

Movie-Netnaija Web Scrapping

Importing all relevant libraries

- About Netnaija

Netnaija is a website designed for blogging ,music and movie downloads and it was created in 2016.

- Website Main Sections
 - Blogging
 - Music
 - Movies
- Website section scrapped
 - Movies
- Aims of the project

The project aims at scraping data from the movies section that has been uploaded from 2016 till date(March 2023)

```
In [1]: 1 import pandas as pd
        2 from bs4 import BeautifulSoup#Beautiful Soup is a Python library u
        3 import requests#used to send GET and POST requests to websites and
        4 from tqdm.notebook import tqdm_notebook
```

Creating reusable fuctions

```
In [2]: 1 #function to load url and parse its function
        2 def parse_html(url):
        3     #this will get url from the web
        4     response=requests.get(url)
        5     soup=BeautifulSoup(response.content,'lxml')#used for processin
        6     return soup
        7
        8
        9 #Function to print next page link
       10 def nextlink(soup):
       11     try:
       12         next_link=soup.find('a', class_="next page-numbers", href=
       13         return next_link
       14     except:
       15         return #return noting if at last page
       16
       17
```

Get the last page of the website(netnaija)

```
In [3]: 1 url='https://www.thenetnaija.net/videos/movies'
        2 soup=parse_html(url)
        3 paginating=soup.find('ul',class_='pagination')
        4 x=int(paginating.findAll('li')[5].text)
```

```
In [4]: 1 x #the number of pages is stored in variable 'x'
```

```
Out[4]: 226
```

Loop to get link of all pages

Base url

```
In [5]: 1 url='https://www.thenetnaija.net/videos/movies' #you may open on y
        2
        3 pagelink_storage=[]
        4
        5 for i in tqdm_notebook(range(x), desc='Loading....'):
        6     soup=parse_html(url)
        7     url=nextlink(soup)
        8     if not url:
        9         break
        10    pagelink_storage.append(url)
        11
        12
        13
        14
```

Loading....:

225/226 [01:31<00:00,

100%

3.85it/s]

```
In [6]: 1 len(pagelink_storage)
```

```
Out[6]: 225
```

Loop to get the links to all movies on each page

```
In [7]: 1 movie_links=[] # empty list to store movie links form each page
2 for page in tqdm_notebook(pagelink_storage, desc='Loading...'): #
3     soup= parse_html(page) # load each page and parse
4
5     # this series of code get all link to movies on each page and
6     video_files=soup.find("div", class_="video-files")
7     class_info=video_files.findAll("div", class_="info")
8     for x in class_info:
9         link=x.find("a", href=True)['href']
10        movie_links.append(link)
11
```

Loading...:

225/225 [01:09<00:00,

100%

4.45it/s]

```
In [8]: 1 len(movie_links)
```

```
Out[8]: 4040
```

```
In [ ]: 1
```

Empty list to store data we need about each movie

```
In [9]: 1 titles=[] # movie titles
2 movie_linkss=[] # movie links
3 movie_types=[] #video type
4 time_of_uplos=[] # date of upload
5 movie_lengths=[] # lenght of movie
6 num_of_comments=[] #numbers of comment
7 mo_summarys=[] # moive summary
8 Genres=[] #movie genre
9 Release_Dates=[] #release date
10 Starss=[] # actors and actress
11 Languages=[] #movie language
12 Subtitles=[] #available subtitle
13 imdb_links=[] #imdb link
```

```
In [10]: 1 for link in tqdm_notebook(movie_links, desc='Loading'):
2         soup= parse_html(link) # browse movie link and parse
3
4         #This series of code get the requiried data and append to the o
5         try:
6             title=soup.find('h1', class_="page-h1").text
7             titles.append(title)
8         except:
9             titles.append(' ')
10        post_meta=soup.find("div", class_="post-meta")
11        try:
12            movie_link=post_meta.find('a', href=True)['href']
13            movie_linkss.append(movie_link)
14        except:
15            movie_links.append(' ')
16        meta_one=soup.findAll('span', class_='meta-one')
17        try:
18            movie_type=meta_one[0].text.split()
19            movie_types.append(movie_type)
20        except:
21            movie_types.append(' ')
22        #x=meta_one[1].text.split()
23        try:
24            x=meta_one[1].text.split()
25            time_of_uplo=' '.join(x)
26            time_of_uplos.append(time_of_uplo)
27        except:
28            time_of_uplos.append(' ')
29        try:
30            movie_length=meta_one[2].text.split()
31            movie_lengths.append(movie_length)
32        except:
33            movie_lengths.append(' ')
34        try:
35            num_of_comment=meta_one[3].text.split()
36            num_of_comments.append(num_of_comment)
37        except:
38            num_of_comments.append('0')
39        try:
40            mo_summary=soup.find('p').next_element
41            mo_summarys.append(mo_summary)
42        except:
43            mo_summarys.append(' ')
44        try:
45            block=soup.find('blockquote', class_='quote-content')
46            y=block.findAll('p')
47
48            try:
49                Genre=y[1].text.split(':')[1:]
50                Genres.append(Genre)
51            except:
52                Genres.append('missing')
53            try:
54                Release_Date=y[2].text
55                Release_Dates.append(Release_Date)
```

```
56         except:
57             Release_Dates.append('missing')
58
59         try:
60             Stars=y[3].text.split(':')[1:]
61             Starss.append(Stars)
62         except:
63             Starss.append('missing')
64
65         try:
66             Language=y[5].text.split(':')[1:]
67             Languages.append(Language)
68         except:
69             Languages.append('missing')
70
71         try:
72             Subtitle=y[6].text.split(':')[1:]
73             Subtitles.append(Subtitle)
74         except:
75             Subtitles.append('missing')
76
77         except:
78             Genres.append('missing')
79             Release_Dates.append('missing')
80             Starss.append('missing')
81             Languages.append('missing')
82             Subtitles.append('missing')
83
84         try:
85             imdb_link=block.find('a', href=True)['href']
86             imdb_links.append(imdb_link)
87         except:
88             imdb_links.append('missing')
```

Loading:

4040/4040 [34:06<00:00,

100%

1.90it/s]

Creating a table of all data with pandas dataframe

```
In [11]: 1 df=pd.DataFrame({"titles":titles,
2                       "movie_types":movie_types,
3                       "time_of_uplos":time_of_uplos,
4                       "movie_lengths":movie_lengths,
5                       "num_of_comments":num_of_comments,
6                       "Genres":Genres,
7                       "Release_Dates":Release_Dates,
8                       "Starss":Starss,
9                       "Languages":Languages,
10                      "Subtitles":Subtitles,
11                      "movie_linkss":movie_linkss,
12                      "imdb_links":imdb_links,
13                      "mo_summaries":mo_summaries,
14                      })
```

save data in cvs and excel format

```
In [12]: 1 df.to_csv('netnaija.csv')
2 df.to_excel('netnaija.xlsx')
```

```
In [13]: 1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4040 entries, 0 to 4039
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
---  -
0   titles                4040 non-null   object
1   movie_types           4040 non-null   object
2   time_of_uplos         4040 non-null   object
3   movie_lengths         4040 non-null   object
4   num_of_comments       4040 non-null   object
5   Genres                4040 non-null   object
6   Release_Dates         4040 non-null   object
7   Starss                4040 non-null   object
8   Languages             4040 non-null   object
9   Subtitles             4040 non-null   object
10  movie_linkss          4040 non-null   object
11  imdb_links            4040 non-null   object
12  mo_summaries          4040 non-null   object
dtypes: object(13)
memory usage: 410.4+ KB
```

In [15]: 1 df.head(8)

Out[15]:

	titles	movie_types	time_of_uplos	movie_lengths	num_of_comments	Genres	Relea
0	Kuttey (2023) [Indian]	[Movies]	Mar 20	[01:48:51]	[68]	[Action, Comedy, Crime, Thriller]	Rele Jar
1	Bad City (2022) [Japanese]	[Movies]	Mar 17	[01:58:04]	[50]	[Action, Crime]	Rele Jar
2	In His Shadow (2023) [French]	[Movies]	Mar 17	[01:29:56]	[24]	[Crime, Drama, Family, Thriller]	Rele Mai
3	Noise (2023) [Dutch]	[Movies]	Mar 17	[01:30:34]	[6]	[Drama, Mystery, Thriller]	Rele Mai
4	Boston Strangler (2023)	[Movies]	Mar 17	[01:52:18]	[17]	[Crime, Drama, History, Thriller]	Rele Mai (Unit
5	Vaathi (2023) [Indian]	[Movies]	Mar 17	[05:00:52]	[38]	[Action, Comedy, Drama, Romance]	Rele Fel
6	Haunted Universities 2 (2022) [Thai]	[Movies]	Mar 15	[02:03:59]	[28]	[Comedy, Horror, Thriller]	Rele Mai
7	The Lake (2022) [Thai]	[Movies]	Mar 15	[01:44:20]	[55]	[Drama, Horror, Sci-Fi, Thriller]	Rele Aug