| PROJECT TITLE | Indian Human Image Classification Model |
|---|---|
| GROUP NO. | B8 |
| AGENDA | Tools & Strategy |
| DATE | 07/04/2021 |
| DESCRIPTION | Different Methods for collecting data |

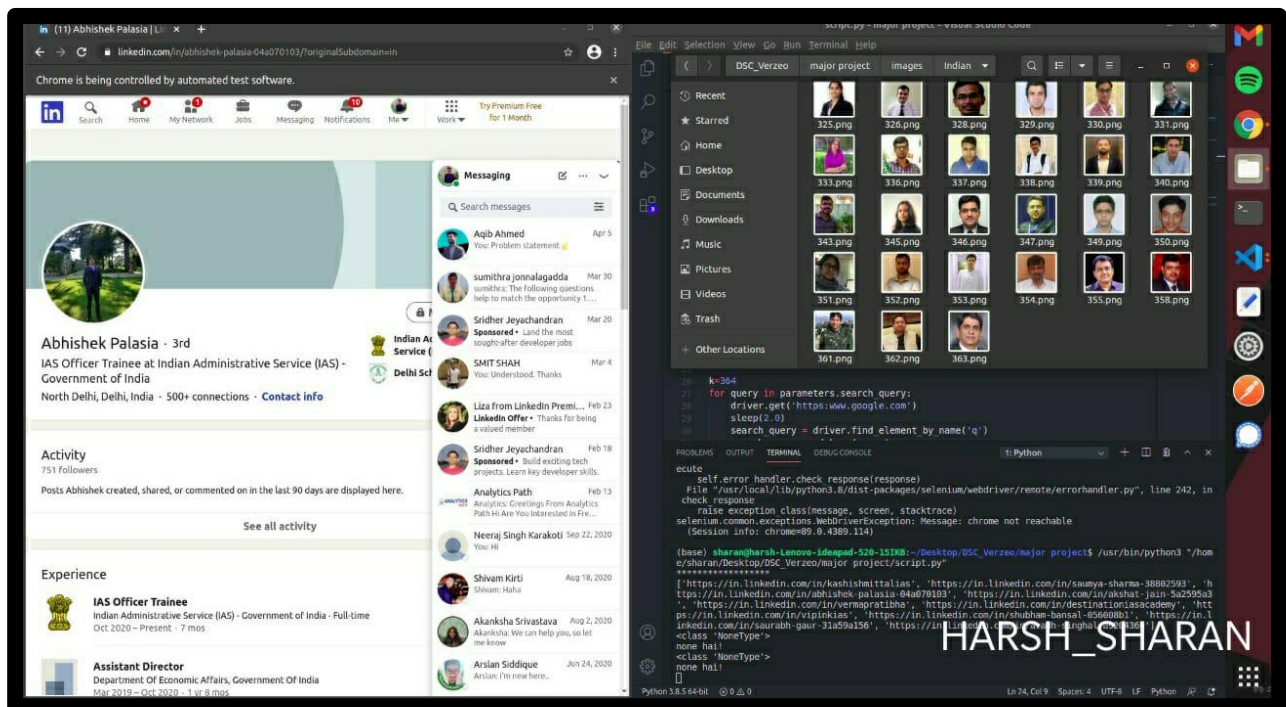–> OUR TEAM HAS COME UP WITH SEVERAL IDEAS TO EXTRACT THE REQUIRED DATASET FOR THIS PROJECT.

–> ALL METHODS MENTIONED HAVE BEEN TRIED AND TESTED BY ONE OR MORE GROUP MEMBERS WITH THE SCRIPTS RUNNING PERFECTLY(the screenshots have been attached for reference ).

# 1. Manually:

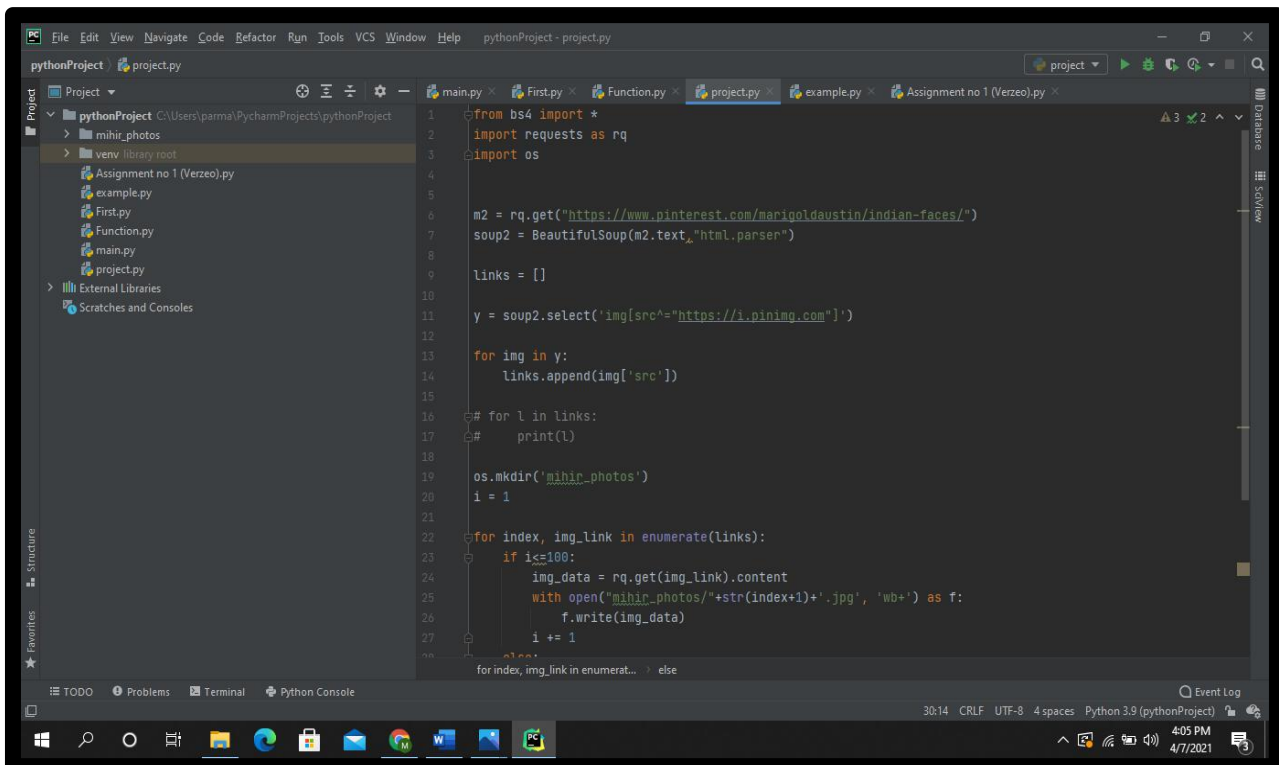We will search from Google, Youtube, and many different resources and we will be able to get images

# 2. Selenium (Automation):

By using a selenium webdriver we will be able to get images from a given URL and after applying a keyword we will able to get more images.

## **3.** Crawler:

Crawler can be used to find links to an image and after obtaining the link we will make a directory to our local system and get all images based on the path which we have given.

```python
from bs4 import *
import requests as rq
import os


m2 = rq.get("https://www.pinterest.com/marigoldaustin/indian-faces/")
soup2 = BeautifulSoup(m2.text, "html.parser")

links = []

y = soup2.select('img[src^="https://i.pinimg.com"]')

for img in y:
    links.append(img['src'])

# for l in links:
#     print(l)

os.mkdir('mihir_photos')
i = 1

for index, img_link in enumerate(links):
    if i<=100:
        img_data = rq.get(img_link).content
        with open("mihir_photos/"+str(index+1)+'.jpg', 'wb+') as f:
            f.write(img_data)
        i += 1
```

# **4.** From Softwares:

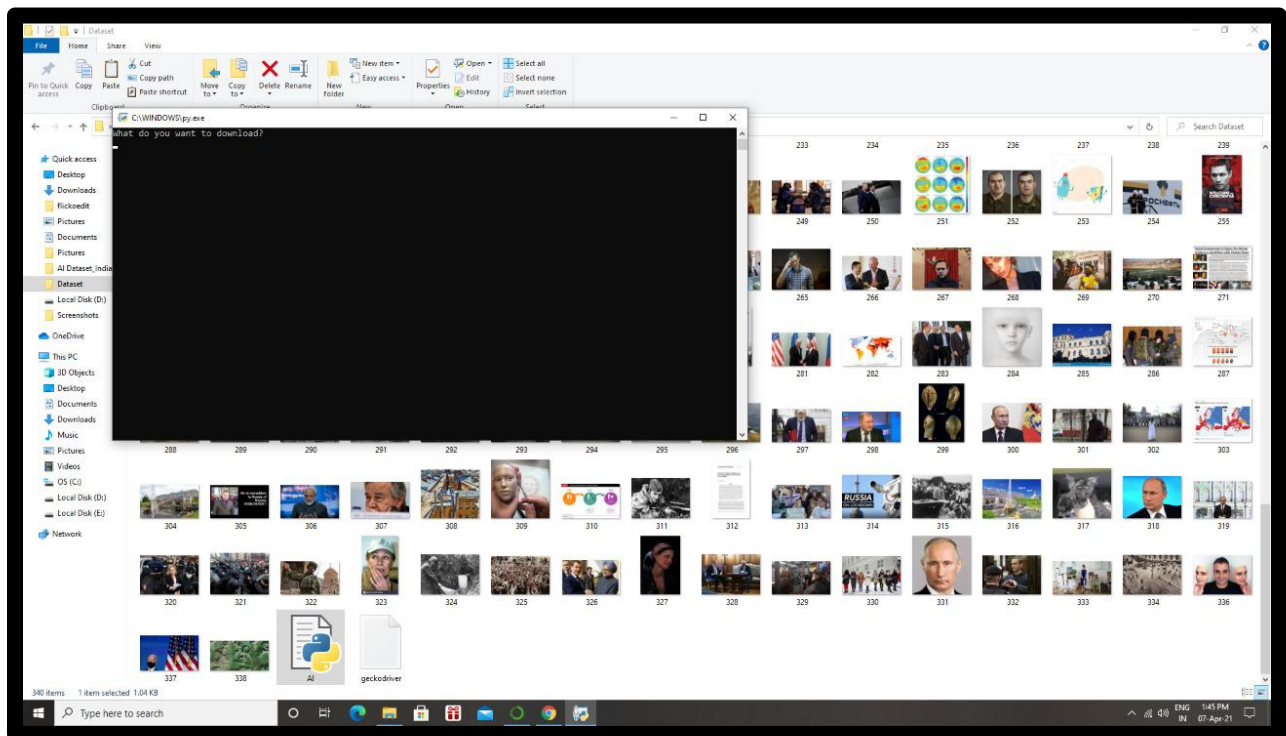(like Octoparse)
By simply providing the URL
of any website and applying
some settings it will fetch
all the images from the
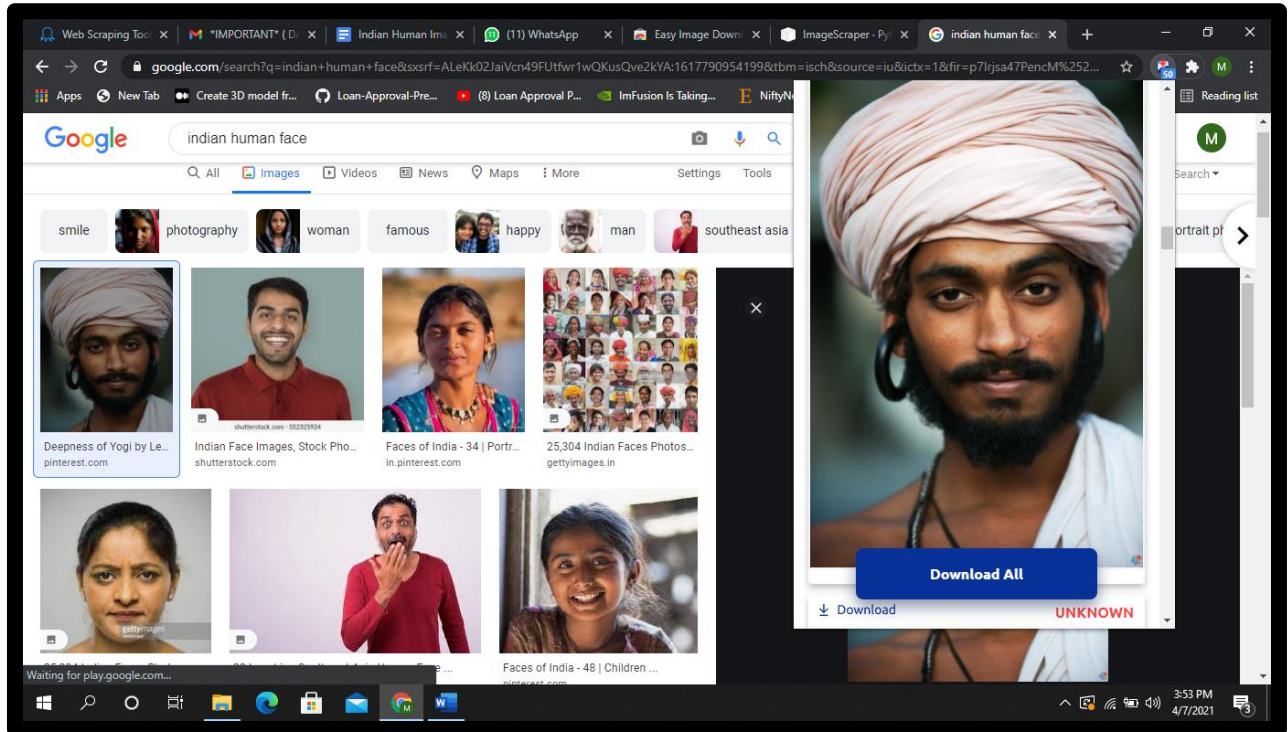website.

# 5. ImageScraper:

It is a python library. Simply giving a command (let "image-scraper ananth.co.in/test.html") to the command prompt it will give you the desired output.

**6.** From chrome extensions:

We will be using some chrome extensions in which after giving a particular keyword to any search box it will download all images of that particular keyword.



-----------------END-----------------