

## Github Topics Details - Web Scrapping Project



### Introduction:

Web scrapping is the process of parsing and extracting data from the websites. It is very useful when we want to collect data from the websites for our work.

I'm going to scrap Github's topics page(<https://github.com/topics>) to know-

- Different topics present on Github.
- For each topics,
  - topic titles, description and topic URL.
  - popular repository.
- And for each popular repositories-
  - username(created by).
  - number of stars
  - URL of the repository.
  - number of times forked.
  - number of times committed.
  - last committed time.

### Required data format:

Topics	Description	Topic_URL	Popular_Repository	Username	Repository_URL	Number_Of_Stars	Forked_Count	Total_commits
3D	--descr--	<a href="https://github.com/topics/3d">https://github.com/topics/3d</a>	three.js	mrdoob	<a href="https://github.com/mrdoob/three.js">https://github.com/mrdoob/three.js</a>	73300	28708	37892

### Tools to be used:

Python, Pandas, BeautifulSoup, Requests

### Let's scrap it...

#### Part I

First, let's import required libraries.

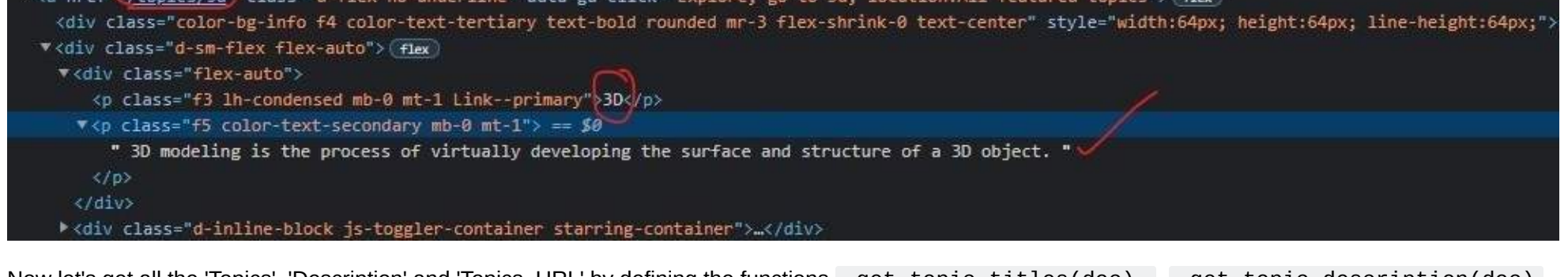
```
In [1]: import requests
from bs4 import BeautifulSoup
import pandas as pd
```

Now let's define a function which loads all the required topics pages and returns BeautifulSoup object.

```
In [2]: def topics_pages_loader():
page_content=''
i=1
while i<7: # i<7 because there are only 6 topics pages
url = 'https://github.com/topics?page={}'.format(i)
r = requests.get(url)
if r.status_code != 200:
i+=1 #if status code!=200, it means failed to load the page. So reloading the page by decrementing the i
else:
page_content += '\n' + r.text
i+=1
doc = BeautifulSoup(page_content, 'html.parser')
return doc
```

We have loaded all the required topics pages successfully.

We can see the required data in the below mentioned pic, we need to extract that from the respective tags.



Now let's get all the 'Topics', 'Description' and 'Topics\_URL' by defining the functions 'get\_topic\_titles(doc)', 'get\_topic\_description(doc)' and 'get\_topic\_urls(doc)' which return topics titles, description and topics URL respectively.

```
In [3]: def get_topic_titles(doc):
topic_selection_class = 'f3 lh-condensed mb-0 mt-1 Link--primary'
topic_p_tags = doc.find_all('p',{'class':topic_selection_class})
topic_titles = []
for tag in topic_p_tags:
topic_titles.append(tag.text)
return topic_titles

def get_topic_description(doc):
desc_selection_class = 'f5 color-text-secondary mb-0 mt-1'
desc_ptags = doc.find_all('p',{'class':desc_selection_class})
topic_descs = []
for tag in desc_ptags:
topic_descs.append(tag.text.strip())
return topic_descs

def get_topic_urls(doc):
link_selection_class = 'd-flex no-underline'
topic_link_tags = doc.find_all('a',{'class':link_selection_class})
topic_urls = []
base_url = 'https://github.com'
for tag in topic_link_tags:
topic_urls.append(base_url + tag['href'])
return topic_urls
```

Now let's define a function 'scrape\_topics()' which returns the dataframe containing the columns 'Topics', 'Description', and 'Topic\_URL'.

```
In [4]: def scrape_topics():
doc = topics_pages_loader()
topics_dict = {
'Topics': get_topic_titles(doc),
'Description': get_topic_description(doc),
'Topic_URL': get_topic_urls(doc)
}
return pd.DataFrame(topics_dict)
```

```
In [5]: df_topics = scrape_topics()
df_topics.head()
```

Out[5]:

	Topics	Description	Topic_URL
0	3D	3D modeling is the process of virtually develo...	<a href="https://github.com/topics/3d">https://github.com/topics/3d</a>
1	Ajax	Ajax is a technique for creating interactive w...	<a href="https://github.com/topics/ajax">https://github.com/topics/ajax</a>
2	Algorithm	Algorithms are self-contained sequences that c...	<a href="https://github.com/topics/algorithm">https://github.com/topics/algorithm</a>
3	Amp	Amp is a non-blocking concurrency framework fo...	<a href="https://github.com/topics/amphtml">https://github.com/topics/amphtml</a>
4	Android	Android is an operating system built by Google...	<a href="https://github.com/topics/android">https://github.com/topics/android</a>

This is how our required semi-dataframe looks like. We have scraped first 3 columns which gives us the basic informations of the topics.

This is not complete CSV file we wanted but we can store this also as another CSV file which basically gives us the basic informations of the topics.

We can store this in a CSV format as shown below.

```
In [6]: df_topics.to_csv('Github_topics.csv',index=False)
```

#### Part II

Now let's gather some data about popular repositories for each topics.

We can utilize the 'Topic\_URL' column from the above dataframe to get some of the required information of the popular repositories.

```
In [7]: topic_urls = df_topics['Topic_URL']

<os class="f3 color-text-secondary text-normal lh-condensed">
<div class="f3 lh-condensed mb-0 mt-1 Link--primary">
<a href="https://github.com/topics/3d" class="d-flex no-underline" data-ga-click="Explore, go to 3d, location:All featured topics"> (New)
<div class="color-bg-info f4 color-text-tertiary text-bold rounded mw-3 flex-shrink-0 text-center" style="width:64px; height:64px; line-height:64px;">
<div class="d-flex flex-auto">
<div class="flex-auto">
<div class="f3 lh-condensed mb-0 mt-1 Link--primary">3D</div>
<div class="f5 color-text-secondary mb-0 mt-1">
3D modeling is the process of virtually developing the surface and structure of a 3D object.
</div>
</div>
<div class="d-inline-block js-toggle-container starring-container"></div>
```

As we can see that repository name and username are present in 'h3' tag, let's define a function which returns us these values. And also let's get the repository URL from the same function only.

```
In [8]: def get_repo_info(h3_tags):
atags = h3_tags.find_all('a')
username = atags[0].text.strip()
repo_name = atags[1].text.strip()
repo_url = 'https://github.com' + atags[1]['href']

return username, repo_name, repo_url
```

The function 'popular\_repo\_info()' loads each topics pages and get the popular repository(most starred) then returns 'Popular\_Repository', 'Username' and 'Repo\_URL'.

```
In [9]: def popular_repo_info():
repo_details = {'Popular_repository(most_starred)':[], 'Repo_Username':[], 'Repo_URL':[]}
i=0
while i<len(topic_urls):
url = topic_urls[i] + '?o=desc&s=stars'
r = requests.get(url)
if r.status_code != 200:
i+=1
else:
doc = BeautifulSoup(r.text, 'html.parser')
h3_tags = doc.find_all('h3',{'class': 'f3 color-text-secondary text-normal lh-condensed'}, limit=1)
repo_info = get_repo_info(h3_tags[0])
repo_details['Repo_Username'].append(repo_info[0])
repo_details['Popular_repository(most_starred)'].append(repo_info[1])
repo_details['Repo_URL'].append(repo_info[2])
i+=1
return pd.DataFrame(repo_details)
```

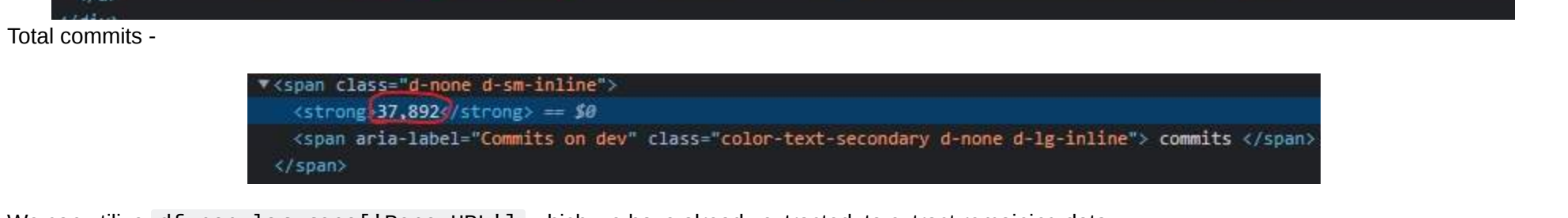
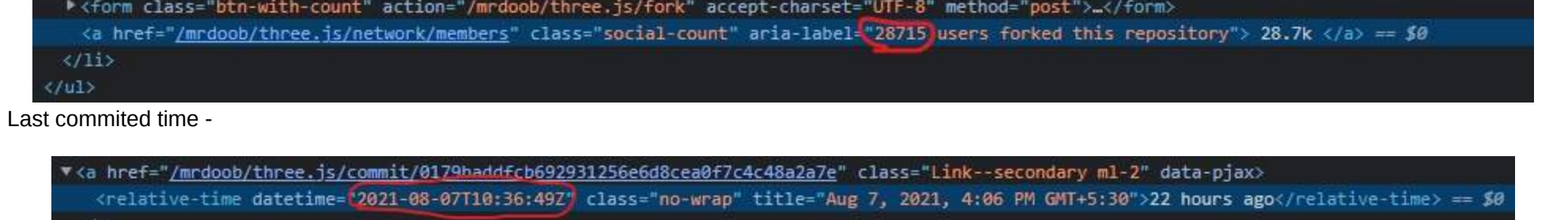
```
In [10]: df_popular_repo = popular_repo_info()
df_popular_repo.head()
```

Out[10]:

	Popular_repository(most_starred)	Repo_Username	Repo_URL
0	three.js	mrdoob	<a href="https://github.com/mrdoob/three.js">https://github.com/mrdoob/three.js</a>
1	infinite-scroll	metafizzy	<a href="https://github.com/metafizzy/infinite-scroll">https://github.com/metafizzy/infinite-scroll</a>
2	coding-interview-university	qwasham	<a href="https://github.com/qwasham/coding-interview-university">https://github.com/qwasham/coding-interview-university</a>
3	amp	ampphp	<a href="https://github.com/ampphtml">https://github.com/ampphtml</a>
4	flutter	flutter	<a href="https://github.com/flutter/flutter">https://github.com/flutter/flutter</a>

Now we have the basic informations of the popular repositories like 'Popular\_Repository', 'Username' and 'Repo\_URL'. Let's get some other informations for these popular repositories like 'Number\_Of\_Stars', 'Forked\_Count', 'Total\_Commits' and 'Last\_Committed'.

First let's have a look where our required data is located.



We can utilize 'df\_popular\_repo['Repo\_URL']' which we have already extracted, to extract remaining data.

```
In [11]: repo_urls = df_popular_repo['Repo_URL']
```

Let's define a function 'get\_repo\_info2(star\_atags, forks\_atags, commit\_span\_tags, last\_commit\_atag)' which returns us all the values we are expecting(as above).

```
In [12]: def get_repo_info2(star_atags, forks_atags, commit_span_tags, last_commit_atag):
stars = int(star_atags[0]['aria-label'].split()[0])
forks = int(forks_atags[1]['aria-label'].split()[0])
commits = int(commit_span_tags[1].strong.text.replace(',','')) if len(commit_span_tags)==2 else int(commit_span_tags[0].strong.text.replace(',',''))
last_updated = last_commit_atag[0].find_all('relative-time')[0]['datetime'] if len(last_commit_atag)==1 else None

return stars, forks, commits, last_updated
```

The function 'popular\_repo\_info2()' which we are going to define, loads the pages to get the required values (Number\_Of\_Stars, Forked\_Count, Total\_Commits and Last\_Committed) and then returns a dataframe of these values.

```
In [13]: def popular_repo_info2():
repo_details2 = {'Number_Of_Stars':[], 'Forked_count':[], 'Total_commits':[], 'Last_committed':[]}
i=0
while i<len(repo_urls):
url = repo_urls[i]
r = requests.get(url)
if r.status_code != 200:
i+=1
else:
doc = BeautifulSoup(r.text, 'html.parser')

star_atags = doc.find_all('a',{'class': 'social-count js-social-count'})
forks_atags = doc.find_all('a',{'class': 'social-count'})
commit_span_tags = doc.find_all('span',{'class': 'd-none d-sm-inline'})
last_commit_atag = doc.find_all('a',{'class': 'Link--secondary ml-2'})

repo_info2 = get_repo_info2(star_atags, forks_atags, commit_span_tags, last_commit_atag)
repo_details2['Number_Of_Stars'].append(repo_info2[0])
repo_details2['Forked_count'].append(repo_info2[1])
repo_details2['Total_commits'].append(repo_info2[2])
repo_details2['Last_committed'].append(repo_info2[3])
i+=1
return pd.DataFrame(repo_details2)
```

```
In [14]: df_pop_repo_info2 = popular_repo_info2()
df_pop_repo_info2.head()
```

Out[14]:

	Number_Of_Stars	Forked_count	Total_commits	Last committed
0	73372	28744	37899	2021-08-10T10:15:21Z
1	6932	1760	471	2021-01-03T19:57:17Z
2	189361	50869	1735	2021-07-31T17:45:08Z
3	3303	184	1385	2021-07-16T20:25:08Z
4	127074	18224	25382	2021-08-09T22:21:17Z

Now let's concat all these dataframes to get the final required outcome.

```
In [15]: df = pd.concat([df_topics, df_popular_repo, df_pop_repo_info2], axis=1)
df.head()
```

Out[15]:

	Topics	Description	Topic_URL	Popular_repository(most_starred)	Repo_Username	Repo_URL	Number_Of_Stars
0	3D	3D modeling is the process of virtually develo...	<a href="https://github.com/topics/3d">https://github.com/topics/3d</a>	three.js	mrdoob	<a href="https://github.com/mrdoob/three.js">https://github.com/mrdoob/three.js</a>	73372
1	Ajax	Ajax is a technique for creating interactive w...	<a href="https://github.com/topics/ajax">https://github.com/topics/ajax</a>	infinite-scroll	metafizzy	<a href="https://github.com/metafizzy/infinite-scroll">https://github.com/metafizzy/infinite-scroll</a>	6932
2	Algorithm	Algorithms are self-contained sequences that c...	<a href="https://github.com/topics/algorithm">https://github.com/topics/algorithm</a>	coding-interview-university	qwasham	<a href="https://github.com/qwasham/coding-interview-university">https://github.com/qwasham/coding-interview-university</a>	189361
3	Amp	Amp is a non-blocking concurrency framework fo...	<a href="https://github.com/topics/amphtml">https://github.com/topics/amphtml</a>	amp	ampphp	<a href="https://github.com/ampphtml">https://github.com/ampphtml</a>	3303
4	Android	Android is an operating system built by Google...	<a href="https://github.com/topics/android">https://github.com/topics/android</a>	flutter	flutter	<a href="https://github.com/flutter/flutter">https://github.com/flutter/flutter</a>	127074

We have extracted all the required data. Now the final work is to store this data into a CSV file.

```
In [16]: df.to_csv('Github_topics_detailed.csv')
```

Here we come to an end of this project.

Thank you...