

Scraping Movie Information

USING SCRAPY, BEAUTIFUL SOUP AND SELENIUM

Joblo.com

The first step is to create a csv file to write to. We need to specify the name of the csv file we would like to save, set writing mode to “w”, which is going to overwrite the old file with the same name. If we the writing mode to “a”, it’s going to append the new rows to the existing file. Once we create the file and open it for writing, we write the first row as the column names. Then, we send html request to the url we would like to parse and use BeautifulSoup to extract the information we would like to store in the csv file.

```
1. import requests
2. from bs4 import BeautifulSoup
3. import csv
4.
5. # Create a file to write to, add headers row
6. f = csv.writer(open('joblo_poster.csv', 'w', newline=''))
7. f.writerow(['Title', 'Poster'])
8.
9. page = requests.get('https://www.joblo.com/movie-posters/archives/ALL/')
10. # Create a BeautifulSoup object
11. soup = BeautifulSoup(page.text, 'html.parser')
12. last_links = soup.find_all("li", attrs={"class":"vertical"})
13. number_of_poster = len(last_links)
14.
15. for i in range(0,number_of_poster):
16.     Title= last_links[i].find("div",class_="bottom-poster").a.string
17.     Poster = last_links[i].div.a.img.get("src")
18.     Poster = 'https://www.joblo.com'+Poster
19.     print(i)
20.     f.writerow([Title,Poster])
21.
22. del f
```

The table below shows the result of the above code. The column “Poster” is the link to the image that we will download in the next step.

Title	Poster
The Butterfly Effect	https://www.joblo.com/assets/images/oldsite/posters/images/full/2004-butterfly_effect-1_thumb.jpg
The Butterfly Effect	https://www.joblo.com/assets/images/oldsite/posters/images/full/2004-butterfly_effect-2_thumb.jpg
Agent Cody Banks 2	https://www.joblo.com/assets/images/oldsite/posters/images/full/agent-cody-banks-two-poster_thumb.jpg
Blade: Trinity	https://www.joblo.com/assets/images/oldsite/posters/images/full/2004-blade_trinity-1_thumb.jpg
Blade: Trinity	https://www.joblo.com/assets/images/oldsite/posters/images/full/2004-blade_trinity-2_thumb.jpg
Blade: Trinity	https://www.joblo.com/assets/images/oldsite/posters/images/full/2004-blade_trinity-3_thumb.jpg
Blade: Trinity	https://www.joblo.com/assets/images/oldsite/posters/images/full/2004-blade_trinity-4_thumb.jpg
Blade: Trinity	https://www.joblo.com/assets/images/oldsite/posters/images/full/2004-blade_trinity-5_thumb.jpg
Barbershop 2: Back in	https://www.joblo.com/assets/images/oldsite/posters/images/full/barbershop2-poster_thumb.jpg

The next step is to download all the image from the links that we scraped to csv file earlier. First, we create the name of directory that we are going to store the images. After that, we create a list of all image links by using “Pandas” library.

```

1. #####Download from joblo csv file#####
2. #Define the name of directory
3. f_name = "/joblo_poster"
4.
5. # Create the directory name where the images will be saved
6. base_dir = os.getcwd()
7. dir_name = (f_name.split('/')[-1]).split('.')[0]
8. dir_name
9. dir_path = os.path.join(base_dir, dir_name)
10.
11. #Create the directory if already not there
12. if not os.path.exists(dir_path):
13.     os.mkdir(dir_path)
14.
15. # Read the csv with links to all the image pages
16. os.getcwd()
17. df = pd.read_csv("CSV_File\joblo_poster.csv")
18. df.columns
19. links=df.Poster

```

Once we have the list of image links ready, we define a function to download the image save them in the format “Pic_number_MovieName.jpg”. In order to avoid getting blocked by the site you are scraping or downloading, we change user agents name in the header.

```

1. # Function to take an image url and save the image in the given directory
2. def download_image(url,image_number):
3.     print("[INFO] downloading {}".format(url))
4.     name = str(url.split('/')[-1])
5.     name = f"Pic_{image_number}"+name #File name that will be saved
6.     opener=urllib.request.build_opener()#Add header and chage user-
       agent name so we don't get forbidden by the website
7.     opener.addheaders=[('User-
       Agent','Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chro
       me/36.0.1941.0 Safari/537.36')]
8.     urllib.request.install_opener(opener)
9.     urllib.request.urlretrieve(url,os.path.join(dir_path, name))
10.
11. #Download ALL image in URL links list
12. j=1
13. for i in links:
14.     print (j)
15.     download_image(i,j)
16.     j+=1

```

Impawards.com

The information we want from this website are also images, so the process is basically the same as for joblo.com. The difference is just the scraping part where we extract images link which store in different HTML tag format. The code below shows how to get the image links.

```

1. import requests
2. from bs4 import BeautifulSoup
3. import csv

```

```

4. #from time import sleep
5.
6. # Create a file to write to, add headers row
7. f = csv.writer(open('impawards.csv', 'a', newline='')) # 'a' to append not overwrite the
   original file ,, newline='' to write row without blank row between each row
8. f.writerow(['Designer', 'Number_of_poster', "Links", "Poster_name", "Year"])
9.
10.
11. page = requests.get('http://www.impawards.com/designers/index.html')
12. # Create a BeautifulSoup object
13. soup = BeautifulSoup(page.text, 'html.parser')
14. last_links = soup.find_all("div", attrs={"class": "col-md-4"}, limit=3)
15. number = len(last_links)
16. number
17.
18. for i in range(0, number):
19.
20.     Designerssss= last_links[i].ul.find_all("li")
21.     for j in range(0, len(Designerssss)):
22.         #Designer name
23.         Designer = Designerssss[j].find("b").string
24.         #Poster
25.         link_to_poster = Designerssss[j].a.get("href")
26.         link_to_poster = "http://www.impawards.com/designers/"+link_to_poster
27.         newpage = requests.get(link_to_poster) #Go into the link
28.         newsoup = BeautifulSoup(newpage.text, 'html.parser')
29.         new_last_links= newsoup.find("center").find_next_sibling("center")
30.         Posterssss=new_last_links.find_all("a")
31.         #Number of poster of each designer
32.         Designerssss[j].b.decompose()
33.         Number_of_poster = Designerssss[j].get_text()
34.         #Poster (continue)
35.         for k in range(0, len(Posterssss)):
36.
37.             #Get poster link URL
38.             Poster=Posterssss[k].get("href")
39.             if Poster[0:5]=='/intl':
40.                 inter=Poster.split("/")
41.                 MyPoster="http://www.impawards.com"+"/"+inter[1]+"/"+inter[2]+"/"+inter[3]+"/posters/"+inter[4][:4]+"jpg"
42.                 Poster_name = inter[4][:5]
43.                 Year = inter[3]
44.
45.             elif Poster[0:3]=='/tv' :
46.                 inter=Poster.split("/")
47.                 MyPoster="http://www.impawards.com"+"/"+inter[1]+"/posters/"+inter[2][:4]+"jpg"
48.                 Poster_name = Poster[4:-5]
49.                 Year = "na"
50.             else:
51.                 MyPoster="http://www.impawards.com"+Poster[0:5]+"/posters"+Poster[5:-4]+"jpg"
52.                 Poster_name = Poster[6:-5]
53.                 Year = Poster[1:5]
54.
55.             print(Designer)
56.             print(MyPoster)
57.             f.writerow([Designer, Number_of_poster, MyPoster, Poster_name, Year])
58.
59. del f

```

Designer	Number	Links	Poster_name	Year	pic_no	tconst	genres	originalTitle
11:24 Design Advertising	-52	http://www.impawa	maynard	2017	1	tt4375446	Documentary	maynard
11:24 Design Advertising	-52	http://www.impawa	true to the gar	2017	2	tt5116504	Drama	true to the game
11:24 Design Advertising	-52	http://www.impawa	sister code	2015	3	tt3912040	Comedy,Drama,F	sister code
11:24 Design Advertising	-52	http://www.impawa	falcon rising	2014	4	tt2295722	Action,Adventur	falcon rising
11:24 Design Advertising	-52	http://www.impawa	oldboy	2013	5	tt1321511	Action,Drama,M	oldboy

The next step is to download the image from the links we scraped earlier in the CSV file. The code below shows how to achieve this task.

```

1. f_name = "Impawards_updated"
2.
3. # Create the directory name where the images will be saved
4. base_dir = os.getcwd()
5. dir_name = (f_name.split('/')[-1]).split('.')[0]
6. dir_name
7. dir_path = os.path.join(base_dir, dir_name)
8.
9. #Create the directory if already not there
10. if not os.path.exists(dir_path):
11.     os.mkdir(dir_path)
12.
13. # Read the csv with links to all the image pages
14. os.getcwd()
15. df = pd.read_csv("impawards_updated.csv")
16. df.columns
17. links=df.Links
18.
19. # Function to take an image url and save the image in the given directory
20. def download_image(url,image_number):
21.     print("[INFO] downloading {}".format(url))
22.     name = f"Pic_{image_number}.jpg" #File name that will be saved
23.     try:
24.         opener=urllib.request.build_opener()
25.         opener.addheaders=[('User-Agent', 'Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/36.0.1941.0 Safari/537.36')]
26.         urllib.request.install_opener(opener)
27.         urllib.request.urlretrieve(url,os.path.join(dir_path, name))
28.     except Exception as error:
29.         print ("[~] Error Occured with %s : %s" % (name, error))
30.
31. for i in links:
32.     print (j)
33.     download_image(i,j)
34.     j+=1

```

Rottentomatoes.com

We also collect image information, so the process is, again, similar to that of joblo.com. We collect the links save them to csv and download from there. However, this website is back by javascript. We can't send the request regularly and get the full information. We need to

render the javascript part. To be able to render the javascript part, we use selenium. The result in csv file and the code are shown below.

Title	Metascore	Release_date	Image
Shazam!	91%	2-Jul	https://resizing.flixster.com/kuYtu_B1xzuyeXnjzWhEQFtlvMQ=/130x0/v1.bTsxMzAxODU5MTtqOzE4MTk1OzEyMDA7Mjc2NDs0MDk2
Avengers:	94%	30-Jul	https://resizing.flixster.com/7-hNtSRfnHzzQB6U_IKkDgW1pkk=/130x0/v1.bTsxMzAxOTkzMjttqOzE4MTk1OzEyMDA7MTY4ODsyNTAw
John Wick	90%	23-Aug	https://resizing.flixster.com/O_CeZ0_x3j2jnw7dlU4GZE45Grg=/130x0/v1.bTsxMzAyNzExOTtqOzE4MTk1OzEyMDA7MzYwMDs1NTUw
Long Shot	81%	16-Jul	https://resizing.flixster.com/-JUHHrJKNxwvPLg6g3oG_wuV6Wc=/130x0/v1.bTsxMzE1MDg4NDtqOzE4MTk3OzEyMDA7MjAwMDszMDAw
Hale Cour	97%	18-Jun	https://resizing.flixster.com/gt4X968w0EEMqcRY_Xp2NGYyN34=/130x0/v1.bTsxMjgzNzQ5NTtqOzE4MTkzOzEyMDA7Mzc4OzU1OQ

```

1. from bs4 import BeautifulSoup
2. import csv
3. from selenium import webdriver
4.
5. # Create a file to write to, add headers row
6. f = csv.writer(open('Rt_rating.csv', 'w', newline='', encoding='utf-8'))
7. f.writerow(['Title', 'Metascore', 'Release_date', 'Image'])
8.
9. #####Use Selenium to open web browser
10. browser = webdriver.Chrome() #replace with .Firefox(), or with the browser of your choice
11. url = "https://www.rottentomatoes.com/browse/cf-dvd-streaming-all/"
12. browser.get(url) #navigate to page
13.
14. #Get the script from the site
15. innerHTML = browser.execute_script("return document.body.innerHTML")
16. # execute script to scroll down the page
17. browser.execute_script("window.scrollTo(0, document.body.scrollHeight);var lenOfPage=document.body.scrollHeight;return lenOfPage;")
18.
19. #Beautiful Soup
20. soup = BeautifulSoup(innerHTML, 'html.parser')
21. last_links = soup.find_all("div", attrs={"class": "mb-movie"})
22. number_of_movies = len(last_links)
23. number_of_movies
24.
25. for i in range(0,number_of_movies):
26.     Image =last_links[i].div.a.img.get('src')
27.     Title =last_links[i].find("div", attrs={"class": "movie_info"}).h3.string
28.     Metascore =last_links[i].find("span", class_="tMeterScore").get_text(" ", strip=True)
29.     Release_date = last_links[i].find("p",class_="release-date").get_text(" ")[11:]
30.     print(Release_date)
31.     f.writerow([Title, Metascore, Release_date, Image])
32. del f
33.
34. browser.close()

```

Once we have the image links, we can download from them with the same function as we use for joblo.com. We just need to change how we want to save our image files. The code below shows how to download images and save them to the specified directory.

```

1. #Define the name of directory
2. f_name = "Rottentomato"
3.
4. # Create the directory name where the images will be saved
5. base_dir = os.getcwd()
6. dir_name = (f_name.split('/')[ -1]).split('.')[0]

```

```

7. dir_name
8. dir_path = os.path.join(base_dir, dir_name)
9.
10. #Create the directory if already not there
11. if not os.path.exists(dir_path):
12.     os.mkdir(dir_path)
13.
14. # Read the csv with links to all the image pages
15. os.getcwd()
16. df = pd.read_csv("Rottentomato_rating.csv")
17. df.columns
18. links=df.Image
19.
20. # Function to take an image url and save the image in the given directory
21. def download_image(url,image_number):
22.     print("[INFO] downloading {}".format(url))
23.     name = f"Pic_{image_number}.jpg" #File name that will be saved
24.     try:
25.         opener=urllib.request.build_opener()
26.         opener.addheaders=[('User-Agent', 'Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/36.0.1941.0 Safari/537.36')]
27.         urllib.request.install_opener(opener)
28.         urllib.request.urlretrieve(url,os.path.join(dir_path, name))
29.     except Exception as error:
30.         print ("[~] Error Occured with %s : %s" % (name, error))
31.
32. # Print the number of images
33. print ("[INFO] Downloading {} images".format(len(links)))
34. #Download ALL image in URL links list
35. for i in links:
36.     print (j)
37.     download_image(i,j)
38.     j+=1

```

Boxofficemojo.com

In this website, we extract movie news title, contents of the news and the date posted. The title and date are not difficult to get because it includes in the page shows below; however; to get the content for each news, it is a little bit harder. We need to click in the article links. The code below shows how to do that task. We first loop into every page to get every news article and in each page, we loop again to send another request to get the news content.

News

1-20 of 3,636

← Prev 20 Next 20 →

'Frozen II' Delivers \$350 Million Global Debut

November 21, 2019 19:35 PST - By Brad Brevet - Box Office News

SUNDAY AM UPDATE: Disney's [Frozen II](#) delivered a monster debut, bringing in an estimated **\$127 million** domestically, which serves as the largest animated opening ever outside the summer corridor and a record animated opening for the month of November. In fact, the performance is the **fifth largest November opening of all-time**, topping the \$125 million opening for [Harry Potter and the Deathly Hallows: Part](#)

'Ford v Ferrari' Races to #1 While 'Joker' Becomes First R-Rated Film to Ever Top \$1 Billion Globally

November 17, 2019 9:23 PST - By Brad Brevet - Box Office News

Fox's [Ford v Ferrari](#) more than lived up to the most aggressive of pre-weekend expectations, delivering a #1 performance at the domestic weekend box office. However, Sony's [Charlie's Angels](#) struggled mightily in its debut, failing to reach the lowest of expectations, which means the film's third place finish puts the weekend's overall performance into perspective. In better news, WB's [Joker](#) became the fourth DC Comics

'Ford v Ferrari' Racing Toward \$30 Million Debut, 'Charlie's Angels' Far from Heavenly

```

1. # Create a file to write to, add headers row
2. f = csv.writer(open('Boxofficemojo_news.csv', 'w', newline='', encoding="utf-8"))
3. f.writerow(['Date', 'Article', 'Content'])
4.
5.
6. pages = []
7. for i in range(1, 56):
8.     url = 'https://www.boxofficemojo.com/news/?view=&page=' + str(i) + '&p=.htm'
9.     pages.append(url)
10.
11. for item in pages:
12.     page = requests.get(item)
13.     # Create a BeautifulSoup object
14.     soup = BeautifulSoup(page.text, 'html.parser')
15.     last_links = soup.find("td", attrs={"colspan": "3"})
16.     last_links=last_links.find_all('tr')
17.
18.     number_of_article = len(last_links)
19.
20.     for i in range(1,number_of_article):
21.
22.         Date= last_links[i].find_all("td")[0].string
23.
24.         Article =str(last_links[i].find_all("td")[1].string)
25.
26.         insidelink = last_links[i].find_all("td")[1].a.get('href')
27.         insidelink = 'https://www.boxofficemojo.com'+insidelink
28.         pageinside = requests.get(insidelink)
29.         soupinside = BeautifulSoup(pageinside.text, 'html.parser')
30.         soupinside = soupinside.find_all("p", attrs={"align":"justify"})
31.         soupinside = BeautifulSoup(str(soupinside[0:-
1]),'lxml').get_text(" ",strip=True).replace(" , ", " ").replace("[ ", "").replace(" ]
", "")

```



```

32.
33.         print(Date)
34.         f.writerow([Date,Article,soupinside])

```

We also extract revenue information from this website. The code below shows how to achieve the task. The information we scraped are also shown in the table below.

Rank	Title	Studio	Worldwide_revenue	Domestic_revenue	Domestic_revenue_percent	Oversea_revenue	Oversea_revenue_percent	Year
1	engers: Endg	BV	\$2,796.30	\$858.40	30.70%	\$1,937.90	69.30%	2019
2	Avatar	Fox	\$2,789.70	\$760.50	27.30%	\$2,029.20	72.70%	2009
3	Titanic	Par.	\$2,187.50	\$659.40	30.10%	\$1,528.10	69.90%	1997
4	The Force	BV	\$2,068.20	\$936.70	45.30%	\$1,131.60	54.70%	2015

```

1. # Create a file to write to, add headers row
2. f = csv.writer(open('Boxofficemojo_revenue.csv', 'w', newline=''))
3. f.writerow(['Rank', 'Title', 'Studio', 'Worldwide_revenue', 'Domestic_revenue'
4.             , 'Domestic_revenue_percent', 'Oversea_revenue', 'Oversea_revenue_percent'
5.             , 'Year'])
6. pages = []
7. for i in range(1, 9):
8.     url = 'https://www.boxofficemojo.com/alltime/world/?pagenum=' + str(i) + '&p=.htm'
9.     pages.append(url)
10.
11. for item in pages:
12.     page = requests.get(item)
13.     # Create a BeautifulSoup object
14.     soup = BeautifulSoup(page.text, 'html.parser')
15.     last_links = soup.table.next_sibling.next_sibling.find_next('table')
16.     number_of_movie = len(last_links.find_all("tr"))
17.
18.     for i in range(1,number_of_movie):
19.
20.         Rank= last_links.find_all("tr")[i].find_all("td")[0].string
21.         Title =last_links.find_all("tr")[i].find_all("td")[1].string
22.         Studio = last_links.find_all("tr")[i].find_all("td")[2].string
23.         Worldwide_revenue = last_links.find_all("tr")[i].find_all("td")[3].string
24.         Domestic_revenue = last_links.find_all("tr")[i].find_all("td")[4].string
25.         Domestic_revenue_percent = last_links.find_all("tr")[i].find_all("td")[5].string
26.         Oversea_revenue = last_links.find_all("tr")[i].find_all("td")[6].string
27.         Oversea_revenue_percent = last_links.find_all("tr")[i].find_all("td")[7].string
28.         Year = last_links.find_all("tr")[i].find_all("td")[8].string[0:4]# [0:4] is to remove
29.         print(Year)
30.
31.         f.writerow([Rank, Title,Studio,Worldwide_revenue,Domestic_revenue,Domestic_revenue_per
32. cent,Oversea_revenue,Oversea_revenue_percent,Year])
33. del f

```

OMDB

This site allows us to send API request with the limit of 1,000 a day. It will take IMBD id to search for information. First, we need to have a list of IMDB ID. After we decide which movie

id we would like to find, we then create list of OMDb URL to send request to. The output in CSV file will be the information of the movie in json format.

[illegible]

```

1. import pandas as pd
2. import requests
3. from bs4 import BeautifulSoup
4. import csv
5. import json
6.
7. #Create link for OMDB API Using imdb id
8. list_imdbid='http://www.omdbapi.com/?i='+list_imdbid+'&apikey=3e46a856'
9. list_imdbid
10. list_imdbid=list(list_imdbid)
11. len(list_imdbid)
12.
13. list_imdbid[0][26:35]
14.
15. # Create a file to write to, add headers row
16. f = csv.writer(open('OMDB_api.csv', 'w', newline='', encoding="utf-8"))
17. f.writerow(['IMDB_ID', 'Info_json_form'])
18.
19. for item in list_imdbid:
20.     page = requests.get(item)
21.     # Create a BeautifulSoup object
22.     soup = BeautifulSoup(page.text, 'html.parser')
23.     soup=str(soup) #Change to string
24.     json_soup = json.loads(soup) #Change to JSON
25.     IMDB_ID=item[26:35] # TO get id
26.
27.     print(IMDB_ID)
28.     f.writerow([IMDB_ID,json_soup])
29.
30. del f

```

Scraping Top grossing movies of 2019 from the-numbers.com using Scrapy:

The task assigned was to collect top grossing movies in the year 2019. In order to achieve the task, we considered extracting information from the-numbers.com website. The data in tabular format can easily be extracted with scrapy. Links to scrape is given as an argument and scrapy gives you the response variable with all the information of the web page. The data from the web page can be accessed with xpath command of scrapy and tabular data can be obtained by using tr and td elements of the web page and can be parsed through easily in python as a list. The code for it is as follows:

```
1. import scrapy
2. import csv
```

```

3. import logging
4. import time
5. import csv
6. pages_length=0
7. class tutorial(scrapy.Spider):
8.     name='tutori'
9.     def start_requests(self):
10.         #links for the data to be scrapped
11.         urls = [
12.             'https://www.the-numbers.com/box-office-records/worldwide/all-
movies/cumulative/released-in-2019'
13.         ]
14.         for url in urls:
15.             #headers are changed to avoid getting 403 response
16.             yield scrapy.Request(url=url, headers={'User-
Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Geck
o) Chrome/77.0.3865.90 Safari/537.36'},callback=self.parse)
17.     def parse(self, response):
18.         global pages_length
19.         #iterating through multiple pages
20.         pages_length=pages_length+1
21.         data=[]
22.         #response body of the table to be scraped
23.         body=response.xpath("//tbody")
24.         #iterating through each rows in the table
25.         for row in body.xpath('.//tr'):
26.             #taking each column values
27.             columns=row.xpath(".//td")
28.             #extracting movie name
29.             movie_name=row.xpath("../../../b/a/text()").extract()
30.             row_data=[]
31.             #extracting other info related to the movie
32.             for col in columns:
33.                 each_column=col.xpath('text()').extract()
34.                 if len(each_column)==1:
35.                     row_data.append(each_column[0])
36.                 else:
37.                     row_data.append('')
38.             row_data[1]=movie_name[0]
39.             #appending movie info to a list
40.             data.append(row_data)
41.             #saving each row to a csv file
42.             with open('movie_data.csv', 'a',newline='',encoding="utf-8") as fd:
43.                 writer = csv.writer(fd)
44.                 for rows in data:
45.                     writer.writerow(rows)
46.             #taking the next page data if exists
47.             next_page=response.xpath('//div[@class="pagination"]//a/@href').extract()
48.             if len(next_page)>0 and len(next_page)!=pages_length:
49.                 next_page='https://www.the-numbers.com'+next_page[pages_length]
50.             yield response.follow(next_page,headers={'User-
Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Geck
o) Chrome/77.0.3865.90 Safari/537.36'}, callback=self.parse,dont_filter=True)

```

Scraping top grossing production houses from the-numbers site using Scrapy:

The task assigned was to collect top grossing production houses for year 2019 . In order to achieve the task, we considered extracting information from the-numbers.com website. The data in tabular format can easily be extracted with scrapy. Links to scrape is given as an argument and scrapy gives you the response variable with all the information of the web page. The data from the web page can be accessed with xpath command of scrapy and tabular data can be obtained by using tr and td elements of the web page and can be parsed through easily in python as a list. The code for it is as follows:

```
1. import scrapy
2. import csv
3. import logging
4. import time
5. import csv
6. class tutorial(scrapy.Spider):
7.     name='tutori'
8.     def start_requests(self):
9.         #links for the data to be scrapped
10.        urls = [
11.            'https://www.the-numbers.com/movies/production-companies/'
12.        ]
13.        for url in urls:
14.            #headers are changed to avoid getting 403 response
15.            yield scrapy.Request(url=url, headers={'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/77.0.3865.90 Safari/537.36'},callback=self.parse)
16.    def parse(self, response):
17.        data=[]
18.        #response body of the table to be scraped
19.        body=response.xpath("//tbody")
20.        #iterating through each rows in the table
21.        for row in body.xpath('.//tr'):
22.            #taking each column values
23.            columns=row.xpath(".//td")
24.            #extracting Production name
25.            production_name=row.xpath(".//td[1]//b/a/text()").extract()
26.            row_data=[]
27.            #extracting other info related to the Production company
28.            for col in columns:
29.                each_column=col.xpath('text()').extract()
30.                if len(each_column)==1:
31.                    row_data.append(each_column[0])
32.                else:
33.                    row_data.append('')
34.            row_data[0]=production_name[0]
35.            #appending production info to a list
36.            data.append(row_data)
37.            #saving each row to a csv file
38.        with open('production_data.csv', 'a',newline='',encoding="utf-8") as fd:
39.            writer = csv.writer(fd)
40.            for rows in data:
41.                writer.writerow(rows)
```

Collecting movie data from IMDB website

Subsets of IMDb data are available for access to customers for personal and non-commercial use

Data Location

The dataset files can be accessed and downloaded from <https://datasets.imdbws.com/>. The data is refreshed daily.

IMDb Dataset Details

Each dataset is contained in a gzipped, tab-separated-values (TSV) formatted file in the UTF-8 character set. The first line in each file contains headers that describe what is in each column. A '\N' is used to denote that a particular field is missing or null for that title/name. The available datasets are as follows:

title.akas.tsv.gz - Contains the following information for titles:

- titleId (string) - a tconst, an alphanumeric unique identifier of the title
- ordering (integer) - a number to uniquely identify rows for a given titleId
- title (string) - the localized title
- region (string) - the region for this version of the title
- language (string) - the language of the title
- types (array) - Enumerated set of attributes for this alternative title. One or more of the following: "alternative", "dvd", "festival", "tv", "video", "working", "original", "imdbDisplay". New values may be added in the future without warning
- attributes (array) - Additional terms to describe this alternative title, not enumerated
- isOriginalTitle (boolean) - 0: not original title; 1: original title

title.basics.tsv.gz - Contains the following information for titles:

- tconst (string) - alphanumeric unique identifier of the title
- titleType (string) - the type/format of the title (e.g. movie, short, tvseries, tvepisode, video, etc)
- primaryTitle (string) - the more popular title / the title used by the filmmakers on promotional materials at the point of release
- originalTitle (string) - original title, in the original language
- isAdult (boolean) - 0: non-adult title; 1: adult title
- startYear (YYYY) - represents the release year of a title. In the case of TV Series, it is the series start year

- endYear (YYYY) – TV Series end year. '\N' for all other title types
- runtimeMinutes – primary runtime of the title, in minutes
- genres (string array) – includes up to three genres associated with the title

title.crew.tsv.gz – Contains the director and writer information for all the titles in IMDb.
Fields include:

- tconst (string) - alphanumeric unique identifier of the title
- directors (array of nconsts) - director(s) of the given title
- writers (array of nconsts) – writer(s) of the given title

title.episode.tsv.gz – Contains the tv episode information. Fields include:

- tconst (string) - alphanumeric identifier of episode
- parentTconst (string) - alphanumeric identifier of the parent TV Series
- seasonNumber (integer) – season number the episode belongs to
- episodeNumber (integer) – episode number of the tconst in the TV series

title.principals.tsv.gz – Contains the principal cast/crew for titles

- tconst (string) - alphanumeric unique identifier of the title
- ordering (integer) – a number to uniquely identify rows for a given titleId
- nconst (string) - alphanumeric unique identifier of the name/person
- category (string) - the category of job that person was in
- job (string) - the specific job title if applicable, else '\N'
- characters (string) - the name of the character played if applicable, else '\N'

title.ratings.tsv.gz – Contains the IMDb rating and votes information for titles

- tconst (string) - alphanumeric unique identifier of the title
- averageRating – weighted average of all the individual user ratings
- numVotes - number of votes the title has received

name.basics.tsv.gz – Contains the following information for names:

- nconst (string) - alphanumeric unique identifier of the name/person
- primaryName (string)– name by which the person is most often credited
- birthYear – in YYYY format
- deathYear – in YYYY format if applicable, else '\N'
- primaryProfession (array of strings)– the top-3 professions of the person
- knownForTitles (array of tconsts) – titles the person is known for

Conclusion

Web Scraping is the technique that is widely used by researchers to collect the data from websites. Many python libraries are built to achieve this task. Scrapy is powerful and fast but need some learning curve at the beginning. BeautifulSoup is easier set up. Selenium is useful for javascript-content website. The next step we would like to suggest is to use Scrapy combined with Splash to scrape javascript-content website. This could be faster than using Selenium.