- 10.0/10.0 MB 8.0 MB/s eta 0:00:01
   ---------------------------------------- 10.0/10.0 MB 7.9 MB/s eta 0:00:00
Downloading trio-0.25.0-py3-none-any.whl (467 kB)
   ---------------------------------------- 0.0/467.2 kB ? eta -:--:--
   ------------------------- ------------- 307.2/467.2 kB 9.6 MB/s eta 0:00:01
   ---------------------------------------- 467.2/467.2 kB 7.4 MB/s eta 0:00:00
Downloading trio_websocket-0.11.1-py3-none-any.whl (17 kB)
Downloading typing_extensions-4.10.0-py3-none-any.whl (33 kB)
```

```
Downloading attrs-23.2.0-py3-none-any.whl (60 kB)
   ---------------------------------------- 0.0/60.8 kB ? eta -:--:--
   ---------------------------------------- 60.8/60.8 kB 3.2 MB/s eta 0:00:00
Downloading sniffio-1.3.1-py3-none-any.whl (10 kB)
Downloading wsproto-1.2.0-py3-none-any.whl (24 kB)
Downloading outcome-1.3.0.post0-py2.py3-none-any.whl (10 kB)
Downloading h11-0.14.0-py3-none-any.whl (58 kB)
   ---------------------------------------- 0.0/58.3 kB ? eta -:--:--
   ---------------------------------------- 58.3/58.3 kB 3.2 MB/s eta 0:00:00
Installing collected packages: typing_extensions, sniffio, h11, attrs, wsproto, outcome, trio, trio-websocke
t, selenium
  Attempting uninstall: typing_extensions
    Found existing installation: typing_extensions 4.7.1
    Uninstalling typing_extensions-4.7.1:
      Successfully uninstalled typing_extensions-4.7.1
  Attempting uninstall: sniffio
    Found existing installation: sniffio 1.2.0
    Uninstalling sniffio-1.2.0:
      Successfully uninstalled sniffio-1.2.0
  Attempting uninstall: attrs
    Found existing installation: attrs 23.1.0
    Uninstalling attrs-23.1.0:
      Successfully uninstalled attrs-23.1.0
Successfully installed attrs-23.2.0 h11-0.14.0 outcome-1.3.0.post0 selenium-4.18.1 sniffio-1.3.1 trio-0.25.0
trio-websocket-0.11.1 typing_extensions-4.10.0 wsproto-1.2.0
Note: you may need to restart the kernel to use updated packages.
```
```
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. Th
is behaviour is the source of the following dependency conflicts.
python-lsp-black 1.2.1 requires black>=22.3.0, but you have black 0.0 which is incompatible.
```

## Image Scraping using Selenium

Note: Run the snippet of code using local jupyter notebook

In [133...

```python
!pip install selenium
import sys
sys.path.insert(0,'/usr/lib/chromium-browser/chromedriver')


from selenium import webdriver
from selenium.webdriver.common.by import By
import time
import requests
import shutil
import os
import getpass
import urllib.request
import io
import time
from PIL import Image
user = getpass.getuser()
chrome_options = webdriver.ChromeOptions()
chrome_options.add_argument('--headless')
chrome_options.add_argument('--no-sandbox')
chrome_options.add_argument('--disable-dev-shm-usage')
driver = webdriver.Chrome()
def scroll_to_end(driver):
    driver.execute_script("window.scrollTo(0, document.body.scrollHeight);")
    time.sleep(5)#sleep_between_interactions


def getImageUrls(name,totalImgs,driver):
    search_url = "https://www.google.com/search?q=cat&tbm=isch&ved=2ahUKEwjNn_Gn7YyFAxU3yDgGHQYQCesQ2-cCegQI
    driver.get(search_url)
    img_urls = set()
    img_count = 0
    results_start = 0
```

```python
    while(img_count+results_start<totalImgs): #Extract actual images now
        scroll_to_end(driver)
        totalResults = driver.find_elements(By.CLASS_NAME,"Q4LuWd")
        print('total results:', len(totalResults))
        print(f"Found: {totalResults} search results. Extracting links from{results_start}:{totalResults}")
        for img in totalResults[results_start:totalImgs]:
            img.click()
            time.sleep(5)
            image = driver.find_element(By.CLASS_NAME,'iPVvYb')
            img_urls.add(image.get_attribute('src'))
            print(img_urls)
            img_count=len(img_urls)
            print(img_count)

    return img_urls

def downloadImages(folder_path,file_name,url):
    try:
        image_content = requests.get(url).content
    except Exception as e:
        print(f"ERROR - COULD NOT DOWNLOAD {url} - {e}")
    try:
        image_file = io.BytesIO(image_content)
        image = Image.open(image_file).convert('RGB')
        file_path = os.path.join(folder_path, file_name)
        with open(file_path, 'wb') as f:
            image.save(f, "JPEG", quality=85)
        print(f"SAVED - {url} - AT: {file_path}")
    except Exception as e:
        print(f"ERROR - COULD NOT SAVE {url} - {e}")

def saveInDestFolder(searchNames,destDir,totalImgs,driver):
    for name in list(searchNames):
        path=os.path.join(destDir,name)
        if not os.path.isdir(path):
            os.mkdir(path)
        print('Current Path',path)
        totalLinks=getImageUrls(name,totalImgs,driver)
        print('totalLinks',totalLinks)

    if totalLinks is None:
        print('images not found for :',name)

    else:
        for i, link in enumerate(totalLinks):
            file_name = f"{i:150}.jpg"
            downloadImages(path,file_name,link)

searchNames=['cat']
destDir=f'C:/Users/Lenovo/Downloads/HOA7.2 Webscraping using BeautifulSoap and Request'
totalImgs=5

saveInDestFolder(searchNames,destDir,totalImgs,driver)
```

```
Requirement already satisfied: selenium in c:\users\lenovo\anaconda3\lib\site-packages (4.18.1)
Requirement already satisfied: urllib3<3,>=1.26 in c:\users\lenovo\anaconda3\lib\site-packages (from urllib3
[socks]<3,>=1.26->selenium) (1.26.18)
Requirement already satisfied: trio~=0.17 in c:\users\lenovo\anaconda3\lib\site-packages (from selenium) (0.2
5.0)
Requirement already satisfied: trio-websocket~=0.9 in c:\users\lenovo\anaconda3\lib\site-packages (from selen
ium) (0.11.1)
Requirement already satisfied: certifi>=2021.10.8 in c:\users\lenovo\anaconda3\lib\site-packages (from seleni
um) (2024.2.2)
Requirement already satisfied: typing_extensions>=4.9.0 in c:\users\lenovo\anaconda3\lib\site-packages (from
selenium) (4.10.0)
Requirement already satisfied: attrs>=23.2.0 in c:\users\lenovo\anaconda3\lib\site-packages (from trio~=0.17-
>selenium) (23.2.0)
Requirement already satisfied: sortedcontainers in c:\users\lenovo\anaconda3\lib\site-packages (from trio~=0.
17->selenium) (2.4.0)
Requirement already satisfied: idna in c:\users\lenovo\anaconda3\lib\site-packages (from trio~=0.17->seleniu
m) (3.4)
Requirement already satisfied: outcome in c:\users\lenovo\anaconda3\lib\site-packages (from trio~=0.17->selen
ium) (1.3.0.post0)
Requirement already satisfied: sniffio>=1.3.0 in c:\users\lenovo\anaconda3\lib\site-packages (from trio~=0.17
->selenium) (1.3.1)
Requirement already satisfied: cffi>=1.14 in c:\users\lenovo\anaconda3\lib\site-packages (from trio~=0.17->se
lenium) (1.16.0)
Requirement already satisfied: exceptiongroup in c:\users\lenovo\anaconda3\lib\site-packages (from trio~=0.17
->selenium) (1.0.4)
Requirement already satisfied: wsproto>=0.14 in c:\users\lenovo\anaconda3\lib\site-packages (from trio-websoc
ket~=0.9->selenium) (1.2.0)
Requirement already satisfied: PySocks!=1.5.7,<2.0,>=1.5.6 in c:\users\lenovo\anaconda3\lib\site-packages (fr
om urllib3[socks]<3,>=1.26->selenium) (1.7.1)
Requirement already satisfied: pycparser in c:\users\lenovo\anaconda3\lib\site-packages (from cffi>=1.14->tri
o~=0.17->selenium) (2.21)
Requirement already satisfied: h11<1,>=0.9.0 in c:\users\lenovo\anaconda3\lib\site-packages (from wsproto>=0.
14->trio-websocket~=0.9->selenium) (0.14.0)
Current Path C:/Users/Lenovo/Downloads/HOA7.2 Webscraping using BeautifulSoap and Request\cat
total results: 100
Found: [<selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element
="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.370")>, <selenium.webdriver.remote.
webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE2
9.d.267976116F5ED131BF63F65BE4CFBFA9.e.384")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6
d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA
9.e.398")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", ele
ment="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.412")>, <selenium.webdriver.rem
ote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C418
8FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.426")>, <selenium.webdriver.remote.webelement.WebElement (session
="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4
CFBFA9.e.440")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638
a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.454")>, <selenium.webdri
ver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041
BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.468")>, <selenium.webdriver.remote.webelement.WebElement (s
ession="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63
F65BE4CFBFA9.e.482")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd4742
2638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.496")>, <selenium.we
bdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40
D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.510")>, <selenium.webdriver.remote.webelement.WebElemen
t (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131
BF63F65BE4CFBFA9.e.524")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd
47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.538")>, <seleniu
m.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFF
AC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.552")>, <selenium.webdriver.remote.webelement.WebEl
ement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5E
D131BF63F65BE4CFBFA9.e.566")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee7
3dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.580")>, <sel
enium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2
AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.600")>, <selenium.webdriver.remote.webelement.W
ebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.26797611
6F5ED131BF63F65BE4CFBFA9.e.614")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22f
eee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.628")>,
<selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.3854
4AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.642")>, <selenium.webdriver.remote.webeleme
nt.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.2679
```

76116F5ED131BF63F65BE4CFBFA9.e.661")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.680")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.699")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.718")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.767")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.786")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.805")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.824")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.843")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.862")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.881")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.900")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.919")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.938")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.957")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.976")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.995")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1014")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1033")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1058")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1077")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1096")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1115")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1134")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1153")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1172")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1191")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1210")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2896")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2872")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2848")>, <selenium.webdriver.remote.webelement.WebE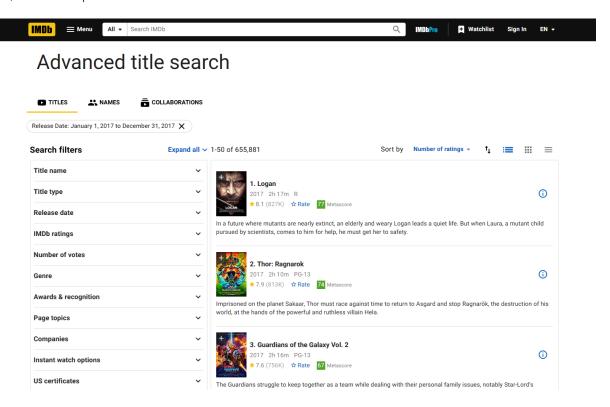lement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2824")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2800")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2776")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2752")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2728")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2698")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2674")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2650")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2626")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2602")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2578")>, <selenium.webdriver.remote.webelement.WebE

lement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2554")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2530")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2506")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2482")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2458")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2434")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2410")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2386")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2362")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2332")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2266")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2242")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2218")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2194")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2170")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2146")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2122")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2098")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2074")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2050")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2026")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2002")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1978")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1954")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1930")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1906")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1882")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1852")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1828")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1804")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1780")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1756")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1732")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1708")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1641")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1617")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1593")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1569")>] search results. Extracting links from0:[<selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.370")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.384")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.398")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22

feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.412")>,
<selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.3854
4AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.426")>, <selenium.webdriver.remote.webeleme
nt.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.2679
76116F5ED131BF63F65BE4CFBFA9.e.440")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34
f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.45
4")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element
="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.468")>, <selenium.webdriver.remote.
webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE2
9.d.267976116F5ED131BF63F65BE4CFBFA9.e.482")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6
d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA
9.e.496")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", ele
ment="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.510")>, <selenium.webdriver.rem
ote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C418
8FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.524")>, <selenium.webdriver.remote.webelement.WebElement (session
="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4
CFBFA9.e.538")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638
a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.552")>, <selenium.webdri
ver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041
BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.566")>, <selenium.webdriver.remote.webelement.WebElement (s
ession="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63
F65BE4CFBFA9.e.580")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd4742
2638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.600")>, <selenium.we
bdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40
D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.614")>, <selenium.webdriver.remote.webelement.WebElemen
t (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131
BF63F65BE4CFBFA9.e.628")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd
47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.642")>, <seleniu
m.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFF
AC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.661")>, <selenium.webdriver.remote.webelement.WebEl
ement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5E
D131BF63F65BE4CFBFA9.e.680")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee7
3dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.699")>, <sel
enium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2
AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.718")>, <selenium.webdriver.remote.webelement.W
ebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.26797611
6F5ED131BF63F65BE4CFBFA9.e.767")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22f
eee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.786")>,
<selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.3854
4AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.805")>, <selenium.webdriver.remote.webeleme
nt.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.2679
76116F5ED131BF63F65BE4CFBFA9.e.824")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34
f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.84
3")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element
="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.862")>, <selenium.webdriver.remote.
webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE2
9.d.267976116F5ED131BF63F65BE4CFBFA9.e.881")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6
d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA
9.e.900")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", ele
ment="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.919")>, <selenium.webdriver.rem
ote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C418
8FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.938")>, <selenium.webdriver.remote.webelement.WebElement (session
="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4
CFBFA9.e.957")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638
a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.976")>, <selenium.webdri
ver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041
BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.995")>, <selenium.webdriver.remote.webelement.WebElement (s
ession="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63
F65BE4CFBFA9.e.1014")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd474
22638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1033")>, <selenium.
webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC
40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1058")>, <selenium.webdriver.remote.webelement.WebEle
ment (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED
131BF63F65BE4CFBFA9.e.1077")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee7
3dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1096")>, <se
lenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF
2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1115")>, <selenium.webdriver.remote.webelemen
t.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.26797
6116F5ED131BF63F65BE4CFBFA9.e.1134")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34
f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.115

3")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1172")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1191")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1210")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2896")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2872")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2848")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2824")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2800")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2776")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2752")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2728")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2698")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2674")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2650")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2626")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2602")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2578")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2554")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2530")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2506")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2482")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2458")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2434")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2410")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2386")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2362")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2332")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2266")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2242")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2218")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2194")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2170")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2146")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2122")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2098")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2074")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2050")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2026")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.2002")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1978")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1954")>, <selenium.webdriver.remote.webelement.WebElement (session

="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4
CFBFA9.e.1930")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638
a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1906")>, <selenium.webdr
iver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D04
1BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1882")>, <selenium.webdriver.remote.webelement.WebElement
(session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF
63F65BE4CFBFA9.e.1852")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd4
7422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1828")>, <seleniu
m.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFF
AC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1804")>, <selenium.webdriver.remote.webelement.WebE
lement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5
ED131BF63F65BE4CFBFA9.e.1780")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22fee
e73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1756")>, <
selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544
AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1732")>, <selenium.webdriver.remote.webeleme
nt.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.2679
76116F5ED131BF63F65BE4CFBFA9.e.1708")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a3
4f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.164
1")>, <selenium.webdriver.remote.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element
="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1617")>, <selenium.webdriver.remot
e.webelement.WebElement (session="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188F
E29.d.267976116F5ED131BF63F65BE4CFBFA9.e.1593")>, <selenium.webdriver.remote.webelement.WebElement (session
="ec6d6be7a34f22feee73dbd47422638a", element="f.38544AF2AEFFAC40D041BD2C4188FE29.d.267976116F5ED131BF63F65BE4
CFBFA9.e.1569")>]
{'https://i.natgeofe.com/n/548467d8-c5f1-4551-9f58-6817a8d2c45e/NationalGeographic_2572187_square.jpg'}
1
{'https://i.natgeofe.com/n/548467d8-c5f1-4551-9f58-6817a8d2c45e/NationalGeographic_2572187_square.jpg', 'http
s://cdn.britannica.com/70/234870-050-D4D024BB/Orange-colored-cat-yawns-displaying-teeth.jpg'}
2
{'https://i.natgeofe.com/n/548467d8-c5f1-4551-9f58-6817a8d2c45e/NationalGeographic_2572187_square.jpg', 'http
s://upload.wikimedia.org/wikipedia/commons/thumb/1/15/Cat_August_2010-4.jpg/1200px-Cat_August_2010-4.jpg', 'h
ttps://cdn.britannica.com/70/234870-050-D4D024BB/Orange-colored-cat-yawns-displaying-teeth.jpg'}
3
{'https://cdn.britannica.com/34/235834-050-C5843610/two-different-breeds-of-cats-side-by-side-outdoors-in-the
-garden.jpg', 'https://i.natgeofe.com/n/548467d8-c5f1-4551-9f58-6817a8d2c45e/NationalGeographic_2572187_squar
e.jpg', 'https://upload.wikimedia.org/wikipedia/commons/thumb/1/15/Cat_August_2010-4.jpg/1200px-Cat_August_20
10-4.jpg', 'https://cdn.britannica.com/70/234870-050-D4D024BB/Orange-colored-cat-yawns-displaying-teeth.jpg'}
4
{'https://media.4-paws.org/5/b/4/b/5b4b5a91dd9443fa1785ee7fca66850e06dcc7f9/VIER%20PFOTEN_2019-12-13_209-2890
x2000-1920x1329.jpg', 'https://upload.wikimedia.org/wikipedia/commons/thumb/1/15/Cat_August_2010-4.jpg/1200px
-Cat_August_2010-4.jpg', 'https://i.natgeofe.com/n/548467d8-c5f1-4551-9f58-6817a8d2c45e/NationalGeographic_25
72187_square.jpg', 'https://cdn.britannica.com/70/234870-050-D4D024BB/Orange-colored-cat-yawns-displaying-tee
th.jpg', 'https://cdn.britannica.com/34/235834-050-C5843610/two-different-breeds-of-cats-side-by-side-outdoor
s-in-the-garden.jpg'}
5
totalLinks {'https://media.4-paws.org/5/b/4/b/5b4b5a91dd9443fa1785ee7fca66850e06dcc7f9/VIER%20PFOTEN_2019-12-
13_209-2890x2000-1920x1329.jpg', 'https://upload.wikimedia.org/wikipedia/commons/thumb/1/15/Cat_August_2010-
4.jpg/1200px-Cat_August_2010-4.jpg', 'https://i.natgeofe.com/n/548467d8-c5f1-4551-9f58-6817a8d2c45e/NationalG
eographic_2572187_square.jpg', 'https://cdn.britannica.com/70/234870-050-D4D024BB/Orange-colored-cat-yawns-di
splaying-teeth.jpg', 'https://cdn.britannica.com/34/235834-050-C5843610/two-different-breeds-of-cats-side-by-
side-outdoors-in-the-garden.jpg'}
SAVED - https://media.4-paws.org/5/b/4/b/5b4b5a91dd9443fa1785ee7fca66850e06dcc7f9/VIER%20PFOTEN_2019-12-13_20
9-2890x2000-1920x1329.jpg - AT: C:/Users/Lenovo/Downloads/HOA7.2 Webscraping using BeautifulSoap and Request
\cat\
0.jpg
SAVED - https://upload.wikimedia.org/wikipedia/commons/thumb/1/15/Cat_August_2010-4.jpg/1200px-Cat_August_201
0-4.jpg - AT: C:/Users/Lenovo/Downloads/HOA7.2 Webscraping using BeautifulSoap and Request\cat\
1.jpg
SAVED - https://i.natgeofe.com/n/548467d8-c5f1-4551-9f58-6817a8d2c45e/NationalGeographic_2572187_square.jpg -
AT: C:/Users/Lenovo/Downloads/HOA7.2 Webscraping using BeautifulSoap and Request\cat\
2.jpg
SAVED - https://cdn.britannica.com/70/234870-050-D4D024BB/Orange-colored-cat-yawns-displaying-teeth.jpg - AT:
C:/Users/Lenovo/Downloads/HOA7.2 Webscraping using BeautifulSoap and Request\cat\
3.jpg
SAVED - https://cdn.britannica.com/34/235834-050-C5843610/two-different-breeds-of-cats-side-by-side-outdoors-
in-the-garden.jpg - AT: C:/Users/Lenovo/Downloads/HOA7.2 Webscraping using BeautifulSoap and Request\cat\
4.jpg

## Web Scraping of Movies Information using BeautifulSoup

We want to analyze the distributions of IMDB and Metacritic movie ratings to see if we find anything interesting. To do this, we'll first scrape data for over 2000 movies.



## Identifying URL structure

In the image above you can see that the URL has several parameters after the question mark:

- release_date - Shows only the movies released in a specific year.
- sort - Sorts the movies on the page. sort = num_votes desc translates to sort by number of votes in a descending order.
- page - Specifies the page number
- ref_ - Take us to the next or previous page. The reference is the page we are currently on. adv_nxt and adv_prv are two possible values. They translate to advance to the next page, and advance to the previous page respectively

```
In [3]: from requests import get
        url = 'https://www.imdb.com/search/title/?release_date=2017-01-01,2017-12-31&sort=num_votes,desc&page=1'
        agent = {"User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chr
        response = get(url,headers=agent)
        print(response.text[:500])
```

```
<!DOCTYPE html><html lang="en-US" xmlns:og="http://opengraphprotocol.org/schema/" xmlns:fb="http://www.facebo
ok.com/2008/fbml"><head><meta charSet="utf-8"/><meta name="viewport" content="width=device-width"/><script>if
(typeof uet === 'function'){ uet('bb', 'LoadTitle', {wb: 1}); }</script><script>window.addEventListener('loa
d', (event) => {
        if (typeof window.csa !== 'undefined' && typeof window.csa === 'function') {
            var csaLatencyPlugin = window.csa('Content', {
```

Understanding the HTML structure of a single page

Using BeautifulSoup to parse the HYML content

To parse our HTML document and extract the 50 containers, we'll use a python module called BeautifulSoup, the most common web scraping module for Python.

In the following code we will:

- Import the BeautifulSoup class creator from the package bs4.
- Parse response.text by creating a BeautifulSoup object, and assign this object to html_soup. The 'html.parse' argument indicates that we want to do the parsing using Python's built-in HTML parser.

In [4]:
```python
from bs4 import BeautifulSoup
html_soup = BeautifulSoup(response.text,'html.parser')
headers = {'Accept-Language': 'en-US,en;q=0.8'}
type(html_soup)
```

Out[4]:  bs4.BeautifulSoup

Before extracting the 50 div containers, we need to figure out what distinguishes them from other div elements on that page. Often, the distinctive mark resides in the class attribute. If you inspect the HTML lines of the containers o interest, you'll notice that the class attribute has two values: lister-item and mode-advanced. This combination i unique to these div containers. We can see that's true by doing a quick search (Ctrl + F). We have 50 such containers, o we expect to see only 50 matches:

Now let's use the find_all() method to extract all the div containers that have a class attribute of lister-item mode-advanced:

In [5]:
```python
movie_containers = html_soup.find_all('li',class_='ipc-metadata-list-summary-item')
print(type(movie_containers))
print(len(movie_containers))
```
```
<class 'bs4.element.ResultSet'>
50
```

find_all() returned a ResultSet object which is a list containing all the 50 divs we are interested in.

Now we'll select only the first container, and extract, by turn, each item of interest:

-

The name of the movie

- The year of release.
- The IMDB rating.
- The Metascore.
- The number of votes..otes.

Extracting the data for a single movie

We can access the first container, which contains information about a single movie, by using list notation on movie_containers.

In [6]:
```python
first_movie = movie_containers[0]
first_movie
```

Out[6]: `<li class="ipc-metadata-list-summary-item"><div class="ipc-metadata-list-summary-item__c"><div class="ipc-metadata-list-summary-item__tc"><span aria-disabled="false" class="ipc-metadata-list-summary-item__t"></span><div class="sc-ab6fa25a-3 bVYfLY dli-parent"><div class="sc-ab6fa25a-2 gOsifL"><div class="sc-e5a25b0f-0 jQjDIb dli-poster-container"><div class="ipc-poster ipc-poster--base ipc-poster--dynamic-width ipc-sub-grid-item ipc-sub-grid-item--span-2" role="group"><div aria-label="add to watchlist" class="ipc-watchlist-ribbon ipc-focusable ipc-watchlist-ribbon--s ipc-watchlist-ribbon--base ipc-watchlist-ribbon--loading ipc-watchlist-ribbon--onImage ipc-poster__watchlist-ribbon" role="button" tabindex="0"><svg class="ipc-watchlist-ribbon__bg" height="34px" role="presentation" viewbox="0 0 24 34" width="24px" xmlns="http://www.w3.org/2000/svg"><polygon class="ipc-watchlist-ribbon__bg-ribbon" fill="#000000" points="24 0 0 0 32 12.2436611 26.2926049 24 31.7728343"></polygon><polygon class="ipc-watchlist-ribbon__bg-hover" points="24 0 0 0 32 12.2436611 26.2926049 24 31.7728343"></polygon><polygon class="ipc-watchlist-ribbon__bg-shadow" points="24 31.7728343 24 33.7728343 12.2436611 28.2926049 0 34 0 32 12.2436611 26.2926049"></polygon></svg><div class="ipc-watchlist-ribbon__icon" role="presentation"><svg class="ipc-loader ipc-loader--circle ipc-watchlist-ribbon__loader" data-testid="watchlist-ribbon-loader" height="48px" role="presentation" version="1.1" viewbox="0 0 48 48" width="48px" xmlns="http://www.w3.org/2000/svg"><g class="ipc-loader__container" fill="currentColor"><circle class="ipc-loader__circle ipc-loader__circle--one" cx="24" cy="9" r="4"></circle><circle class="ipc-loader__circle ipc-loader__circle--two" cx="35" cy="14" r="4"></circle><circle class="ipc-loader__circle ipc-loader__circle--three" cx="39" cy="24" r="4"></circle><circle class="ipc-loader__circle ipc-loader__circle--four" cx="35" cy="34" r="4"></circle><circle class="ipc-loader__circle ipc-loader__circle--five" cx="24" cy="39" r="4"></circle><circle class="ipc-loader__circle ipc-loader__circle--six" cx="13" cy="34" r="4"></circle><circle class="ipc-loader__circle ipc-loader__circle--seven" cx="9" cy="24" r="4"></circle><circle class="ipc-loader__circle ipc-loader__circle--eight" cx="13" cy="14" r="4"></circle></g></svg></div></div><div class="ipc-media ipc-media--poster-27x40 ipc-image-media-ratio--poster-27x40 ipc-media--base ipc-media--poster-m ipc-poster__poster-image ipc-media__img" style="width:100%"><img alt="Hugh Jackman in Logan (2017)" class="ipc-image" loading="lazy" sizes="50vw, (min-width: 480px) 34vw, (min-width: 600px) 26vw, (min-width: 1024px) 16vw, (min-width: 1280px) 16vw" src="https://m.media-amazon.com/images/M/MV5BYzc5MTU4N2EtYTkyMi00NjdhLTg3NWEtMTY4OTEyMzJhZTAzXkEyXkFqcGdeQXVyNjc1NTYyMjg@._V1_QL75_UX140_CR0,1,140,207_.jpg" srcset="https://m.media-amazon.com/images/M/MV5BYzc5MTU4N2EtYTkyMi00NjdhLTg3NWEtMTY4OTEyMzJhZTAzXkEyXkFqcGdeQXVyNjc1NTYyMjg@._V1_QL75_UX140_CR0,1,140,207_.jpg 140w, https://m.media-amazon.com/images/M/MV5BYzc5MTU4N2EtYTkyMi00NjdhLTg3NWEtMTY4OTEyMzJhZTAzXkEyXkFqcGdeQXVyNjc1NTYyMjg@._V1_QL75_UX210_CR0,2,210,311_.jpg 210w, https://m.media-amazon.com/images/M/MV5BYzc5MTU4N2EtYTkyMi00NjdhLTg3NWEtMTY4OTEyMzJhZTAzXkEyXkFqcGdeQXVyNjc1NTYyMjg@._V1_QL75_UX280_CR0,3,280,414_.jpg 280w" width="140"/></div><a aria-label="View title page for Logan" class="ipc-lockup-overlay ipc-focusable" href="/title/tt3315342/?ref_=sr_i_1"><div class="ipc-lockup-overlay__screen"></div></a></div></div><div class="sc-b0691f29-0 jbYPfh"><div class="ipc-title ipc-title--base ipc-title--title ipc-title-link-no-icon ipc-title--on-textPrimary sc-b0691f29-9 klOwFB dli-title"><a class="ipc-title-link-wrapper" href="/title/tt3315342/?ref_=sr_t_1" tabindex="0"><h3 class="ipc-title__text">1. Logan</h3></a></div><div class="sc-b0691f29-7 hrgukm dli-title-metadata"><span class="sc-b0691f29-8 ilsLEX dli-title-metadata-item">2017</span><span class="sc-b0691f29-8 ilsLEX dli-title-metadata-item">2h 17m</span><span class="sc-b0691f29-8 ilsLEX dli-title-metadata-item">R-16</span></div><span class="sc-b0691f29-1 grHDBY"><div class="sc-e2dbc1a3-0 ajrIH sc-b0691f29-2 bhhtyj dli-ratings-container" data-testid="ratingGroup--container"><span aria-label="IMDb rating: 8.1" class="ipc-rating-star ipc-rating-star--base ipc-rating-star--imdb ratingGroup--imdb-rating" data-testid="ratingGroup--imdb-rating"><svg class="ipc-icon ipc-icon--star-inline" fill="currentColor" height="24" role="presentation" viewbox="0 0 24 24" width="24" xmlns="http://www.w3.org/2000/svg"><path d="M12 20.1l5.82 3.682c1.066.675 2.37-.322 2.09-1.584l-1.543-6.926 5.146-4.667c.94-.85.435-2.465-.799-2.567l-6.773-.602L13.29.89a1.38 1.38 0 0 0-2.581 0l-2.65 6.53-6.774.602C.052 8.126-.453 9.74.486 10.59l5.147 4.666-1.542 6.926c-.28 1.262 1.023 2.26 2.09 1.585L12 20.099z"></path></svg>8.1<span class="ipc-rating-star--voteCount"> (<!-- -->827K<!-- -->)</span></span><button aria-label="Rate Logan" class="ipc-rate-button sc-e2dbc1a3-1 jboOQc ratingGroup--user-rating ipc-rate-button--unrated ipc-rate-button--base" data-testid="rate-button"><span class="ipc-rating-star ipc-rating-star--base ipc-rating-star--rate"><svg class="ipc-icon ipc-icon--star-border-inline" fill="currentColor" height="24" role="presentation" viewbox="0 0 24 24" width="24" xmlns="http://www.w3.org/2000/svg"><path d="M22.724 8.217l-6.786-.587-2.65-6.22c-.477-1.133-2.103-1.133-2.58 0l-2.65 6.234-6.772.573c-1.234.098-1.739 1.636-.8 2.446l5.146 4.446-1.542 6.598c-.28 1.202 1.023 2.153 2.09 1.511l5.818-3.495 5.819 3.509c1.065.643 2.37-.308 2.089-1.51l-1.542-6.612 5.145-4.446c.94-.81.45-2.348-.785-2.446zm-10.726 8.89l-5.272 3.174 1.402-5.983-4.655-4.026 6.141-.531 2.384-5.634 2.398 5.648 6.14.531-4.654 4.026 1.402 5.983-5.286-3.187z"></path></svg><span class="ipc-rating-star--rate">Rate</span></span></button></div><span class="sc-b0691f29-11 TmkKM"><span class="sc-b0901df4-0 bcQdDJ metacritic-score-box" style="background-color:#54A72A">77</span><span class="metacritic-score-label">Metascore</span></span></span></div><div class="sc-ab6fa25a-4 ggHbBR dli-post-element"><button aria-disabled="false" aria-label="See more information about Logan" class="ipc-icon-button dli-info-icon ipc-icon-button--base ipc-icon-button--onAccent2" role="button" tabindex="0" title="See more information about Logan"><svg class="ipc-icon ipc-icon--info" fill="currentColor" height="24" role="presentation" viewbox="0 0 24 24" width="24" xmlns="http://www.w3.org/2000/svg"><path d="M0 0h24v24H0V0z" fill="none"></path><path d="M11 7h2v2h-2zm0 4h2v6h-2zm1-9C6.48 2 2 6.48 2 12s4.48 10 10 10 10-4.48 10-10S17.52 2 12 2zm0 18c-4.41 0-8-3.59-8-8s3.59-8 8-8 8 3.59 8 8-3.59 8-8 8z"></path></svg></button></div></div><div class="sc-ab6fa25a-1 bBwFsP"><div class="ipc-html-content ipc-html-content--base sc-ab6fa25a-0 bhexuD dli-plot-container" role="presentation"><div class="ipc-html-content-inner-div">In a future where mutants are nearly extinct, an elderly and weary Logan leads a quiet life. But when Laura, a mutant child pursued by scientists, comes to him for help, he must get her to safety.</div></div></div></div></div></div></li>`

The name of the Movie

```
In [7]: first_movie.div
```

```
Out[7]:  <div class="ipc-metadata-list-summary-item__c"><div class="ipc-metadata-list-summary-item__tc"><span aria-d
         isabled="false" class="ipc-metadata-list-summary-item__t"></span><div class="sc-ab6fa25a-3 bVYfLY dli-paren
         t"><div class="sc-ab6fa25a-2 gOsifL"><div class="sc-e5a25b0f-0 jQjDIb dli-poster-container"><div class="ipc
         -poster ipc-poster--base ipc-poster--dynamic-width ipc-sub-grid-item ipc-sub-grid-item--span-2" role="grou
         p"><div aria-label="add to watchlist" class="ipc-watchlist-ribbon ipc-focusable ipc-watchlist-ribbon--s ipc
         -watchlist-ribbon--base ipc-watchlist-ribbon--loading ipc-watchlist-ribbon--onImage ipc-poster__watchlist-r
         ibbon" role="button" tabindex="0"><svg class="ipc-watchlist-ribbon__bg" height="34px" role="presentation" v
         iewbox="0 0 24 34" width="24px" xmlns="http://www.w3.org/2000/svg"><polygon class="ipc-watchlist-ribbon__bg
         -ribbon" fill="#000000" points="24 0 0 0 0 32 12.2436611 26.2926049 24 31.7728343"></polygon><polygon class
         ="ipc-watchlist-ribbon__bg-hover" points="24 0 0 0 0 32 12.2436611 26.2926049 24 31.7728343"></polygon><pol
         ygon class="ipc-watchlist-ribbon__bg-shadow" points="24 31.7728343 24 33.7728343 12.2436611 28.2926049 0 34
         0 32 12.2436611 26.2926049"></polygon></svg><div class="ipc-watchlist-ribbon__icon" role="presentation"><sv
         g class="ipc-loader ipc-loader--circle ipc-watchlist-ribbon__loader" data-testid="watchlist-ribbon-loader"
         height="48px" role="presentation" version="1.1" viewbox="0 0 48 48" width="48px" xmlns="http://www.w3.org/2
         000/svg"><g class="ipc-loader__container" fill="currentColor"><circle class="ipc-loader__circle ipc-loader_
         _circle--one" cx="24" cy="9" r="4"></circle><circle class="ipc-loader__circle ipc-loader__circle--two" cx
         ="35" cy="14" r="4"></circle><circle class="ipc-loader__circle ipc-loader__circle--three" cx="39" cy="24" r
         ="4"></circle><circle class="ipc-loader__circle ipc-loader__circle--four" cx="35" cy="34" r="4"></circle><c
         ircle class="ipc-loader__circle ipc-loader__circle--five" cx="24" cy="39" r="4"></circle><circle class="ipc
         -loader__circle ipc-loader__circle--six" cx="13" cy="34" r="4"></circle><circle class="ipc-loader__circle i
         pc-loader__circle--seven" cx="9" cy="24" r="4"></circle><circle class="ipc-loader__circle ipc-loader__circl
         e--eight" cx="13" cy="14" r="4"></circle></g></svg></div></div><div class="ipc-media ipc-media--poster-27x4
         0 ipc-image-media-ratio--poster-27x40 ipc-media--base ipc-media--poster-m ipc-poster__poster-image ipc-medi
         a__img" style="width:100%"><img alt="Hugh Jackman in Logan (2017)" class="ipc-image" loading="lazy" sizes
         ="50vw, (min-width: 480px) 34vw, (min-width: 600px) 26vw, (min-width: 1024px) 16vw, (min-width: 1280px) 16v
         w" src="https://m.media-amazon.com/images/M/MV5BYzc5MTU4N2EtYTkyMi00NjdhLTg3NWEtMTY4OTEyMzJhZTAzXkEyXkFqcGd
         eQXVyNjc1NTYyMjg@._V1_QL75_UX140_CR0,1,140,207_.jpg" srcset="https://m.media-amazon.com/images/M/MV5BYzc5MT
         U4N2EtYTkyMi00NjdhLTg3NWEtMTY4OTEyMzJhZTAzXkEyXkFqcGdeQXVyNjc1NTYyMjg@._V1_QL75_UX140_CR0,1,140,207_.jpg 14
         0w, https://m.media-amazon.com/images/M/MV5BYzc5MTU4N2EtYTkyMi00NjdhLTg3NWEtMTY4OTEyMzJhZTAzXkEyXkFqcGdeQXV
         yNjc1NTYyMjg@._V1_QL75_UX210_CR0,2,210,311_.jpg 210w, https://m.media-amazon.com/images/M/MV5BYzc5MTU4N2EtY
         TkyMi00NjdhLTg3NWEtMTY4OTEyMzJhZTAzXkEyXkFqcGdeQXVyNjc1NTYyMjg@._V1_QL75_UX280_CR0,3,280,414_.jpg 280w" wid
         th="140"/></div><a aria-label="View title page for Logan" class="ipc-lockup-overlay ipc-focusable" href="/t
         itle/tt3315342/?ref_=sr_i_1"><div class="ipc-lockup-overlay__screen"></div></a></div></div><div class="sc-b
         0691f29-0 jbYPfh"><div class="ipc-title ipc-title--base ipc-title--title ipc-title-link-no-icon ipc-title--
         on-textPrimary sc-b0691f29-9 klOwFB dli-title"><a class="ipc-title-link-wrapper" href="/title/tt3315342/?re
         f_=sr_t_1" tabindex="0"><h3 class="ipc-title__text">1. Logan</h3></a></div><div class="sc-b0691f29-7 hrgukm
         dli-title-metadata"><span class="sc-b0691f29-8 ilsLEX dli-title-metadata-item">2017</span><span class="sc-b
         0691f29-8 ilsLEX dli-title-metadata-item">2h 17m</span><span class="sc-b0691f29-8 ilsLEX dli-title-metadata
         -item">R-16</span></div><span class="sc-b0691f29-1 grHDBY"><div class="sc-e2dbc1a3-0 ajrIH sc-b0691f29-2 bh
         htyj dli-ratings-container" data-testid="ratingGroup--container"><span aria-label="IMDb rating: 8.1" class
         ="ipc-rating-star ipc-rating-star--base ipc-rating-star--imdb ratingGroup--imdb-rating" data-testid="rating
         Group--imdb-rating"><svg class="ipc-icon ipc-icon--star-inline" fill="currentColor" height="24" role="prese
         ntation" viewbox="0 0 24 24" width="24" xmlns="http://www.w3.org/2000/svg"><path d="M12 20.1l5.82 3.682c1.0
         66.675 2.37-.322 2.09-1.584l-1.543-6.926 5.146-4.667c.94-.85.435-2.465-.799-2.567l-6.773-.602L13.29.89a1.38
         1.38 0 0 0-2.581 0l-2.65 6.53-6.774.602C.052 8.126-.453 9.74.486 10.59l5.147 4.666-1.542 6.926c-.28 1.262
         1.023 2.26 2.09 1.585L12 20.099z"></path></svg>8.1<span class="ipc-rating-star--voteCount"> (<!-- -->827K<!
         -- -->)</span></span><button aria-label="Rate Logan" class="ipc-rate-button sc-e2dbc1a3-1 jboOQc ratingGrou
         p--user-rating ipc-rate-button--unrated ipc-rate-button--base" data-testid="rate-button"><span class="ipc-r
         ating-star ipc-rating-star--base ipc-rating-star--rate"><svg class="ipc-icon ipc-icon--star-border-inline"
         fill="currentColor" height="24" role="presentation" viewbox="0 0 24 24" width="24" xmlns="http://www.w3.or
         g/2000/svg"><path d="M22.724 8.217l-6.786-.587-2.65-6.22c-.477-1.133-2.103-1.133-2.58 0l-2.65 6.234-6.772.5
         73c-1.234.098-1.739 1.636-.8 2.446l5.146 4.446-1.542 6.598c-.28 1.202 1.023 2.153 2.09 1.5l5.818-3.495 5.8
         19 3.509c1.065.643 2.37-.308 2.089-1.5l-1.542-6.612 5.145-4.446c.94-.81.45-2.348-.785-2.446zm-10.726 8.89l
         -5.272 3.174 1.402-5.983-4.655-4.026 6.141-.531 2.384-5.634 2.398 5.648 6.14.531-4.654 4.026 1.402 5.983-5.
         286-3.187z"></path></svg><span class="ipc-rating-star--rate">Rate</span></span></button></div><span class
         ="sc-b0691f29-11 TmkKM"><span class="sc-b0901df4-0 bcQdDJ metacritic-score-box" style="background-color:#54
         A72A">77</span><span class="metacritic-score-label">Metascore</span></span></span></div><div class="sc-ab6f
         a25a-4 ggHbBR dli-post-element"><button aria-disabled="false" aria-label="See more information about Logan"
         class="ipc-icon-button dli-info-icon ipc-icon-button--base ipc-icon-button--onAccent2" role="button" tabind
         ex="0" title="See more information about Logan"><svg class="ipc-icon ipc-icon--info" fill="currentColor" he
         ight="24" role="presentation" viewbox="0 0 24 24" width="24" xmlns="http://www.w3.org/2000/svg"><path d="M0
         0h24v24H0V0z" fill="none"></path><path d="M11 7h2v2h-2zm0 4h2v6h-2zm1-9C6.48 2 2 6.48 2 12s4.48 10 10 10 10
         -4.48 10-10S17.52 2 12 2zm0 18c-4.41 0-8-3.59-8-8s3.59-8 8-8 8 3.59 8 8-3.59 8-8 8z"></path></svg></button>
         </div></div><div class="sc-ab6fa25a-1 bBwFsP"><div class="ipc-html-content ipc-html-content--base sc-ab6fa2
         5a-0 bhexuD dli-plot-container" role="presentation"><div class="ipc-html-content-inner-div">In a future whe
         re mutants are nearly extinct, an elderly and weary Logan leads a quiet life. But when Laura, a mutant chil
         d pursued by scientists, comes to him for help, he must get her to safety.</div></div></div></div></div></d
         iv>
```

```
In [8]:  first_movie.a
```

```
Out[8]:  <a aria-label="View title page for Logan" class="ipc-lockup-overlay ipc-focusable" href="/title/tt3315342/?
         ref_=sr_i_1"><div class="ipc-lockup-overlay__screen"></div></a>
```

```
In [9]:  first_movie.h3
```

```
Out[9]:  <h3 class="ipc-title__text">1. Logan</h3>
```

```
In [27]:  first_name = first_movie.h3.text[3:]
          first_name
```

```
Out[27]:  'Logan'
```

The year of movie's release

```
In [52]:  first_movie.find('span',class_ ='sc-b0691f29-8 ilsLEX dli-title-metadata-item')
```

```
Out[52]:  <span class="sc-b0691f29-8 ilsLEX dli-title-metadata-item">2017</span>
```

```
In [49]:  first_year = first_movie.find('span', class_ ='sc-b0691f29-8 ilsLEX dli-title-metadata-item')
          first_year.text
```

```
Out[49]:  '2017'
```

The IMDB rating

```
In [57]:  first_imdb = first_movie.find('span', class_ ="ipc-rating-star ipc-rating-star--base ipc-rating-star--imdb r
          first_imdb
```

```
Out[57]:  <span aria-label="IMDb rating: 8.1" class="ipc-rating-star ipc-rating-star--base ipc-rating-star--imdb rati
          ngGroup--imdb-rating" data-testid="ratingGroup--imdb-rating"><svg class="ipc-icon ipc-icon--star-inline" fi
          ll="currentColor" height="24" role="presentation" viewbox="0 0 24 24" width="24" xmlns="http://www.w3.org/2
          000/svg"><path d="M12 20.1l5.82 3.682c1.066.675 2.37-.322 2.09-1.584l-1.543-6.926 5.146-4.667c.94-.85.435-
          2.465-.799-2.567l-6.773-.602L13.29.89a1.38 1.38 0 0 0-2.581 0l-2.65 6.53-6.774.602C.052 8.126-.453 9.74.486
          10.59l5.147 4.666-1.542 6.926c-.28 1.262 1.023 2.26 2.09 1.585L12 20.099z"></path></svg>8.1<span class="ipc
          -rating-star--voteCount"> (<!-- -->827K<!-- -->)</span></span>
```

```
In [58]:  first_imdb.text
```

```
Out[58]:  '8.1\xa0(827K)'
```

```
In [60]:  import re
          first_imdb.find(string=re.compile("."))
```

```
Out[60]:  '8.1'
```

The Metascore

```
In [65]:  first_mscore = first_movie.find('span',  class_ ='sc-b0901df4-0 bcQdDJ metacritic-score-box')
          first_mscore.text
```

```
Out[65]:  '77'
```

The number of votes

```
In [74]:  first_votes = first_movie.find('span', class_='ipc-rating-star--voteCount').find(string=re.compile("K"))
          first_votes
```

```
Out[74]:  '827K'
```

```
In [88]:  #List to store the scrapped data in
          names=[]
          years= []
          imdb_ratings = []
          metascores = []
```

```
votes = []

#Extract data from individual movie container

for container in movie_containers:
#if the movie has metascore, then extract:
    if container.find('span', class_ = 'sc-b0691f29-11 TmkKM') is not None:
    #The name
        name = container.h3.text[3:]
        names.append(name)
    #The year
        year = container.find('span', class_ ='sc-b0691f29-8 ilsLEX dli-title-metadata-item').text
        years.append(year)
    #The IMDB rating
        imdb = container.find('span', class_ ="ipc-rating-star ipc-rating-star--base ipc-rating-star--imdb r
        imdb_ratings.append(float(imdb))
    # The Metascore
        m_score = container.find('span',  class_ ='sc-b0901df4-0 bcQdDJ metacritic-score-box').text
        metascores.append(int(m_score))
    #The number of votes
        vote =container.find('span', class_='ipc-rating-star--voteCount').find(string=re.compile("K"))
        votes.append(vote)
```

In [89]:
```
#To Test the following result of the above code:
for i in range(10):
    print(names[i])
    print(years[i])
    print(imdb_ratings[i])
    print(metascores[i])
    print(votes[i])
    print("\n")
```

Logan
2017
8.1
77
827K


Thor: Ragnarok
2017
7.9
74
813K


Guardians of the Galaxy Vol. 2
2017
7.6
67
756K


Dunkirk
2017
7.8
94
736K


Spider-Man: Homecoming
2017
7.4
73
716K


Wonder Woman
2017
7.3
76
698K


Get Out
2017
7.8
85
691K


Star Wars: Episode VIII - The Last Jedi
2017
6.9
84
670K


Blade Runner 2049
2017
8.0
81
658K


 Baby Driver
2017
7.5
86
605K

```python
In [91]:  import pandas as pd
          test_df = pd.DataFrame({'movie':names,
                                  'year':years,
                                  'imdb':imdb_ratings,
                                  'metascore':metascores,
                                  'votes':votes
                                  })
          print(test_df.info())
          test_df
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 41 entries, 0 to 40
Data columns (total 5 columns):
 #   Column     Non-Null Count  Dtype
---  ------     --------------  -----
 0   movie      41 non-null     object
 1   year       41 non-null     object
 2   imdb       41 non-null     float64
 3   metascore  41 non-null     int64
 4   votes      41 non-null     object
dtypes: float64(1), int64(1), object(3)
memory usage: 1.7+ KB
None
```

|    | movie | year | imdb | metascore | votes |
|----|-------|------|------|-----------|-------|
| 0 | Logan | 2017 | 8.1 | 77 | 827K |
| 1 | Thor: Ragnarok | 2017 | 7.9 | 74 | 813K |
| 2 | Guardians of the Galaxy Vol. 2 | 2017 | 7.6 | 67 | 756K |
| 3 | Dunkirk | 2017 | 7.8 | 94 | 736K |
| 4 | Spider-Man: Homecoming | 2017 | 7.4 | 73 | 716K |
| 5 | Wonder Woman | 2017 | 7.3 | 76 | 698K |
| 6 | Get Out | 2017 | 7.8 | 85 | 691K |
| 7 | Star Wars: Episode VIII - The Last Jedi | 2017 | 6.9 | 84 | 670K |
| 8 | Blade Runner 2049 | 2017 | 8.0 | 81 | 658K |
| 9 | Baby Driver | 2017 | 7.5 | 86 | 605K |
| 10 | It | 2017 | 7.3 | 69 | 603K |
| 11 | Coco | 2017 | 8.4 | 81 | 586K |
| 12 | Three Billboards Outside Ebbing, Missouri | 2017 | 8.1 | 88 | 553K |
| 13 | John Wick: Chapter 2 | 2017 | 7.4 | 75 | 509K |
| 14 | Justice League | 2017 | 6.1 | 45 | 477K |
| 15 | The Shape of Water | 2017 | 7.3 | 87 | 446K |
| 16 | Jumanji: Welcome to the Jungle | 2017 | 6.9 | 58 | 435K |
| 17 | Kingsman: The Golden Circle | 2017 | 6.7 | 44 | 361K |
| 18 | Kong: Skull Island | 2017 | 6.7 | 62 | 345K |
| 19 | Pirates of the Caribbean: Salazar's Revenge | 2017 | 6.5 | 39 | 344K |
| 20 | Beauty and the Beast | 2017 | 7.1 | 65 | 333K |
| 21 | Lady Bird | 2017 | 7.4 | 93 | 326K |
| 22 | Call Me by Your Name | 2017 | 7.8 | 94 | 313K |
| 23 | The Greatest Showman | 2017 | 7.5 | 48 | 310K |
| 24 | Alien: Covenant | 2017 | 6.4 | 65 | 302K |
| 25 | Murder on the Orient Express | 2017 | 6.5 | 52 | 295K |
| 26 | War for the Planet of the Apes | 2017 | 7.4 | 82 | 280K |
| 27 | Wind River | 2017 | 7.7 | 73 | 279K |
| 28 | Fast & Furious 8 | 2017 | 6.6 | 56 | 253K |
| 29 | Life | 2017 | 6.6 | 54 | 252K |
| 30 | Mother! | 2017 | 6.6 | 76 | 249K |
| 31 | The Hitman's Bodyguard | 2017 | 6.9 | 47 | 246K |
| 32 | I, Tonya | 2017 | 7.5 | 77 | 242K |
| 33 | King Arthur: Legend of the Sword | 2017 | 6.7 | 41 | 232K |
| 34 | Ghost in the Shell | 2017 | 6.3 | 52 | 227K |
| 35 | Darkest Hour | 2017 | 7.4 | 75 | 220K |
| 36 | American Made | 2017 | 7.1 | 65 | 207K |
| 37 | Atomic Blonde | 2017 | 6.7 | 63 | 206K |

| | movie | year | imdb | metascore | votes |
|---|---|---|---|---|---|
| **38** | The Mummy | 2017 | 5.4 | 34 | 206K |
| **39** | Baywatch | 2017 | 5.5 | 37 | 201K |
| **40** | Bright | 2017 | 6.3 | 29 | 201K |

The Script for Multiple pages

```python
from time import time
from time import sleep
from random import randint
from IPython.display import clear_output

years_url = [ '2004','2005','2006','2007','2008',
              '2009','2010','2011','2012','2013',
              '2014','2015','2016','2017', '2018',
              '2019', '2020','2021','2022','2023']

names = []
years = []
imdb_ratings = []
metascores = []
votes = []

start_time = time()
requests = 0


agent = {"User-Agent": 'Mozilla/5.0 (Linux; Android 6.0; Nexus 5 Build/MRA58N) AppleWebKit/537.36 (KHTML, li
# For every year in the interval 2004-2023
for year_url in years_url:
    #make a get request
    url = f"https://www.imdb.com/search/title/?release_date={year_url}-01-01,{year_url}-12-31&sort=num_votes
    print(url)
    response = get(url, headers=agent)
    print(response.url)
    sleep(randint(8,15))

    requests += 1
    elapsed_time = time() - start_time
    print('Request:{}; Frequency: {} requests/s'.format(requests, requests/elapsed_time))
    clear_output(wait = True)

    if response.status_code != 200:
        print('Request: {}; Status code: {}'.format(requests, response.status_code))

    if requests > 72:
        print('Number of requests was greater than expected.')
        break

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

    mv_containers = page_html.find_all('li', class_ = 'ipc-metadata-list-summary-item')

    for container in mv_containers:
        if container.find('span', class_="sc-b0691f29-11 TmkKM") is not None:
            name = container.h3.text[3:]
            names.append(name)

            year = container.find('span', class_ = 'sc-b0691f29-8 ilsLEX dli-title-metadata-item').text
            years.append(year)

            imdb = float(container.find('span', class_='ipc-rating-star ipc-rating-star--base ipc-rating-sta
            imdb_ratings.append(imdb)

            m_score = int(container.find('span', class_ = 'sc-b0901df4-0 bcQdDJ metacritic-score-box').text)
            metascores.append(m_score)
```

```
            vote = container.find('span', class_="ipc-rating-star--voteCount").find(string= re.compile("[KM]
            votes.append(vote)

    del response
```

https://www.imdb.com/search/title/?release_date=2023-01-01,2023-12-31&sort=num_votes,desc
https://www.imdb.com/search/title/?release_date=2023-01-01,2023-12-31&sort=num_votes,desc
Request:20; Frequency: 0.06436068578389736 requests/s

In [118... 
```python
movie_ratings = pd.DataFrame({'movie':names,
                              'year':years,
                              'imdb':imdb_ratings,
                              'metascore':metascores,
                              'votes':votes
                             })
print(movie_ratings.info())
movie_ratings
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 845 entries, 0 to 844
Data columns (total 5 columns):
 #   Column     Non-Null Count  Dtype
---  ------     --------------  -----
 0   movie      845 non-null    object
 1   year       845 non-null    object
 2   imdb       845 non-null    float64
 3   metascore  845 non-null    int64
 4   votes      845 non-null    object
dtypes: float64(1), int64(1), object(3)
memory usage: 33.1+ KB
None
```

Out[118...

|     | movie | year | imdb | metascore | votes |
|-----|-------|------|------|-----------|-------|
| 0   | Eternal Sunshine of the Spotless Mind | 2004 | 8.3 | 89 | 1.1M |
| 1   | The Incredibles | 2004 | 8.0 | 90 | 804K |
| 2   | Kill Bill: Vol. 2 | 2004 | 8.0 | 83 | 801K |
| 3   | Million Dollar Baby | 2004 | 8.1 | 86 | 719K |
| 4   | Spider-Man 2 | 2004 | 7.5 | 83 | 705K |
| ... | ... | ... | ... | ... | ... |
| 840 | Sound of Freedom | 2023 | 7.7 | 36 | 111K |
| 841 | Asteroid City | 2023 | 6.5 | 75 | 110K |
| 842 | A Haunting in Venice | 2023 | 6.5 | 63 | 109K |
| 843 | The Hunger Games: The Ballad of Songbirds & S... | 2023 | 6.8 | 54 | 108K |
| 844 | The Equalizer 3 | 2023 | 6.8 | 58 | 107K |

845 rows × 5 columns

In [119... 
```python
movie_ratings.head(10)
```

| | movie | year | imdb | metascore | votes |
|---|---|---|---|---|---|
| 0 | Eternal Sunshine of the Spotless Mind | 2004 | 8.3 | 89 | 1.1M |
| 1 | The Incredibles | 2004 | 8.0 | 90 | 804K |
| 2 | Kill Bill: Vol. 2 | 2004 | 8.0 | 83 | 801K |
| 3 | Million Dollar Baby | 2004 | 8.1 | 86 | 719K |
| 4 | Spider-Man 2 | 2004 | 7.5 | 83 | 705K |
| 5 | Harry Potter and the Prisoner of Azkaban | 2004 | 7.9 | 82 | 688K |
| 6 | The Notebook | 2004 | 7.8 | 53 | 617K |
| 7 | Shaun of the Dead | 2004 | 7.9 | 76 | 592K |
| 8 | I, Robot | 2004 | 7.1 | 59 | 573K |
| 9 | Troy | 2004 | 7.3 | 56 | 568K |

In [120...
```
movie_ratings.tail(10)
```

| | movie | year | imdb | metascore | votes |
|---|---|---|---|---|---|
| 835 | La sociedad de la nieve | 2023 | 7.8 | 72 | 122K |
| 836 | The Marvels | 2023 | 5.6 | 50 | 119K |
| 837 | Scream VI | 2023 | 6.5 | 61 | 118K |
| 838 | Fast X | 2023 | 5.8 | 56 | 117K |
| 839 | Knock at the Cabin | 2023 | 6.1 | 63 | 114K |
| 840 | Sound of Freedom | 2023 | 7.7 | 36 | 111K |
| 841 | Asteroid City | 2023 | 6.5 | 75 | 110K |
| 842 | A Haunting in Venice | 2023 | 6.5 | 63 | 109K |
| 843 | The Hunger Games: The Ballad of Songbirds & S... | 2023 | 6.8 | 54 | 108K |
| 844 | The Equalizer 3 | 2023 | 6.8 | 58 | 107K |

In [123...
```
movie_ratings.to_csv('C:/Users/Lenovo/Downloads/HOA7.2 Webscraping using BeautifulSoap and Request/movie_rat
```

## Data Preparation

- Collected data may not be compatible or formatted correctly
- 

Data must be prepared before it can be added to a data se

- Extract, Transform and Load (ETL)

  > process for collecting data from a variety of sources, transforming the data, and then loading the data into
  > a database

)

## Data preprocessing

Data Processing is a process of cleaning the raw data i.e. the data is collected in the real world and is converted to a clean data set. In other words, whenever the data is gathered from different sources it is collected in a raw format and this data

isn't feasible for the analysis. Therefore, certain steps are executed to convert the data into a small clean data set, this part of the process is called as data preprocessing.

Most of the real-world data is messy, some of these types of data are:

1. **Missing data**: Missing data can be found when it is not continuously created or due to technical

issues in the application (IOT system)

2. **Noisy Data**: This type of data is also called outliners, this can occur due to human errors (human manually gathering the data) o

some technical problem of the device at the time of collection of dat

3. **Inconsistent Data**: : This type of data might be collected due to human errors (mistakes with the

name or values) or duplication of data.

These are some of the basic pre processing techniques that can be used to convert raw data.

1. **Conversion of data**: As we know that Machine Learning models can only

handle numeric features, hence categorical and ordinal data must be somehow converted into numeric features

2. **Ignoring the missing values**s: Whenever we encounte

missing data in the data set then we can remove the row or column of data depending on our need. This method is known to be efficient but it shouldn't be performed if the e are a lot of missing values in the datas

3. **Filling the missing values**: es: Whenever we encounter missing data in the data set then we can fill the missing data manually, m st

commonly the mean, median or highest frequency value is u

4. **Machine learning:** If we have some missing data then we can predict what data shall be present at the empty position by using the existing data.

5. **Outliers

detection** :Some error datat might be present in our data set that deviates drastically from other observations in a data set. [Example: human weight = 80 Kg; due to mistyping of extra] 0sed.

## Example of Data Preparation of movie_rating.csv

```
In [125... movie_ratings['year'].unique()
```

```
Out[125... array(['2004', '2005', '2006', '2007', '2008', '2009', '2010', '2011',
               '2012', '2013', '2014', '2015', '2016', '2017', '2018', '2019',
               '2020', '2021', '2022', '2023'], dtype=object)
```

```
In [126... movie_ratings.dtypes
```

```
Out[126... movie          object
         year           object
         imdb          float64
         metascore       int64
         votes          object
         dtype: object
```

```
In [127... movie_ratings.head(10)
```

| | movie | year | imdb | metascore | votes |
|---|---|---|---|---|---|
| 0 | Eternal Sunshine of the Spotless Mind | 2004 | 8.3 | 89 | 1.1M |
| 1 | The Incredibles | 2004 | 8.0 | 90 | 804K |
| 2 | Kill Bill: Vol. 2 | 2004 | 8.0 | 83 | 801K |
| 3 | Million Dollar Baby | 2004 | 8.1 | 86 | 719K |
| 4 | Spider-Man 2 | 2004 | 7.5 | 83 | 705K |
| 5 | Harry Potter and the Prisoner of Azkaban | 2004 | 7.9 | 82 | 688K |
| 6 | The Notebook | 2004 | 7.8 | 53 | 617K |
| 7 | Shaun of the Dead | 2004 | 7.9 | 76 | 592K |
| 8 | I, Robot | 2004 | 7.1 | 59 | 573K |
| 9 | Troy | 2004 | 7.3 | 56 | 568K |

```
movie_ratings.tail(10)
```

| | movie | year | imdb | metascore | votes |
|---|---|---|---|---|---|
| 835 | La sociedad de la nieve | 2023 | 7.8 | 72 | 122K |
| 836 | The Marvels | 2023 | 5.6 | 50 | 119K |
| 837 | Scream VI | 2023 | 6.5 | 61 | 118K |
| 838 | Fast X | 2023 | 5.8 | 56 | 117K |
| 839 | Knock at the Cabin | 2023 | 6.1 | 63 | 114K |
| 840 | Sound of Freedom | 2023 | 7.7 | 36 | 111K |
| 841 | Asteroid City | 2023 | 6.5 | 75 | 110K |
| 842 | A Haunting in Venice | 2023 | 6.5 | 63 | 109K |
| 843 | The Hunger Games: The Ballad of Songbirds & S... | 2023 | 6.8 | 54 | 108K |
| 844 | The Equalizer 3 | 2023 | 6.8 | 58 | 107K |

```
movie_ratings
```

| | movie | year | imdb | metascore | votes |
|---|---|---|---|---|---|
| 0 | Eternal Sunshine of the Spotless Mind | 2004 | 8.3 | 89 | 1.1M |
| 1 | The Incredibles | 2004 | 8.0 | 90 | 804K |
| 2 | Kill Bill: Vol. 2 | 2004 | 8.0 | 83 | 801K |
| 3 | Million Dollar Baby | 2004 | 8.1 | 86 | 719K |
| 4 | Spider-Man 2 | 2004 | 7.5 | 83 | 705K |
| ... | ... | ... | ... | ... | ... |
| 840 | Sound of Freedom | 2023 | 7.7 | 36 | 111K |
| 841 | Asteroid City | 2023 | 6.5 | 75 | 110K |
| 842 | A Haunting in Venice | 2023 | 6.5 | 63 | 109K |
| 843 | The Hunger Games: The Ballad of Songbirds & S... | 2023 | 6.8 | 54 | 108K |
| 844 | The Equalizer 3 | 2023 | 6.8 | 58 | 107K |

845 rows × 5 columns