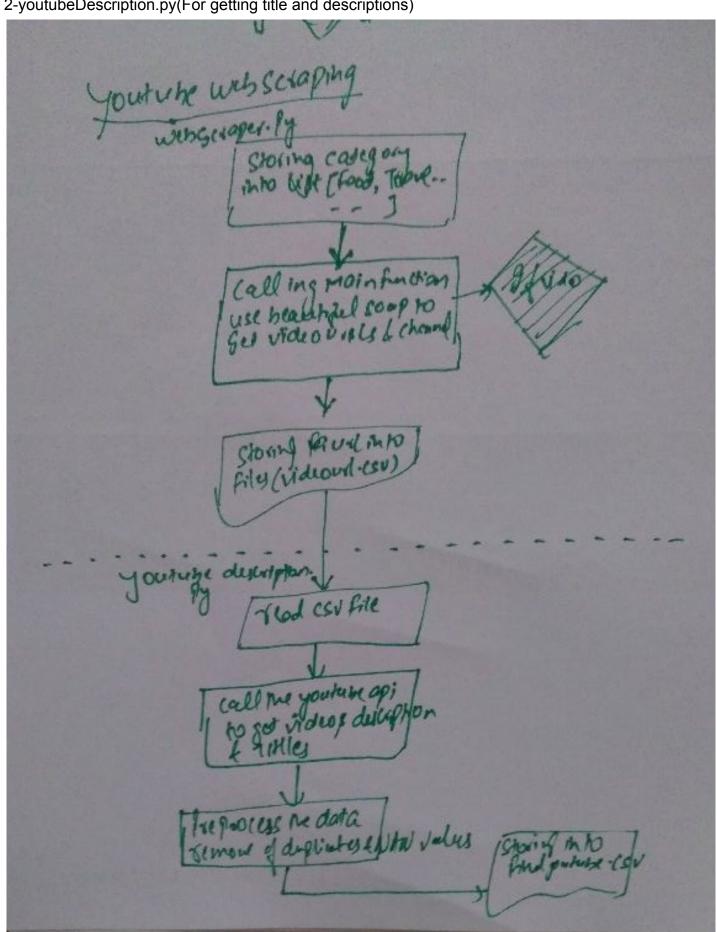# Youtube Webscraper

Youtube webscraper cantain 2 python file
1-WebScraper.py(For getting urls)
2-youtubeDescription.py(For getting title and descriptions)

## 1.WebScraper.py

Webscraper contain one function mainFunction()

If we search Food on youtube we will get video,playlist and channel related to foods.so we have to get url of videos and channel names.

We create a list searchValue and store all categories and call the mainFunction()

```
50  searchValue=['Travel+Blogs','Science+and+Technology','Food','Manufacturing','History','Art+and+Music']
51  mainFunction(searchValue)
```

Initilization lists
youtubeIDs is used to store video ids
youtubeChannel is used to store channel name
videoCate is used to store categories

```
12  youtubeIDs=[]
13  youtubeChannel=[]
14  videoCate=[]
```

Inside mainfunction()
Initializing the youtubeurl,page and count(starting page) and pages(end page)

```
16      youtubeUrl="https://www.youtube.com/results?search_query="
17      page = "&page="
18      count=1
19      pages = 20
20      searchQuery=searchValue
```

```python
15  def mainFunction(searchValue,channel=1):
16      youtubeUrl="https://www.youtube.com/results?search_query="
17      page = "&page="
18      count=1
19      pages = 20
20      searchQuery=searchValue
21      for category in searchQuery:
22          count=1
23          while count <= pages:
24              scrapeURL = youtubeUrl + str(category) + page + str(count)
25              print(category)
26              source = requests.get(scrapeURL).text
27              soup = BeautifulSoup(source, 'lxml')
28              #getting the div yt-lockup-content
29              for content in soup.find_all('div', class_ = "yt-lockup-content"):
30                  try:
31                      ID=content.h3.a
32                      matching=bool('/watch' in ID.get('href'))
33                      if(matching):
34                          youtubeIDs.append(ID.get('href'))
35                          videoCate.append(category)
36                      else:
37                          if(channel):
38                              youtubeChannel.append(channelTitle(content))
39                  except Exception as e:
40                      print(e)
41                      print("Exception")
42                      description = None
43              #increasing the count
44          count=count+1
```

For every category in list we search on youtube using beautiful soap.then find all the div tag containing
yt-lockup-content and store ID in hyper reference
If Id is video then

        1-Video link store in YoutubeIDs list

        2-videoCate store category related to video

Else

        Call the function channelTitle

        Store the channel name in youtubeChannel list

After mainFunction finish then we again call mainFunction(youtubeChannel) to get video from youtube channel

```
48    searchValue=['Travel+Blogs','Science+and+Technology','Food','Manufacturing','History','Art+and+Music']
49    mainFunction(searchValue)
50    #Getting video of youtubeChannel
51    mainFunction(youtubeChannel,channel=0)
52    df = {'Videourl': youtubeIDs,'Category':videoCate}
53    df2=pd.DataFrame(df)
54    #storing Youtube videos link into csv file
55    df2.to_csv("Videourl.csv",index=False)
```

Storing the dataFrame into videourl.csv

## 2.youtubeDescription.py
Main function of of youtubeDescription is to get video youtubeDescription and title.i use youtubeapi to get description of videos
Initilized DEVELOPER_KEY

```
1     #importing library
2     from apiclient.discovery import build
3     import pprint
4     import pandas as pd
5
6     #develper keys
7     DEVELOPER_KEY = "KEY"
8     YOUTUBE_API_SERVICE_NAME = "youtube"
9     YOUTUBE_API_VERSION = "v3"
10    youtube = build(YOUTUBE_API_SERVICE_NAME,YOUTUBE_API_VERSION,developerKey = DEVELOPER_KEY)
11
12    #function to get videos description and title
```

Reading the VideoUrl.csv file which contain videos url
Initilization lists
Description is used to store video description
Title is used to store video titles
video_ids store video is where there is no description
For very videoIds we call video_details function

```
24    dflink=pd.read_csv("Videourl.csv")
25    Description=[]
26    Title=[]
27    video_ids=[]
28    #removing the videourl where no description founf
29    for x in dflink["Videourl"]:
30        newstr = x.replace("/watch?v=", "")
31        video_details(newstr)
32
```

video_detail call youtubeapi and get result of every video url
If result not found store the video url into video_ids list

```
12   #function to get videos description and title
13   def video_details(video_id):
14       list_videos_byid = youtube.videos().list(id = video_id,part = "snippet").execute()
15
16       results = list_videos_byid.get("items", [])
17       if(results):
18           for result in results:
19               Description.append(result["snippet"]["description"])
20               Title.append(result["snippet"]["title"])
21       else:
22           video_ids.append("/watch?v="+video_id)
23
```

Removing all the url where there is no description and storing the data frame into FinalYoutube.csv

```
33   list1 = [item for item in video_ids if item not in dflink["Videourl"]]
34   for y in list1:
35       indexs=dflink[(dflink["Videourl"]==y)]
36       dflink=dflink.drop(indexs.index[0])
37
38   #storing title into DataFrame
39   dflink['Title']=Title
40   #storing description into DataFrame
41   dflink['Description']=Description
42   dflink.drop_duplicates(subset='Videourl', inplace=True)
43   #storing data into csv
44   dflink.to_csv("finalYoutube.csv",index=False)
45
```