



VIT[®]

Vellore Institute of Technology

(Deemed to be University under section 3 of UGC Act, 1956)

GDG – 2 CREDIT COURSE REPORT

Name – Anirudh Batra

Reg No – 17BCE2145

Topic –

SCRAP THE GSOC WEBSITE

Abstract-

This project uses flask which is a python based microframework to create the API and we use beautiful soup4 to scrape the website and bring out the important information that is needed.

Inroduction-

Creating an API using flask framework and scraping the GSOC website to fetch details of all the organisations in Google summer of code and in the API send their name, link to their website, description, the technologies they use and their contact email.

Methodology-

```
1 import requests
2 import json
3 from bs4 import BeautifulSoup as bsp
4 from flask import Flask
5
6 main_org_url = 'https://summerofcode.withgoogle.com/archive/2017/organizations/'
7 base_url = "https://summerofcode.withgoogle.com"
8 app = Flask(__name__)
```

In this part of the code we have imported all the libraries that were needed for this project. Also we have defined the URLs required.

```
11 def org_info():
12     fetch = requests.get(main_org_url)
13     html = fetch.content
14     b_soup1 = bsp(html, "html.parser")
15     fetch_org = b_soup1.findAll("li", {'class': 'organization-card__container'})
16     final_result = get_details(fetch_org)
17     print(final_result)
18     return json.dumps(final_result)
19
```

In this part we access the gsoc website and fetch the details by putting a request and storing it in fetch variable. Then we scrape it using BeautifulSoup and save it in b_soup1 (Name of the organisation) and we store the data in the final_result by calling the function get_details which will be defined in the next segment.

```
21 def get_details(fetch_org):
22     extracted_result = list()
23     count = 0
24     for item in fetch_org:
25         purl = item.find('a', {'class': 'organization-card__link'})
26         organisation_name = item['aria-label']
27
28         information = item.find('div', {'class': 'organization-card__tagline font-black-54'})
29         information = information.text
30         p_link = base_url + purl['href']
31         page = requests.get(p_link)
32         if page.status_code != 200:
33             break
34         p_link = base_url + purl['href']
35         responsel = requests.get(p_link)
36         html1 = responsel.content
37         b_soup2 = bsp(html1, "html.parser")
38         organisation_link = b_soup2.find("a", {"class": "org__link"})
39         organisation_link = organisation_link.text
40         tech_info = b_soup2.findAll("li", {"class": "organization__tag organization__tag--technology"})
41         technology = []
42         for t_tech in tech_info:
43             technology.append(t_tech.text)
44         t_topics = b_soup2.findAll("li", {"class": "organization__tag organization__tag--topic"})
45         topics = []
46         for i in t_topics:
47             topics.append(i.text)
48         count += 1
49         print(count)
50         extracted_result.append({
51             'organization_name': organisation_name,
52             'description': information,
53             'link': organisation_link,
54             'technologies': technology,
55             'topics': topics
56         })
57     #Extracting Only 15 organization details Change below count value to get details of more.
58     if count == 15:
59         return extracted_result
60
```

In the function `get_details` first of we have stored the result in list format to the variable `extracted_result`. Now we have started a loop for extracting the organisation details which include the technology, link to their home site and the description of the company. The result is in Json format we need to convert it to dictionary format.

Result-

```
anirudh@the3d3n: ~/Downloads/pyppy$ python gdg_gsoc1.py
* Serving Flask app "gdg_gsoc1" (lazy loading)
* Environment: production
WARNING: Do not use the development server in a production environment.
Use a production WSGI server instead.
* Debug mode: on
* Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)
* Restarting with stat
* Debugger is active!
* Debugger PIN: 299-845-447
127.0.0.1 ~ - [24/Mar/2019 19:48:16] "GET / HTTP/1.1" 404 -
1
2
[{"organization_name": 'u'52\\xb0North Initiative for Geospatial Open Source Software GmbH', 'technologies': [u'web services', u'ogc standards', u'java', u'javascript', u'web'], 'link':
'u\\Nv\\t\\t\\http://52north.org/\\N\\t\\t\\t', 'description': 'u'52\\xb0North works on innovative ideas and technologies in geoinformatics.', 'topics': [u'geoinformatics', u'sensor web', u'
web-based geoprocessing', u'spatial data infrastructures', u'spatial information']}, {"organization_name": 'u>AboutCode', 'technologies': [u'python', u'C/C++', u'javascript', u' shell
script', u'static analysis'], 'link': 'u'\\Nv\\t\\t\\t\\http://aboutcode.org\\N\\t\\t\\t', 'description': 'u'Open Source for Open Source software license, origin and packages discovery', 'topic
s': [u'free and open source software license and origin', u'package and dependencies licensing and origin', u'package vulnerabilities and security', u'code scan and matching', u'code
analysis and spdx']}]]
127.0.0.1 ~ - [24/Mar/2019 19:48:29] "GET /org_info HTTP/1.1" 200 -
1
2
[{"organization_name": 'u'52\\xb0North Initiative for Geospatial Open Source Software GmbH', 'technologies': [u'web services', u'ogc standards', u'java', u'javascript', u'web'], 'link':
'u\\Nv\\t\\t\\http://52north.org/\\N\\t\\t\\t', 'description': 'u'52\\xb0North works on innovative ideas and technologies in geoinformatics.', 'topics': [u'geoinformatics', u'sensor web', u'
web-based geoprocessing', u'spatial data infrastructures', u'spatial information']}, {"organization_name": 'u>AboutCode', 'technologies': [u'python', u'C/C++', u'javascript', u' shell
script', u'static analysis'], 'link': 'u'\\Nv\\t\\t\\t\\http://aboutcode.org\\N\\t\\t\\t', 'description': 'u'Open Source for Open Source software license, origin and packages discovery', 'topic
s': [u'free and open source software license and origin', u'package and dependencies licensing and origin', u'package vulnerabilities and security', u'code scan and matching', u'code
analysis and spdx']}]]
127.0.0.1 ~ - [24/Mar/2019 19:48:50] "GET /org_info HTTP/1.1" 200 -
```

This is the terminal output of the code



This is the json format output of all the 15 organisation we requested for.

