

Name: - Harsh Patel

A20369913

CS 429 Project Report

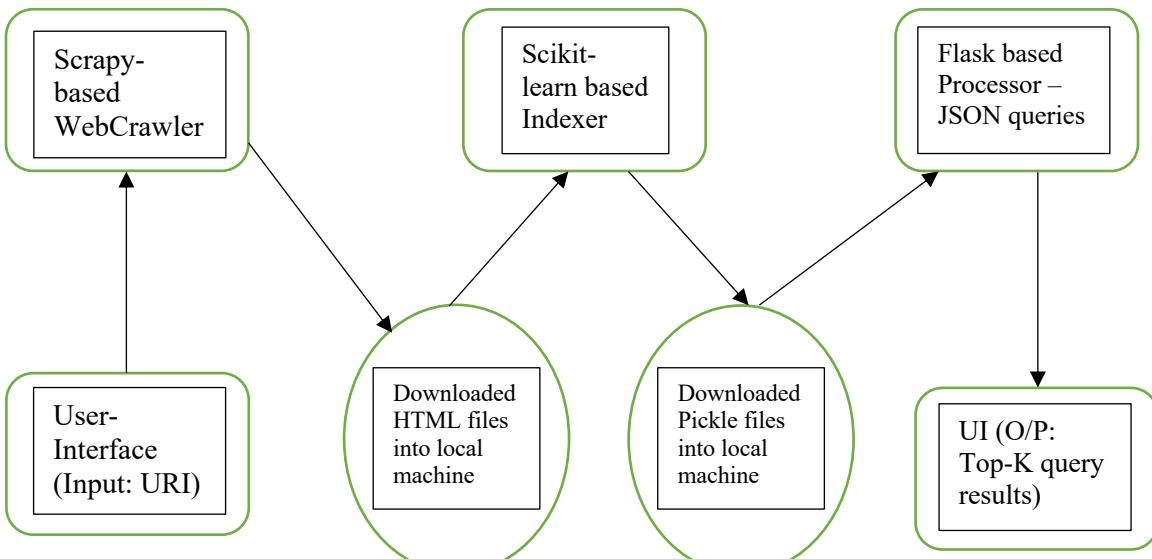
hpatel100@hawk.iit.edu

1. Abstract: -

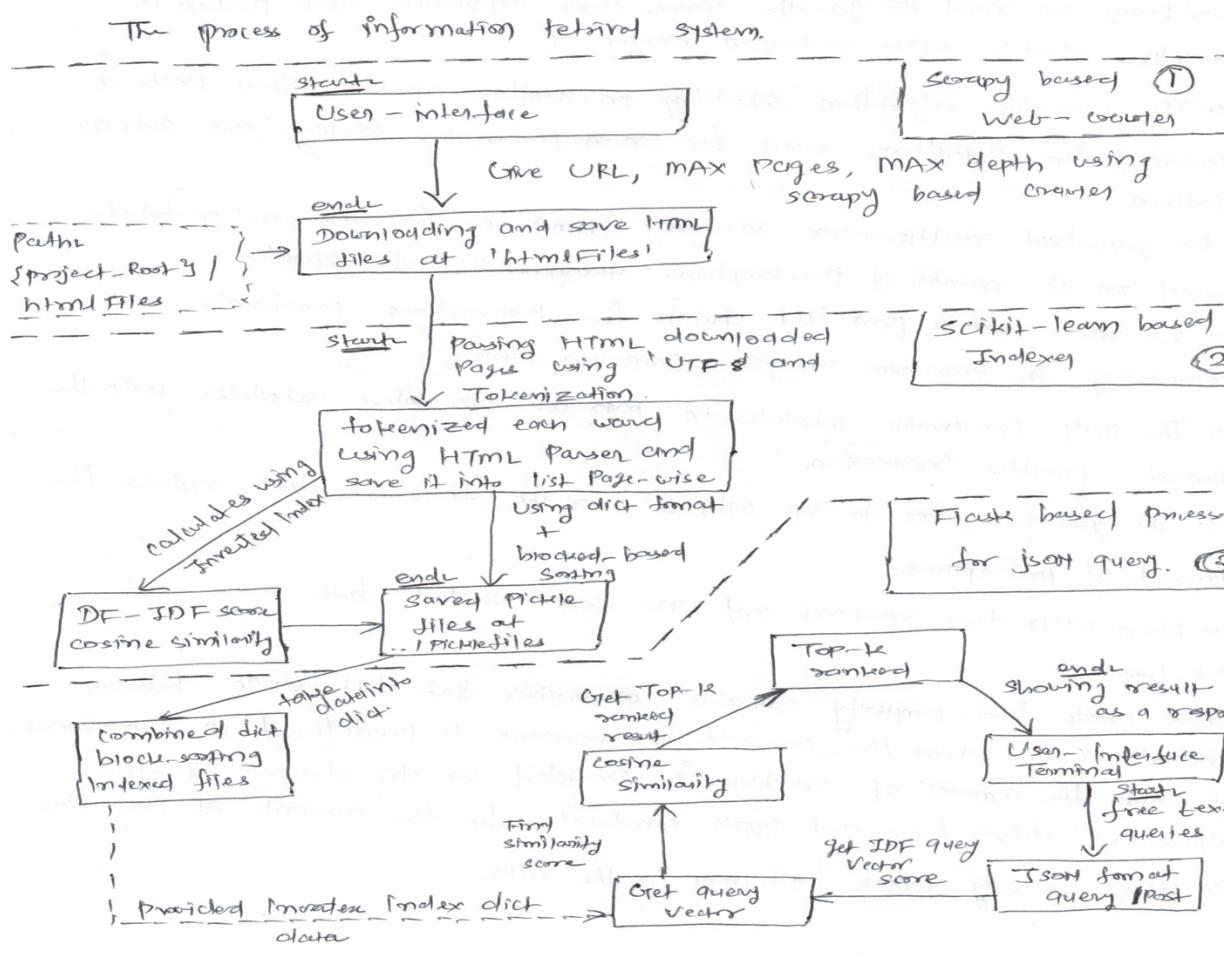
⇒ Finding a Top-k results from User defined query from mapped(user-defined) seeds crawler is the core of this project. A closely related problem refers to searching some specific questions on google, yahoo, Bing, etc. and get back response on Top-K query results that matches with millions of webpages or URLs. We study the local machine information retrieval system that based on scrapy crawling system, saved an inverted index of {word, Pages that includes word} into local machine pickles, and user-defined query in Flask based JSON processor format. We give a detailed analysis and show formally how all these three tasks are Information retrieval in general.

2. Overview: -

- ⇒ Information retrieval (IR) system in computing and information science is the process of obtaining in resources that are relevant to information need from a collection of resources like database, web crawling and information resources, and also can be metadata like described web data, or texts, sounds, or images from resources. Searches can be based on full-text or other content-based indexing.
- ⇒ We are using “Introduction to Information Retrieval” as a reference for building this information retrieval system.
- ⇒ We are building some basic information retrieval system using scrapy and flask library. We analyzed how the scrapy based content crawler for downloading web documents in html format into local machine using Max Pages and Max depth of seed URL/Domain, a scikit-based learn based indexer for constructing an inverted index in pickle format into local machine, and Flask based processor for handling free text queries in JSON format and get Top-K results from html inverted index based on downloaded HTML pages.
- ⇒ The proposed architecture is given below. ‘Loc’ is i.e. location in project directory.
- ⇒ Scrapy-based crawler: - (Loc: '{project_root_directory}/spiders/mainspider.py')
 - We have a user-interface which needs URL, Max Pages, and Max depth as a user-input that can map max depth and max pages to crawl URL seeds pages, using scrapy and saved downloaded HTML pages into local machine at ‘htmlFiles’ directory using renaming it using their domain name. For safe references, and test-cases, I am copying it to ‘htmlFilesToTest’ directory too.
- ⇒ Scikit-Learn based Indexer: - (Loc: '{project_root_directory}/indexer/main.py')
 - We then take all downloaded saved HTML files into consideration and parsed it with HTML parser. And, using block-sorting based indexing, I saved an inverted index of {word, Pages that includes word} as pickle files to the ‘pickleFiles’ directory. We find DF-IDF score and cosine similarity in this block too.
- ⇒ Flask based Processor: - (Loc: '{project_root_directory}/processor/app.py')
 - We then take all saved pickle files and merged/combined for the final one dictionary. I did that for inverted_index, tf_idf, idf pickle files. The flask app is running and identify queries are in valid JSON format or not, try to get query vector and then cosine similarity score. And then we can map Top-K based search result for queries.
- ⇒ Basic Block-Diagram of the system: -



Detailed version: -



3. Design: -

Project Directories structure: -

1. cs429crawler	-> Root Directory
a. cs429crawler	-> Scrapy based Crawler directory
i. spiders	
1. __init__.py	-> spider init class
2. mainspider.py	-> Main python class of spider
3. test.py	-> some test cases.
ii. __init__.py	
iii. middlewares.py	-> spider Middleware class
iv. pipelines.py	-> Spider pipeline class
v. settings.py	-> Spider settings class
b. htmlFiles	-> Spider HTML downloaded files
c. htmlFilesToTest	-> Backup of htmlFiles directory.
d. indexer	-> Scikit learn based indexer
i. __init__.py	-> Indexer init class
ii. main.py	-> Index operation python class
e. pickleFiles	-> Saved pickle files from main.py
f. processor	-> Flask based processor
i. template	-> Template class
1. __init__.py	-> Init class of template
2. index.html	-> html page of Flask
ii. __init__.py	
iii. app.py	-> Main python class of Flask app.
iv. config.py	-> Configuration class
v. logger.py	-> logger class
g. iitEdu.csv	-> Use for HTML parsing (UTF-8)
h. iitEdu_1.csv	-> Use for HTML parsing (UTF-8)
i. scrapy.cfg	-> scrapy config

2. As we have seen here, “cs429crawler” contains the main root directory. We run this program by using command-line using terminal input commands.
3. “cs429crawler/ cs429crawler” is the root directory of the scrapy based crawler, in which we run it by using *mainspider.py* class, and downloaded html files will be get saved into “cs429crawler/ htmlFiles” directory. Middleware.py, pipelines.py and settings.py are the scrapy architectural configuration files.
4. “cs429crawler/ indexer” is the root directory of the Scikit-Learn based Indexer, in which we run it by using *main.py* class, and downloaded inverted index directories pickle files will be get saved into “cs429crawler/ pickleFiles” directory. We simply use block-sorting based algorithm.
5. “cs429crawler/processor” is the root directory of the Flask based Processor, in which we run it by using *app.py* (*our main Flask app main class*), we try to post JSON queries from command-line input and get response as a query result, first it will validate the query, and then get the top-k results. *index.html* is template, and *config.py* is config class.

4 Architecture:-

- ⇒ As we discussed, we have 3 components in this IR system.
- ⇒ Right now, in Scrapy based Crawling in mainspider.py class, we include allowed domains are ['en.wikipedia.org', 'www.iit.edu', 'www.phdessay.com']. We can add any domain under it. For test purposes, I tested scrapy based crawling on base URL of 'en.wikipedia.org'. The spider class name is "MainSpider".
- ⇒ The snippet of the code is given below. It also included the 'DEPTH_LIMIT': '2'. We can also do it by passing it. 'DEPTH_LIMIT' from command line or pre-defined in 'cs429crawler/settings.py'.

```
10  custom_settings = {
11      'DEPTH_LIMIT': '2',
12  }
13
14  def __init__(self, url='', **kwargs):
15      self.links = []
16      self.allowed_domains = ['en.wikipedia.org', 'www.iit.edu', 'www.phdessay.com']
17      self.start_urls = [url] # py36
18      self.base_url = 'en.wikipedia.org'
19  super(MainSpider, self).__init__(**kwargs) # python3
```

- ⇒ Else you can set limit of 'DEPTH LIMIT' in 'cs429crawler/settings.py' as well like below.

```
BOT_NAME = 'cs429crawler'

SPIDER_MODULES = ['cs429crawler.spiders']
NEWSPIDER_MODULE = 'cs429crawler.spiders'

DEPTH_LIMIT = 2
SCHEDULER_DISK_QUEUE = 'scrapy.squeues.PickleFifoDiskQueue'
SCHEDULER_MEMORY_QUEUE = 'scrapy.squeues.FifoMemoryQueue'

SPIDER_MIDDLEWARES = {'cs429crawler.middlewares.Cs429CrawlerSpiderMiddleware': 0}

# FEED_FORMAT="csv"
# FEED_URI="iitEdu.csv"

# Crawl responsibly by identifying yourself (and your website) on the user-agent
#USER_AGENT = 'cs429crawler (+http://www.yourdomain.com)'

# Obey robots.txt rules
ROBOTSTXT_OBEY = True
```

- ⇒ The following snippet shows how every time we scrawl new URL checking "status code == 200", old files get deleted and new files pops up into 'htmlFiles' directory.

```
39  if requests.get('https://{}'.format(filePath)).status_code == 200:
40      for f in os.listdir(save_path):
41          if init == 0:
42              if not f.endswith(".html"):
43                  continue
44              os.remove(os.path.join(save_path, f))
45          if os.path.join('htmlFiles', filePath_upd):
46              with open(os.path.join('htmlFiles', filePath_upd), 'w'):
47                  pass
48              with open(os.path.join('htmlFiles', filePath_upd), "wb") as file:
49                  # response = get(self.base_url + link.xpath('.//@href').get())
50                  response = requests.get('https://{}'.format(filePath))
51                  file.write(response.content)
```

⇒ We also set depth limit into spider Middleware.py class by using filter function and set MAX_DEPTH as requested from response from requested fetched URL's meta.

```

48     def process_spider_output(self, response, result, spider):
49         def _filter(request):
50             if isinstance(request, Request):
51                 depth = response.meta['depth'] + 1
52                 request.meta['depth'] = depth
53                 if self.prio:
54                     request.priority -= depth * self.prio
55                 if self.maxdepth and depth > self.maxdepth:
56                     logging.debug(
57                         "Ignoring link (depth > %(maxdepth)d): %(requrl)s",
58                         {'maxdepth': self.maxdepth, 'requrl': request.url},
59                         extra={'spider': spider})
60             return False
61         else:
62             if self.verbose_stats:
63                 self.stats.inc_value(f'request_depth_count/{depth}', spider=spider)
64                 self.stats.max_value('request_depth_max', depth, spider=spider)
65             return True
66
67     # base case (depth=0)
68     if 'depth' not in response.meta:
69         response.meta['depth'] = 0
70         if self.verbose_stats:
71             self.stats.inc_value('request_depth_count/0', spider=spider)
72
73     return (r for r in result or () if _filter(r))
74
75
76

```

⇒ When we parse user argument given webpage, we try to parse it by, using to identify href link for next depth pages via “link.xpath(‘//@href’).get()” function. We also try to get html webpage name from the URL link that we got from extraction of URL links and saved them into ‘htmlFiles’ directory.

```

23     def parse(self, response, dest_path = save_path, save_path=save_path):
24         init = 0
25         for link in response.xpath('//div/p/a'):
26             # if init == 2:
27             #     break
28             yield {
29                 "link": self.base_url + link.xpath('.//@href').get()
30             }
31             print(self.base_url + link.xpath('.//@href').get())
32             filePath = self.base_url + link.xpath('.//@href').get()
33             words = ['tel', 'phone', 'mob', 'mailto', 'Special']
34             filePath_upd = filePath.rsplit('/', 1)[-1] + '.html'
35             for word in words:
36                 if word in filePath:
37                     filePath = self.base_url
38             print(requests.get('https://{}{}'.format(filePath)))

```

⇒ As a result, All URL links we got from scrapy, downloaded as HTML pages into “htmlFiles” directory automatically deleting last saved files. (in “htmlFiles” directory)

⇒ Now, to start Scikit-learn based Indexer process, we get all the files ends with “.html” extension in the htmlFiles directory and parse with HTML parser encoding with “utf-8”.

```

34     if filename.endswith(".html"):
35         filePath = ROOT_PATH + "/htmlFiles/" + filename
36         # print(filePath)
37         HtmlFile = open(filePath, 'r', encoding='utf-8')
38         source_code = HtmlFile.read()
39         soup = BeautifulSoup(source_code, "html.parser")

```

⇒ Then, we create a tokenized wordlist from each document and add them into list of html parsing list, index wise, by removing stop-words, in “indexer/main.py” file.

```

46     wordList = re.sub("[^\\w]", " ", str(soup.text)).split()
47     stops = ["", "-", "i", "me", "my", "myself", "we", "our", "ours", "ourselves", "you", "your", "yours", "yourself", "yourselves", "he", "him", "his", "hi
48     tokenized_word = [i for i in wordList if str(i).lower() not in stops]
49     # print(tokenized_word)
50     # print("length of tokenized_word= ", len(tokenized_word))
51     html_parsing_list.append(tokenized_word)
52
53

```

⇒ Then, we create the DF dictionary using html_parsing_list and tokenized_words.

```

55     DF = {}
56     upd_DF = {}
57     upd_upd_DF = {}
58     for w in tokenized_word:
59         if w in DF:
60             DF[w].append(w)
61         else:
62             DF[w] = [w]
63     for i in DF:
64         upd_DF[i] = len(DF.get(i))
65     for i in DF:
66         upd_upd_DF[i] = upd_DF.get(i) / len(tokenized_word)
67     df_list.append(upd_upd_DF)
68

```

⇒ We the created inverted index dictionary and save them into pickle files in “pickleFiles”.

⇒ e.g. like, using block-sorted index pickle files.

- save 0.pkl, save 1.pkl, save 2.pkl, etc.

```

76     # print("Number of unique words is present : ")
77     for item in html_parsing_list[init]:
78         if item in check:
79             if item not in dict:
80                 dict[item] = []
81             if item in dict and (init + 1) not in dict.get(item):
82                 dict[item].append(init + 1)
83
84             # Convert Index pickling of dictionary index pages pickle files.
85             fileName = "save_" + str(init) + ".pkl"
86             if not os.path.isfile(os.path.join('picklefiles', fileName)):
87                 with open(os.path.join('picklefiles', fileName), 'wb') as f:
88                     pickle.dump(dict, f, protocol=pickle.HIGHEST_PROTOCOL)
89             else:
90                 os.remove(os.path.join('picklefiles', fileName))
91                 with open(os.path.join('picklefiles', fileName), 'wb') as f:
92                     pickle.dump(html_parsing_list, f, protocol=pickle.HIGHEST_PROTOCOL)
93
94             init += 1

```

⇒ We then get IDF score from html_parsing_list and save them into “pickleFiles”.

```
111     # IDF_score logic
112     for key, values in dict.items():
113         idf_score = np.log(len(html_parsing_list) / len(values))
114         if key not in idf_dict:
115             idf_dict[key] = list()
116             for value in values:
117                 idf_dict[key].append([value, idf_score])
118             print("idf_dict:= ", idf_dict["Binary"])
119
120     # Pickle file - IDF
121     if not os.path.isfile('picklefiles/idf.pkl'):
122         with open('picklefiles/idf.pkl', 'wb') as f:
123             pickle.dump(idf_dict, f, protocol=pickle.HIGHEST_PROTOCOL)
124     else:
125         os.remove('picklefiles/idf.pkl')
126         with open('picklefiles/idf.pkl', 'wb') as f:
127             pickle.dump(idf_dict, f, protocol=pickle.HIGHEST_PROTOCOL)
```

⇒ We then try to get TF-IDF score/weight representation, save it into tf_idf_dict dictionary.

```
129     # Find TF-IDF score
130     init_token = 0
131     for token_item in range(len(html_parsing_list)):
132         # if init_token == 1:
133         #     break
134         for item in df_list[token_item]:
135             if item not in tf_idf_dict:
136                 tf_idf_dict[item] = {}
137             if item in tf_idf_dict:
138                 tf_idf_dict_opt = {}
139                 doc_ids = [docs[0] for docs in idf_dict[item]]
140                 doc_ids_idf_score = [docs[1] for docs in idf_dict[item]]
141                 # print("item:= ", item)
142                 # print("doc_ids:= ", doc_ids)
143                 # print("doc_ids_idf_score:= ", doc_ids_idf_score)
144                 count = 0
145                 for number in doc_ids:
146                     if number not in tf_idf_dict_opt:
147                         tf_idf_dict_opt[number] = []
148                     if number not in tf_idf_dict_opt_ext:
149                         tf_idf_dict_opt_ext[number] = []
150                     if number in tf_idf_dict_opt:
151                         # print("doc_ids_idf_score[count]:= ", doc_ids_idf_score[count])
152                         # print("df_list[number - 1][item]:= ", df_list[number - 1][item])
153                         tf_idf_dict_opt[number].append(doc_ids_idf_score[count] * df_list[number - 1][item])
154                         tf_idf_dict_opt_ext[number].append(doc_ids_idf_score[count] * df_list[number - 1][item])
155                         count = count + 1
156             tf_idf_dict[item] = tf_idf_dict_opt
157             init_token += 1
158             #elementsArticles
159             # print("tf_idf_dict_opt_ext:= ", tf_idf_dict_opt_ext)
160             print(tf_idf_dict["Binary"])
161             print("tf_idf_dict_opt_ext:= ", len(tf_idf_dict_opt_ext))
```

⇒ We then use IDF score, TF-IDF score, we try to implement query vector score using idf_dict dictionary and tokenization into queries.

```

181     # The vector should be a dict with query terms as keys and IDF scores as values.
182     query_to_check = "Binary Electronic"
183     data = re.sub('[^\w|_]', ' ', query_to_check) # only keep numbers and letters and spaces
184     data = data.lower()
185     data = re.sub(r'[\x00-\x7f]', r'', data) # remove non ascii texts
186     tokens = [data_upt for data_upt in data.split(' ') if data_upt] # split the words and remove empty words
187     for token in tokens:
188         doc_ids_idf_score = [idf_dict[token][0][1]]
189         query_vector_dict[token] = doc_ids_idf_score
190
191     print(query_vector_dict)

```

⇒ After that, we get query-vector score, we get the cosine similarity score and even sorted similarity scores using the following snippet of code.

```

192     # Cosine similarity
193     # The result of the dot product between the query vector and the documents (per term) need to be converted into a final sorted list of pairs containing the document identifier and the score.
194     # The sorting is based on highest score to lowest.
195     terms = [docs for docs in query_vector_dict]
196     terms_idf_val = [query_vector_dict.get(docs)[0] for docs in query_vector_dict]
197     for term_num in range(len(terms)):
198         doc_length = []
199         cosine_similarity_dict_clone_res = {}
200         dict_keys = list(tf_idf_dict[terms[term_num]].keys())
201         cosine_similarity_dict = {}
202         wt = terms_idf_val[term_num]
203         for idx in range(len(html_parsing_list)):
204             if idx not in dict_keys:
205                 cosine_similarity_dict[idx] = 0
206             else:
207                 cosine_similarity_dict[idx] = terms_idf_val[0] * tf_idf_dict[terms[term_num]].get(idx)[0]
208                 doc_length.append(tf_idf_dict_opt_ext.get(idx)[0])
209         doc_square = [i ** 2 for i in doc_length]
210         doc_length_sqrt_val = pow(sum(doc_square), 0.5)
211         for key, values in cosine_similarity_dict.items():
212             cosine_similarity_dict_clone_res[key] = cosine_similarity_dict.get(key) / doc_length_sqrt_val
213         for key, values in cosine_similarity_dict_clone_res.items():
214             if key not in cosine_similarity_dict_parent and key != 0:
215                 cosine_similarity_dict_parent[key] = []
216             if key in cosine_similarity_dict_parent:
217                 cosine_similarity_dict_parent[key].append(values)
218             for key, values in cosine_similarity_dict_parent.items():
219                 if key not in cosine_similarity_dict_parent:
220                     cosine_similarity_dict_parent_2[key] = 0
221                 if key in cosine_similarity_dict_parent:
222                     cosine_similarity_dict_parent_2[key] = sum(cosine_similarity_dict_parent.get(key))
223                     cosine_similarity_dict_parent_term_list[terms[term_num]] = cosine_similarity_dict_clone_res
224         print("cosine_similarity_dict_parent_term_list:= ", cosine_similarity_dict_parent_term_list)
225         cosine_similarity_dict_parent_final = sorted(cosine_similarity_dict_parent_2.items(), key=lambda t: t[1],
226                                                       reverse=True)
227
228         # "document identifier, score" for query and document pages for combine query vector and document vector.
229         print(cosine_similarity_dict_parent_final)

```

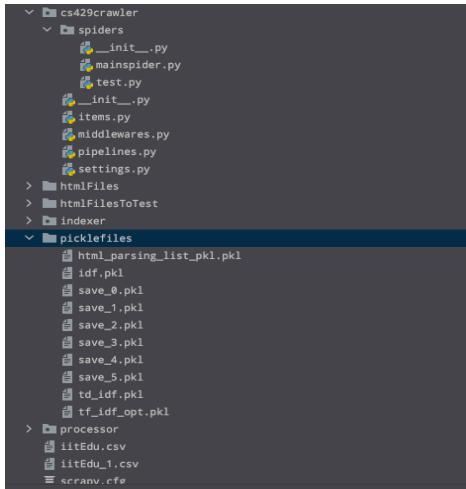
Result of Scikit-Learn based Indexer (TF-IDF/Weight, Cosine similarity): -

```

Terminal: Local : Local (2) x + v
(PycharmProjects) harshpatel@Harshs-MacBook-Pro cs429crawler % python indexer/main.py
dict score of term 'Binary': [2, 6]
idf_dict score of term 'Binary': [[2, 1.0986122886681098], [6, 1.0986122886681098]]
tf_idf_dict score of term 'Binary': [2: [0.0007851426756534723], 6: [0.0001053739315367428]]
tf_idf_dict_opt_ext score:= 6
sample (query vector score): {'binary': [0.6931471885599453], 'electronic': [0.6931471885599453]}
cosine_similarity_dict_parent_term_list:= {'binary': {0: 0.0, 1: 0.070108823467774237, 2: 0.09275166451717752, 3: 0.0, 4: 0.0, 5: 0.03208416814337556}, 'electronic': {0: 0.0, 1: 0.0, 2: 0.5618946278438862, 3: 0.0, 4: 0.0, 5: 0.728878633823944}}
[(0, 0.7689628811757698), (2, 0.65446462923610638), (1, 0.070108823467774237), (3, 0.0), (4, 0.0)]
(PycharmProjects) harshpatel@Harshs-MacBook-Pro cs429crawler %

```

⇒ We got the result of pickle files saved into “pickleFiles” directory. The snapshot is below.



⇒ Now, after that I work on Flask based processor. I have app.py as a main app in “processor” directory and config.py as a configuration file, and logger.py as debug logger file.

⇒ The configuration of flask processor is in “processor/config.py” file. We defined log_level, log_format, application_root, flask_run_mode, and web_port in config.py file.

```
1 import os
2
3 LOG_LEVEL = 'DEBUG'
4 LOG_FORMAT = '%(asctime)s %(levelname)s : %(message)s'
5 APPLICATION_ROOT = '/cs429crawler'
6 FLASK_RUN_MODE = os.environ.get('MODE') or 'PROD'
7 WEB_PORT = os.environ.get('PORT') or 5000
```

⇒ I run flask.app by using app.run(debug=True) in if __name__ == '__main__': class, and app = Flask(__name__).

⇒ We got multiple queries defined in JSON Post processing by given below snippet, in which we check whether the post request has JSON query object or not, and put all the queries into list.

```
34     @app.route('/', methods=['POST'])
35     def process_json():
36         global cosine_similarity_dict_parent_final
37         content_type = request.headers.get('Content-Type')
38         if (content_type == 'application/json'):
39             query_list = []
40             json = request.json
41             # query_1 = json['query1']
42             # print("query_1:= ", query_1)
43             for idx in range(len(json)):
44                 init = idx + 1
45                 queryName = "query" + str(init)
46                 query_list.append(json[queryName])
47             print("Query list(from Json format):= ", query_list)
```

⇒ Query example: - (I used command-line to use it). So, far we have only Flask Post request model.

```
curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"Binary Electronic\"}" "127.0.0.1:5000/"
```

response:

```
[
  [
    [
      5,
      0.7609628011757698
    ],
    [
      2,
      0.6546462923610638
    ],
    [
      1,
      0.07010823467774237
    ],
    [
      3,
      0.0
    ],
    [
      4,
      0.0
    ]
]
```

⇒ We then using block-sorting algorithm, merged/combined together for the final one inverted index.

```
52 # Using block-sorting algorithm, merged/combined together for the final one inverted index.
53 myDir = Path('picklefiles/')
54 fileNames = [file.name for file in myDir.iterdir() if file.name.startswith('save')]
55 for file in fileNames:
56     with open(os.path.join(myDir, file), 'rb') as fileName:
57         dict_new_inverted_index = pickle.load(fileName)
58         for k, v in dict_new_inverted_index.items():
59             try:
60                 main_res_dict_inverted_index[k].extend(v)
61             except KeyError:
62                 main_res_dict_inverted_index[k] = v
63 res_main_res_dict = {key: list(set(value)) for key, value in main_res_dict_inverted_index.items()}
64 print("Inverted index length: = ", len(res_main_res_dict))
```

⇒ We then calculate TF-IDF, and then calculate the IDF score, TF-IDF score, we try to implement query vector score using idf_dict dictionary and tokenization into queries. We run loop through query list as we might have multiple queries. The image is given below of how to calculate query vector.

```
101 #
102 # # Query vector and similarity
103 for idx in range(len(query_list)):
104     query_to_check = query_list[idx]
105     data = re.sub('[\W]+', ' ', query_to_check) # only keep numbers and letters and spaces
106     data = data.lower()
107     data = re.sub(r'^[\x00-\x7f]', r'', data) # remove non ascii texts
108     tokens = [data_up for data_up in data.split(' ') if data_up] # split the words and remove empty words
109     for token in tokens:
110         doc_ids_idf_score = [idf_dict[token][0][1]]
111         query_vector_dict[token] = doc_ids_idf_score
112
```

⇒ And, from then, we need to calculate similarity score, in this case, I chose to do cosine similarity. The code snippet is given below.

```

113     # Cosine similarity
114     # The result of the dot product between the query vector and the documents (per term) need to be converted into a final sorted list of pairs containing t
115     # The sorting is based on highest score to lowest.
116     terms = [docs for docs in query_vector_dict]
117     terms_idf_val = [query_vector_dict.get(docs)[0] for docs in query_vector_dict]
118     for term_num in range(len(terms)):
119         doc_length = []
120         cosine_similarity_dict_clone_res = {}
121         dict_keys = list(tf_idf_dict[terms[term_num]].keys())
122         cosine_similarity_dict = {}
123         for idx in range(len(html_parsing_list)):
124             if idx not in dict_keys:
125                 cosine_similarity_dict[idx] = 0
126             else:
127                 cosine_similarity_dict[idx] = terms_idf_val[0] * tf_idf_dict[terms[term_num]].get(idx)[0]
128                 doc_length.append(tf_idf_dict_opt_ext.get(idx)[0])
129         doc_square = [i ** 2 for i in doc_length]
130         doc_length_sqrt_val = pow(sum(doc_square), 0.5)
131         for key, values in cosine_similarity_dict.items():
132             cosine_similarity_dict_clone_res[key] = cosine_similarity_dict.get(key) / doc_length_sqrt_val
133         for key, values in cosine_similarity_dict_clone_res.items():
134             if key not in cosine_similarity_dict_parent and key != 0:
135                 cosine_similarity_dict_parent[key] = []
136             if key in cosine_similarity_dict_parent:
137                 cosine_similarity_dict_parent[key].append(values)
138             for key, values in cosine_similarity_dict_parent.items():
139                 if key not in cosine_similarity_dict_parent:
140                     cosine_similarity_dict_parent_2[key] = 0
141                 if key in cosine_similarity_dict_parent:
142                     cosine_similarity_dict_parent_2[key] = sum(cosine_similarity_dict_parent.get(key))
143         cosine_similarity_dict_parent_term_list[terms[term_num]] = cosine_similarity_dict_clone_res
144         print("\n\n\nTop-K ranked results(Query word -> similarity result):= ", cosine_similarity_dict_parent_term_list)
145         cosine_similarity_dict_parent_final = sorted(cosine_similarity_dict_parent_2.items(), key=lambda t: t[1],
146                                                       reverse=True)
147         print("Top-K ranked results(Top-K Document Numbers):= ", cosine_similarity_dict_parent_final)
148         # return json
149         jsonify("The result is given below (HTML page saved and their K-top results using similarity")
150         return jsonify(cosine_similarity_dict_parent_final), 200
151     else:
152         return 'Content-Type not supported!'
153

```

And, the we get following results of Top-k: -

```

[(base) harshpatel@Harshs-MacBook-Pro ~ % curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"Binary Electronic\"}" "127.0.0.1:5000/"
[
  [
    5,
    0.7609628011757698
  ],
  [
    2,
    0.6546462923610638
  ],
  [
    1,
    0.07010823467774237
  ],
  [
    3,
    0.0
  ],
  [
    4,
    0.0
  ]
]

```

5 Operations: - (Steps to how to work on local machine)

- ⇒ Please extract the CS429_Harsh_Patel_Project.zip and extract it.
- ⇒ I am running it in PyCharm IDE PyCharm 2021.2.3 (Professional Edition), Build #PY-212.5457.59, built on October 19, 2021
- ⇒ To check, Python 3.0+/sklearn 1.0+/Scrapy 2.0+/Flask 2.0+,

Python: -

- ⇒ First check if python installed in your local machine or not using

```
Python --version
```

- ⇒ If not, use the following command:

```
brew install python
```

- ⇒ I have Python 3.9.12 in my local machine [<https://www.python.org/downloads/release/python-3912/>], released on March 23, 2022.
- ⇒ Again, check python is installed or not by using following commands,

```
Python --version
```

Scrapy: -

- ⇒ Scrapy requires Python 3.6+, either the CPython implementation (default) or the PyPy 7.2.0+ implementation.
- ⇒ First check if scrapy installed in your local machine or not using,

```
scrapy version
```

- ⇒ If not, use the following command:

```
pip install Scrapy
```

OR

```
conda install -c conda-forge scrapy
```

- ⇒ I have Scrapy 2.6.1 in my local machine [<https://www.python.org/downloads/release/python-3912/>], released on March 23, 2022.
- ⇒ Again, check scrapy is installed or not by using following commands,

```
scrapy version
```

Flask: -

⇒ First check if flask installed in your local machine or not using

```
flask --version
```

⇒ In my machine, I got flask version,

Python 3.9.12

Flask 2.0.3

Werkzeug 2.0.3

⇒ If not, use the following command:

```
pip install Flask
```

⇒ I have Python 3.9.12 in my local machine [<https://www.python.org/downloads/release/python-3912/>], released on March 23, 2022.

⇒ Again, check python is installed or not by using following commands,

```
flask --version
```

After checking these basic requirements,

Check BeautifulSoup installed or not, and if not, use this command, I have beautifulsoup4 4.11.1

```
pip install beautifulsoup4
```

How to run command in this project: -

1. A Scrapy based Crawler for downloading web documents in html format, go to **project root directory ‘cs429crawler’** and run the following command.
2. Please type or give URL as an input argument in command prompt.

```
scrapy crawl MainSpider -a url={please give your URL/domain name}
```

3. Now, the downloaded HTML URLs will get stored in ‘htmlFiles’ directory.
4. Then, to run Scikit-Learn based Indexer, we use following command.
5. **project root directory: cs429crawler / package: indexer**

```
python indexer/main.py
```

6. Then, we got saved pickle files into “picklefiles” directory.
7. Now, to run Flask processor server, please use the following command.

```
python processor/app.py runserver
```

If you get the address already in use **ERROR:**

⇒ Type the following command:

```
ps -fA | grep python
```

and, if you found out, “python processor/app.py runserver” process, kill that process using the port number.

Like example, given below.

```
ps -fA | grep python

501 58091 46777 0 3:13PM ttys007 0:00.21 python processor/app.py
runserver

(PycharmProjects) harshpatel@Harshs-MacBook-Pro cs429crawler % kill -9 58091

(PycharmProjects) harshpatel@Harshs-MacBook-Pro cs429crawler %

[1] + killed python processor/app.py runserver
```

And, try again, you will get following result.

```
(PycharmProjects) harshpatel@Harshs-MacBook-Pro cs429crawler % python
processor/app.py runserver

* Serving Flask app 'app' (lazy loading)

* Environment: production

WARNING: This is a development server. Do not use it in a production
deployment.

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: 740-214-834
```

⇒ You will check on <http://127.0.0.1:5000/>

⇒ **Now, To get query, please type following curl JSON Post command in terminal: -**

```
curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"{\\"type your query\\\"}\"}" "127.0.0.1:5000/"
```

⇒ You will get response result like, **ex., docs 5, 2, 1 are relevant; 3 and 4 are not relevant.**
127.0.0.1 - - [30/Apr/2022 16:23:04] "POST / HTTP/1.1" 200 -

```
[  
  [  
    5,  
     0.7609628011757698  
  ],  
  [  
    2,  
     0.6546462923610638  
  ],  
  [  
    1,  
     0.07010823467774237  
  ],  
  [  
    3,  
     0.0  
  ],  
  [  
    4,  
     0.0  
  ]  
]
```

6 Conclusion: - (Steps to how to work on local machine)

Test case: -

1. Go into project root directory, type the following command.

```
scrapy crawl MainSpider -a url=https://en.wikipedia.org/wiki/George_Minaker
```

⇒ will be downloaded 34 HTML documents into ‘htmlFiles’ into my local disk.

2. Now, after we got html pages, we now start indexing process.

```
python indexer/main.py
```

⇒ will be downloaded 37 pickle files into ‘pickleFiles’ into local machines.

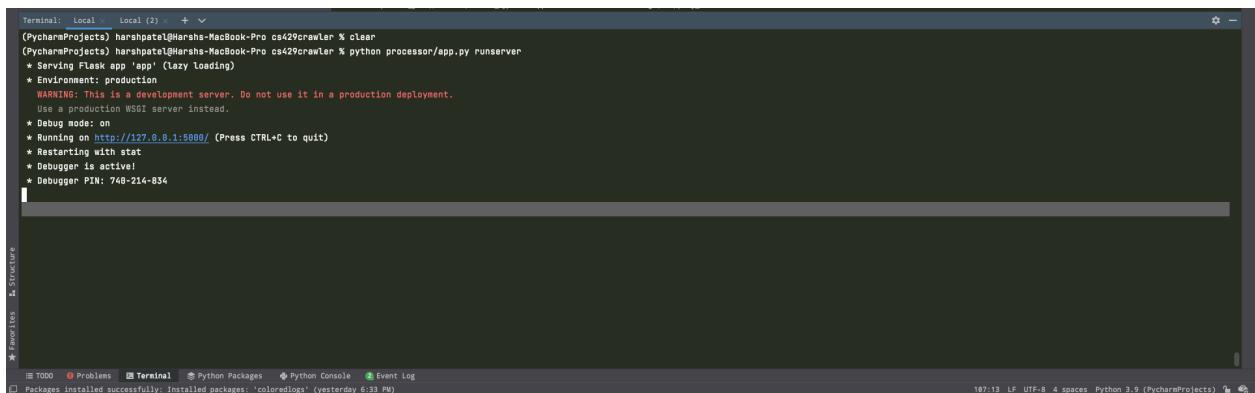
Result: - (TF-IDF score/weight representation, Cosine similarity)

```
[PycharmProjects] harshpatel@Harshs-MacBook-Pro:~/cs429crawler$ 
[PycharmProjects] harshpatel@Harshs-MacBook-Pro:~/cs429crawler$ clear
[PycharmProjects] harshpatel@Harshs-MacBook-Pro:~/cs429crawler$ python indexer/main.py
tf_idf_dict score of term 'Oshawa': [ (0.00824840248650851978, 7), (7. [7.945581045892178e-05], 18): [0.0086185014051916], 20): [0.00805807421480785524], 24): [0.00805285787209233684], 26): [0.00848317839913202316], 30): [0.008114194812867189221)]
tf_idf_dict score of term 'Vancouver': [ (0.0081205465542854982, 3), (8. [8.001286618849220115], 4): [0.00813617254786160483], 6): [6.225732832577231e-05], 7): [0.008329117821532935], 12): [8.948171935634428e-05], 15): [0.008038639413829536], 17): [9.262513248567553e-05], 18): [0.0086353224763778], 20): [0.00887785141616193], 24): [0.0088121647147338426], 26): [0.00815618278822467525], 28): [0.008089440878517156661], 30): [0.0082284957881640583], 32): [0.0085138294625175351]
sample Query vector score: {'Oshawa': [1.588483755688479], 'Vancouver': [8.183103235139513]}
cosine_similarity dict_parameter_term_list: {'Oshawa': [(0.8, 0.1, 0.8, 2.0, 0.8, 0.4, 0.8, 0.5, 0.8, 0.6, 0.1336245862612927, 7, 0.84414962823132146, 8, 0.8, 9.0, 0.8, 10, 0.8, 11, 0.8, 0.12, 0.8, 13, 0.8, 14, 0.8, 15, 0.8, 16, 0.8, 17, 0.8, 18, 0.8, 19, 0.8, 20, 0.8, 21.0, 0.8, 22, 0.8, 23, 0.8, 24, 0.8, 25, 0.8, 26, 0.8, 27, 0.8, 28, 0.8, 29, 0.8, 30, 0.8, 31, 0.8, 32, 0.8, 33, 0.8), 'Vancouver': [(0.8, 0.1, 0.8, 0.8, 0.2, 0.2575637677887657, 3, 4. [4.007612424888787386475, 7, 0.81543995924017414, 8, 0.8, 10, 0.8, 11, 0.8, 12, 0.8, 13, 0.8, 14, 0.8, 15, 0.8, 16, 0.8, 17, 0.8, 18, 0.8, 19, 0.8, 20, 0.8, 21, 0.8, 22, 0.8, 23, 0.8, 24, 0.8, 25, 0.8, 26, 0.8, 27, 0.8, 28, 0.8, 29, 0.8, 30, 0.8, 31, 0.8, 32, 0.8, 33, 0.8)]}
Sorted cosine similarity dict_parameter_final: [(26, 0.73294327892553), (4, 0.7661899783786475), (18, 0.6895988641782), (24, 0.6793215438126684), (2, 0.6515681788924783), (3, 0.6803122491887558), (28, 0.4018863591284689), (2, 0.2575363747788656), (7, 0.198649558643596288), (38, 0.18246464796880133), (15, 0.18143558133576343), (6, 0.18285791408463415), (32, 0.147361178966298), (17, 0.844731588227345394), (12, 0.84281686978161616), (1, 0.8), (5, 0.8), (8, 0.8), (9, 0.8), (11, 0.8), (13, 0.8), (15, 0.8), (17, 0.8), (19, 0.8), (21, 0.8), (22, 0.8), (23, 0.8), (25, 0.8), (27, 0.8), (29, 0.8), (31, 0.8), (33, 0.8)]
[PycharmProjects] harshpatel@Harshs-MacBook-Pro:~/cs429crawler$
```

3. Now, after we got html pages, we now start indexing process.

```
python processor/app.py runserver
```

Result: - (Flask server started)



4. If we send query over in Flask JSON Post server, result would be, 1st query is given below,

```
curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"Oshawa  
Vancouver\"}" "127.0.0.1:5000/"
```

1st Query Result would be: - (In response in pycharm terminal and in terminal)

⇒ NOTE: In Top-K results, it's in sorting decreasing order, 0 means no similarity[irrelevant], positive means relevant that much of similarity in decreasing order, better the score comes first.

```
* Debugger is active!
* Debugger PIN: 740-214-834
Query list[from Json format]:= [ 'Oshawa Vancouver' ]
Inverted index length: = 20395

Top-K ranked results(Query word -> similarity result):= { 'Oshawa': {0: 8.0, 1: 8.0, 2: 8.0, 3: 8.0, 4: 8.0, 5: 8.0, 6: 8.133624586211297, 7: 8.844169628231321486}, 8: 8.0, 9: 8.0, 10: 8.0, 11: 8.0, 12: 8.0, 13: 8.0, 14: 8.0, 15: 8.0, 16: 8.0, 17: 8.0, 18: 8.0, 19: 8.0, 20: 8.0, 21: 8.0, 22: 8.0, 23: 8.0, 24: 8.29378577899521949, 25: 8.0, 26: 8.2684745864176573, 27: 8.0, 28: 8.0, 29: 8.0, 30: 8.0, 31: 8.0, 32: 8.0, 33: 8.0}
Top-K ranked results(Top-K Document Numbers):= [(18, 8.3677471755231222), (24, 8.29378577899521949), (26, 8.2817339833517840864), (26, 8.2684745864176573), (6, 8.133624586211297), (38, 8.08787374385274977), (7, 8.044169628231321486), (1, 8.0), (2, 8.0), (3, 8.0), (4, 8.0), (5, 8.0), (8, 8.0), (10, 8.0), (11, 8.0), (12, 8.0), (13, 8.0), (14, 8.0), (15, 8.0), (16, 8.0), (17, 8.0), (19, 8.0), (21, 8.0), (22, 8.0), (23, 8.0), (25, 8.0), (27, 8.0), (28, 8.0), (29, 8.0), (31, 8.0), (32, 8.0), (33, 8.0)]

Top-K ranked results(Query word -> similarity result):= { 'Oshawa': {0: 8.0, 1: 8.0, 2: 8.0, 3: 8.0, 4: 8.0, 5: 8.0, 6: 8.133624586211297, 7: 8.844169628231321486}, 8: 8.0, 9: 8.0, 10: 8.0, 11: 8.0, 12: 8.0, 13: 8.0, 14: 8.0, 15: 8.0, 16: 8.0, 17: 8.0, 18: 8.0, 19: 8.0, 20: 8.0, 21: 8.0, 22: 8.0, 23: 8.0, 24: 8.29378577899521949, 25: 8.0, 26: 8.2684745864176573, 27: 8.0, 28: 8.0, 29: 8.0, 30: 8.0, 31: 8.0, 32: 8.0, 33: 8.0}, 'Vancouver': {0: 8.0, 1: 8.0, 2: 8.0, 3: 8.0, 4: 8.0, 5: 8.0, 6: 8.0, 7: 8.0, 8: 8.0, 9: 8.0, 10: 8.0, 11: 8.0, 12: 8.0, 13: 8.0, 14: 8.0, 15: 8.0, 16: 8.0, 17: 8.0, 18: 8.0, 19: 8.0, 20: 8.0, 21: 8.0, 22: 8.0, 23: 8.0, 24: 8.0, 25: 8.0, 26: 8.0, 27: 8.0, 28: 8.0, 29: 8.0, 30: 8.0, 31: 8.0, 32: 8.0, 33: 8.0}, 'Vancouver': {0: 8.0, 1: 8.0, 2: 8.0, 3: 8.0, 4: 8.0, 5: 8.0, 6: 8.0, 7: 8.0, 8: 8.0, 9: 8.0, 10: 8.0, 11: 8.0, 12: 8.0, 13: 8.0, 14: 8.0, 15: 8.0, 16: 8.0, 17: 8.0, 18: 8.0, 19: 8.0, 20: 8.0, 21: 8.0, 22: 8.0, 23: 8.0, 24: 8.0, 25: 8.0, 26: 8.0, 27: 8.0, 28: 8.0, 29: 8.0, 30: 8.0, 31: 8.0, 32: 8.0, 33: 8.0}
Top-K ranked results(Top-K Document Numbers):= [(26, 8.9732943748292553), (4, 8.765188978376475), (18, 8.689955988441782), (24, 8.6792351430124684), (28, 8.65156931788924783), (3, 8.68132124288878586), (28, 8.4918643912844689), (2, 8.2957355787864956), (7, 8.1986995883356828), (8, 8.198246617476800133), (15, 8.1843558337854543), (6, 8.16285794108464545), (32, 8.1473611789456298), (17, 8.047451580822743594), (12, 8.0401686778161615), (1, 8.0), (5, 8.0), (8, 8.0), (9, 8.0), (10, 8.0), (11, 8.0), (14, 8.0), (16, 8.0), (19, 8.0), (21, 8.0), (22, 8.0), (23, 8.0), (25, 8.0), (27, 8.0), (29, 8.0), (31, 8.0), (35, 8.0)]
127.0.0.1 - - [30/Jan/2022 17:09:08] "POST /HTTP/1.1" 200 -
```

```
[base] harshpatel@Harshs-MacBook-Pro ~ % curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"Oshawa Vancouver\"}" 127.0.0.1:5000/
[{"lat": 43.7661899783706475, "lon": -79.689559800441782}, {"lat": 43.6792351430126694, "lon": -79.6515601788924783}, {"lat": 43.6013212428057586, "lon": -79.4010063391284489}, {"lat": 43.25753637677886956, "lon": -79.1986895863386288}, {"lat": 43.18240641976000133, "lon": -79.18143558133576343}, {"lat": 43.16285794160463415, "lon": -79.1473611789656298}, {"lat": 43.0644731580227343394, "lon": -79.04281686978161616} ]
```

2nd Query: -

```
curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"Oshawa\", \"query2\" : \"Vancouver\"}" "127.0.0.1:5000/"
```

2nd Query Result would be: - (In response in pycharm terminal and in terminal)

```

Top-K ranked results(Query word -> similarity result):= {'Oshawa': {0: 0.0, 1: 0.0, 2: 0.0, 3: 0.0, 4: 0.0, 5: 0.0, 6: 0.1336245862612927, 7: 0.844149628231321486, 8: 0.0, 9: 0.0, 10: 0.0, 11: 0.0, 12: 0.0, 13: 0.0, 14: 0.0, 15: 0.0, 16: 0.0, 17: 0.0, 18: 0.3677471755231222, 19: 0.0, 20: 0.28173893331784864, 21: 0.0, 22: 0.0, 23: 0.0, 24: 0.2937857789521949, 25: 0.0, 26: 0.26847455864176573, 27: 0.0, 28: 0.0, 29: 0.0, 30: 0.0, 31: 0.0, 32: 0.0, 33: 0.0}
Top-K ranked results(Top-K Document Numbers):= [(18, 0.3677471755231222), (24, 0.2937857789521949), (28, 0.28173893331784864), (26, 0.26847455864176573), (6, 0.1336245862612927), (50, 0.87887374385274877), (7, 0.844149628231321486), (1, 0.0), (2, 0.0), (3, 0.0), (4, 0.0), (5, 0.0), (8, 0.0), (9, 0.0), (10, 0.0), (11, 0.0), (12, 0.0), (13, 0.0), (14, 0.0), (15, 0.0), (16, 0.0), (17, 0.0), (19, 0.0), (21, 0.0), (22, 0.0), (23, 0.0), (25, 0.0), (27, 0.0), (28, 0.0), (29, 0.0), (31, 0.0), (32, 0.0), (33, 0.0)]

Top-K ranked results(Query word -> similarity result):= {'Oshawa': {0: 0.0, 1: 0.0, 2: 0.0, 3: 0.0, 4: 0.0, 5: 0.0, 6: 0.1336245862612927, 7: 0.844149628231321486, 8: 0.0, 9: 0.0, 10: 0.0, 11: 0.0, 12: 0.0, 13: 0.0, 14: 0.0, 15: 0.0, 16: 0.0, 17: 0.0, 18: 0.3677471755231222, 19: 0.0, 20: 0.28173893331784864, 21: 0.0, 22: 0.0, 23: 0.0, 24: 0.2937857789521949, 25: 0.0, 26: 0.26847455864176573, 27: 0.0, 28: 0.0, 29: 0.0, 30: 0.0, 31: 0.0, 32: 0.0, 33: 0.0}}
Top-K ranked results(Top-K Document Numbers):= [(18, 0.735494351642444), (24, 0.587415419858298), (28, 0.5634778665340813), (26, 0.5369491012835515), (6, 0.2672490125228584), (50, 0.15774748610548153), (7, 0.88829925646264297), (1, 0.0), (2, 0.0), (3, 0.0), (4, 0.0), (5, 0.0), (8, 0.0), (9, 0.0), (10, 0.0), (11, 0.0), (12, 0.0), (13, 0.0), (14, 0.0), (15, 0.0), (16, 0.0), (17, 0.0), (19, 0.0), (21, 0.0), (22, 0.0), (23, 0.0), (25, 0.0), (27, 0.0), (28, 0.0), (29, 0.0), (31, 0.0), (32, 0.0), (33, 0.0)]

Top-K ranked results(Query word -> similarity result):= {'Oshawa': {0: 0.0, 1: 0.0, 2: 0.0, 3: 0.0, 4: 0.0, 5: 0.0, 6: 0.1336245862612927, 7: 0.844149628231321486, 8: 0.0, 9: 0.0, 10: 0.0, 11: 0.0, 12: 0.0, 13: 0.0, 14: 0.0, 15: 0.0, 16: 0.0, 17: 0.0, 18: 0.3677471755231222, 19: 0.0, 20: 0.28173893331784864, 21: 0.0, 22: 0.0, 23: 0.0, 24: 0.2937857789521949, 25: 0.0, 26: 0.26847455864176573, 27: 0.0, 28: 0.0, 29: 0.0, 30: 0.0, 31: 0.0, 32: 0.0, 33: 0.0}, 'Vancouver': {0: 0.0, 1: 0.0, 2: 0.0, 3: 0.0, 4: 0.0, 5: 0.0, 6: 0.02923343534334143, 7: 0.15453995240171414, 8: 0.0, 9: 0.0, 10: 0.0, 11: 0.0, 12: 0.0241686978161616, 13: 0.0, 14: 0.0, 15: 0.0, 16: 0.0, 17: 0.0, 18: 0.044731588227343394, 19: 0.0, 20: 0.369821245574377, 21: 0.0, 22: 0.0, 23: 0.0, 24: 0.3855293720597456, 25: 0.0, 26: 0.7848198241874896, 27: 0.0, 28: 0.84018863391284489, 29: 0.0, 30: 0.183532678726056, 31: 0.0, 32: 0.0, 33: 0.0}}
Top-K ranked results(Top-K Document Numbers):= [(26, 0.12471689254718211), (18, 0.875366975964984), (24, 0.9729489139655754), (28, 0.933299112289519), (4, 0.7661899783786475), (3, 0.6813212428857586), (28, 0.84018863391284489), (6, 0.296482478652687), (50, 0.26122881628127421), (2, 0.25753437677886956), (7, 0.24283928886438437), (15, 0.18143558133576343), (32, 0.1473611789656298), (17, 0.844751588227343394), (12, 0.8420168497816161), (1, 0.0), (5, 0.0), (8, 0.0), (9, 0.0), (10, 0.0), (11, 0.0), (13, 0.0), (14, 0.0), (16, 0.0), (19, 0.0), (21, 0.0), (22, 0.0), (23, 0.0), (25, 0.0), (27, 0.0), (29, 0.0), (31, 0.0), (35, 0.0)]
```

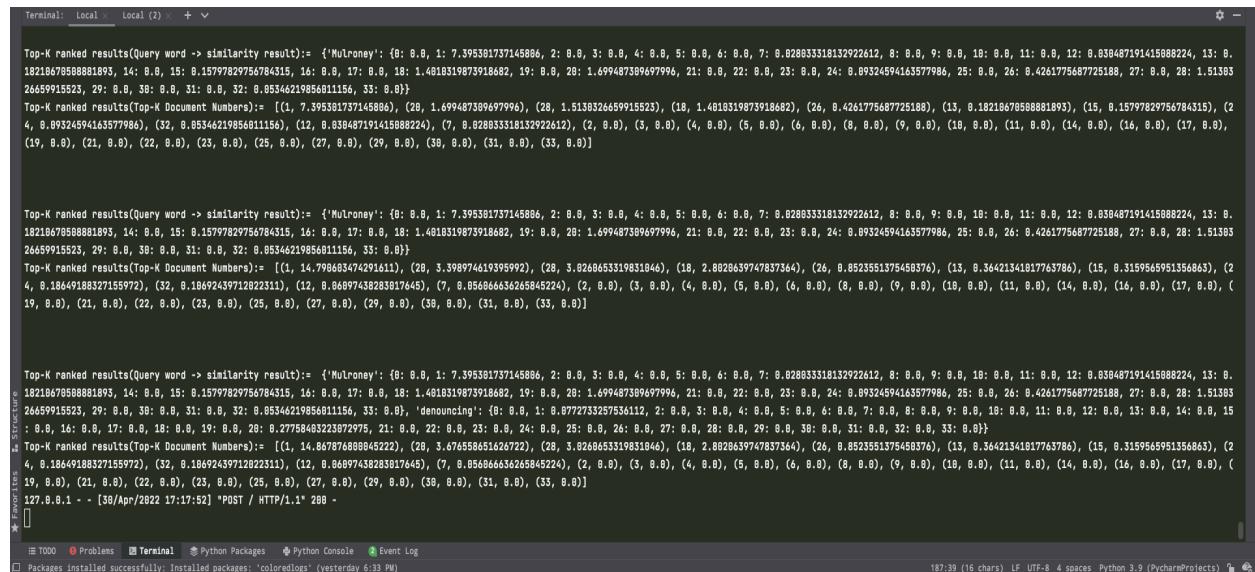
```
(base) harshpatel@Harsha-MacBook-Pro ~ % curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"Oshawa\", \"query2\" : \"Vancouver\"}" "127.0.0.1:5000"
[
  [
    [
      26,
      1.2417689254710211
    ],
    [
      18,
      1.057306975964984
    ],
    [
      24,
      0.9729409139655754
    ],
    [
      28,
      0.933299112209519
    ],
    [
      4,
      0.7661899783706475
    ],
    [
      3,
      0.6013212428057586
    ],
    [
      28,
      0.4010003391284489
    ],
    [
      6,
      0.29648244786592687
    ],
    [
      38,
      0.2612801628127421
    ],
    [
      2,
      0.25753637677806956
    ],
    [
      7,
      0.24283928886438437
    ],
    [
      15,
      0.18143558133576343
    ],
    [
      32,
      0.1473611789656298
    ],
    [
      17,
      0.044731580227343394
    ],
    [
      12,
      0.04201686978161616
    ],
    [

```

3rd Query:-

```
curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"Mulroney\", \"query2\" : \"denouncing\"}" "127.0.0.1:5000/"
```

3rd Query Result would be: - (In response in pycharm terminal and in terminal)



```
Terminal: Local Local (2) + -
Top-K ranked results(Query word -> similarity result):= {'Mulroney': {0: 0.0, 1: 7.395301737145886, 2: 0.0, 3: 0.0, 4: 0.0, 5: 0.0, 6: 0.0, 7: 0.028033318132922612, 8: 0.0, 9: 0.0, 10: 0.0, 11: 0.0, 12: 0.039487191415988224, 13: 0.18218670588881893, 14: 0.0, 15: 0.15797829756784315, 16: 0.0, 17: 0.0, 18: 1.4010319873918682, 19: 0.0, 20: 1.699487389697996, 21: 0.0, 22: 0.0, 23: 0.0, 24: 0.09324594165577986, 25: 0.0, 26: 0.4261775687725188, 27: 0.0, 28: 1.5130326659915523, 29: 0.0, 30: 0.0, 31: 0.0, 32: 0.05346219856811156, 33: 0.0}}
Top-K ranked results(Top-K Document Numbers):= [(1, 7.395301737145886), (28, 1.699487389697996), (28, 1.5130326659915523), (18, 1.4010319873918682), (26, 0.4261775687725188), (13, 0.18218670588881893), (15, 0.15797829756784315), (2, 0.049324594165577986), (32, 0.05346219856811156), (12, 0.058407191415888224), (7, 0.028833318132922612), (2, 0.0), (3, 0.0), (4, 0.0), (5, 0.0), (6, 0.0), (8, 0.0), (9, 0.0), (10, 0.0), (11, 0.0), (14, 0.0), (16, 0.0), (17, 0.0), (19, 0.0), (21, 0.0), (22, 0.0), (23, 0.0), (25, 0.0), (27, 0.0), (29, 0.0), (30, 0.0), (31, 0.0), (33, 0.0)]
```



```
Top-K ranked results(Query word -> similarity result):= {'Mulroney': {0: 0.0, 1: 7.395301737145886, 2: 0.0, 3: 0.0, 4: 0.0, 5: 0.0, 6: 0.0, 7: 0.028033318132922612, 8: 0.0, 9: 0.0, 10: 0.0, 11: 0.0, 12: 0.039487191415988224, 13: 0.18218670588881893, 14: 0.0, 15: 0.15797829756784315, 16: 0.0, 17: 0.0, 18: 1.4010319873918682, 19: 0.0, 20: 1.699487389697996, 21: 0.0, 22: 0.0, 23: 0.0, 24: 0.09324594165577986, 25: 0.0, 26: 0.4261775687725188, 27: 0.0, 28: 1.5130326659915523, 29: 0.0, 30: 0.0, 31: 0.0, 32: 0.05346219856811156, 33: 0.0}}
Top-K ranked results(Top-K Document Numbers):= [(1, 14.798683474291611), (28, 3.39897461939992), (28, 3.0268053319831046), (18, 2.8028639747837364), (26, 0.8523551375450376), (13, 0.36421341017763786), (15, 0.3159565951356863), (2, 0.18649188327155972), (32, 0.18692439712822311), (12, 0.0689743823817645), (7, 0.056866636265845224), (2, 0.0), (3, 0.0), (4, 0.0), (5, 0.0), (6, 0.0), (8, 0.0), (9, 0.0), (10, 0.0), (11, 0.0), (14, 0.0), (16, 0.0), (17, 0.0), (19, 0.0), (21, 0.0), (22, 0.0), (23, 0.0), (25, 0.0), (27, 0.0), (29, 0.0), (30, 0.0), (31, 0.0), (33, 0.0)]
```



```
Top-K ranked results(Query word -> similarity result):= {'Mulroney': {0: 0.0, 1: 7.395301737145886, 2: 0.0, 3: 0.0, 4: 0.0, 5: 0.0, 6: 0.0, 7: 0.028033318132922612, 8: 0.0, 9: 0.0, 10: 0.0, 11: 0.0, 12: 0.039487191415988224, 13: 0.18218670588881893, 14: 0.0, 15: 0.15797829756784315, 16: 0.0, 17: 0.0, 18: 1.4010319873918682, 19: 0.0, 20: 1.699487389697996, 21: 0.0, 22: 0.0, 23: 0.0, 24: 0.09324594165577986, 25: 0.0, 26: 0.4261775687725188, 27: 0.0, 28: 1.5130326659915523, 29: 0.0, 30: 0.0, 31: 0.0, 32: 0.05346219856811156, 33: 0.0}, 'denouncing': {0: 0.0, 1: 0.072733297536112, 2: 0.0, 3: 0.0, 4: 0.0, 5: 0.0, 6: 0.0, 7: 0.0, 8: 0.0, 9: 0.0, 10: 0.0, 11: 0.0, 12: 0.0, 13: 0.0, 14: 0.0, 15: 0.0, 16: 0.0, 17: 0.0, 18: 0.0, 19: 0.0, 20: 0.27758403232072975, 21: 0.0, 22: 0.0, 23: 0.0, 24: 0.0, 25: 0.0, 26: 0.0, 27: 0.0, 28: 0.0, 29: 0.0, 30: 0.0, 31: 0.0, 32: 0.0, 33: 0.0}}
Top-K ranked results(Top-K Document Numbers):= [(1, 14.86787480045222), (28, 3.6765586516262722), (28, 3.0268053319831046), (18, 2.8028639747837364), (26, 0.8523551375450376), (13, 0.36421341017763786), (15, 0.3159565951356863), (2, 0.18649188327155972), (32, 0.18692439712822311), (12, 0.0689743823817645), (7, 0.056866636265845224), (2, 0.0), (3, 0.0), (4, 0.0), (5, 0.0), (6, 0.0), (8, 0.0), (9, 0.0), (10, 0.0), (11, 0.0), (14, 0.0), (16, 0.0), (17, 0.0), (19, 0.0), (21, 0.0), (22, 0.0), (23, 0.0), (25, 0.0), (27, 0.0), (29, 0.0), (30, 0.0), (31, 0.0), (33, 0.0)]
```

127.0.0.1 - - [30/Apr/2022 17:17:52] "POST / HTTP/1.1" 200 -

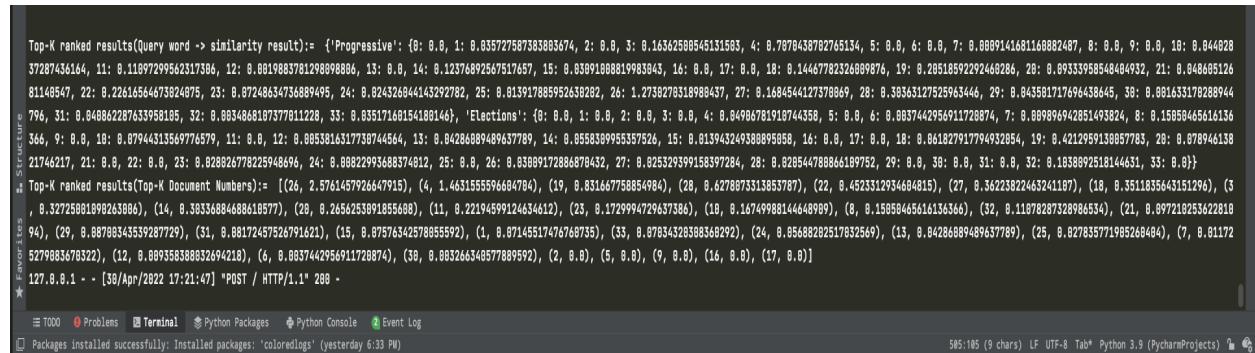
```
(base) harshpatel@Harshs-MacBook-Pro ~ % curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"Mulroney\", \"query2\" : \"denouncing\"}" "127.0.0.1:5000"
[
  [
    1,
    14.867876800045222
  ],
  [
    20,
    3.676558651626722
  ],
  [
    28,
    3.0260653319831046
  ],
  [
    18,
    2.8020639747837364
  ],
  [
    26,
    0.8523551375450376
  ],
  [
    23,
    0.36421341017763786
  ],
  [
    15,
    0.3159565951356863
  ],
  [
    24,
    0.18649188327155972
  ],
  [
    32,
    0.10692439712022311
  ],
  [
    12,
    0.86097438283017645
  ],
  [
    7,
    0.856866636265845224
  ],
  [
    2,
    0.0
  ],
  [
    3,
    0.0
  ],
  [
    4,
    0.0
  ],
  [
    5,
    0.0
  ],
  [

```

4th Query: -

```
curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"Progressive\", \"query2\" : \"Elections\"}" "127.0.0.1:5000/"
```

4th Query Result would be: - (In response in pycharm terminal and in terminal)



```
Top-K ranked results(Query word -> similarity result):= {'Progressive': {0: 0.0, 1: 0.035727587383883674, 2: 0.0, 3: 0.16362500545131503, 4: 0.7879438782765134, 5: 0.0, 6: 0.0, 7: 0.0889141681168882487, 8: 0.0, 9: 0.0, 10: 0.04460283728743614, 11: 0.11897299562317386, 12: 0.00198378129809886, 13: 0.0, 14: 0.12376892547517657, 15: 0.03091088819983843, 16: 0.0, 17: 0.0, 18: 0.14467702326099876, 19: 0.20518592292460286, 20: 0.09333958548404932, 21: 0.0486805126810547, 22: 0.2261654473524075, 23: 0.072486347536889495, 24: 0.024326944143592782, 25: 0.013917885952635202, 26: 1.2738279531998437, 27: 0.16845442127370869, 28: 0.03635127525963446, 29: 0.043581717646438645, 30: 0.091635179288944796, 31: 0.04886228763398105, 32: 0.0034868107377011228, 33: 0.03517168154188146}, 'Elections': {0: 0.0, 1: 0.0, 2: 0.0, 3: 0.0, 4: 0.8498678107443558, 5: 0.0, 6: 0.0, 7: 0.00989642851493824, 8: 0.1589465616136366, 9: 0.0, 10: 0.0, 11: 0.0, 12: 0.0053816317738744564, 13: 0.04268089489637789, 14: 0.0558309955357526, 15: 0.01394324938895858, 16: 0.0, 17: 0.0, 18: 0.0, 19: 0.4212959135887783, 20: 0.078946113621746217, 21: 0.0, 22: 0.0, 23: 0.028826782259484964, 24: 0.00222993688374912, 25: 0.0, 26: 0.03089172886870432, 27: 0.025529399158397284, 28: 0.025529399158397284, 29: 0.0, 30: 0.0, 31: 0.0, 32: 0.18389725181444531, 33: 0.0} }

Top-K ranked results(Top-K Document Numbers):= [(26, 2.5761457926647915), (4, 1.4631555596684794), (19, 0.83167758854984), (28, 0.6228873313853787), (22, 0.4523312934684815), (27, 0.36223822463241107), (18, 0.3511835643151296), (3, 0.3272580198263886), (14, 0.38336884688618577), (28, 0.2656253891855608), (11, 0.22194599124634612), (23, 0.1729994729637386), (18, 0.16749988144648989), (8, 0.15898465616135366), (32, 0.11878287328986554), (21, 0.09721825362281894), (29, 0.087980343535927729), (31, 0.00172457526791621), (15, 0.07576342578855952), (1, 0.07145517476768735), (33, 0.07853432858856292), (24, 0.05688202517832569), (13, 0.04268089489637789), (25, 0.027835771985268484), (7, 0.0117527983678322), (12, 0.00935838832694218), (6, 0.0037442956911728874), (30, 0.003263408577889592), (2, 0.0), (5, 0.0), (9, 0.0), (16, 0.0), (17, 0.0)}]

127.0.0.1 - - [30/Apr/2022 17:21:47] "POST / HTTP/1.1" 200
```

```
((base) harshpatel@Harshs-MacBook-Pro ~ % curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"Progressive\", \"query2\" : \"Elections\"}" "127.0.0.1:5000"
[
  [
    26,
    2.5761457926647915
  ],
  [
    4,
    1.4631555596604784
  ],
  [
    19,
    0.831667758854984
  ],
  [
    28,
    0.6278073313853787
  ],
  [
    22,
    0.4523312934604815
  ],
  [
    27,
    0.36223822463241107
  ],
  [
    18,
    0.3511835643151296
  ],
  [
    3,
    0.32725001090263806
  ],
  [
    14,
    0.30336884688610577
  ],
  [
    20,
    0.2656253091855688
  ],
  [
    11,
    0.22194599124634612
  ],
  [
    23,
    0.1729994729637386
  ],
  [
    10,
    0.16749988144648909
  ],
  [
    8,
    0.15050465616136366
  ],
  [
    32,
    0.11078287328986534
  ],
  [

```

5th Query: -

```
curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"Progressive Affiliation Conservative NDP\"}" "127.0.0.1:5000/"
```

5th Query Result would be: - (In response in pycharm terminal and in terminal)

```
Terminal: Local: Local: 2) + 
Top-K ranked results(Query word -> similarity result): {'Progressive': {0: 0.0, 1: 0.835727587383883674, 2: 0.0, 3: 0.16562588545131593, 4: 0.8708438782765134, 5: 0.0, 6: 0.0, 7: 0.0889141681168882487, 8: 0.0, 9: 0.0, 10: 0.84462837287456164, 11: 0.1189729562317586, 12: 0.081988379129809886, 13: 0.0, 14: 0.12376892567517657, 15: 0.0899188819983943, 16: 0.0, 17: 0.0, 18: 0.14467782326869876, 19: 0.085185722722468286, 20: 0.0933598854846932, 21: 0.0466851268110952, 22: 0.22616564673924975, 23: 0.072468454736889495, 24: 0.02432684143292782, 25: 0.01917885952650262, 26: 1.273827651898437, 27: 0.04356312752963446, 28: 0.0816331792889447, 29: 0.043563171796436465, 30: 0.0816331792889447, 31: 0.0408622876353958105, 32: 0.0834688187377011228, 33: 0.05317160151810146}, 'Affiliation': {0: 0.0, 1: 0.0, 2: 0.0, 3: 0.0, 4: 0.3524416985895989, 5: 0.0, 6: 0.0, 7: 0.0, 8: 0.0, 9: 0.0, 10: 0.0, 11: 0.0, 12: 0.0, 13: 0.0, 14: 0.0, 15: 0.0, 16: 0.0, 17: 0.0, 18: 0.0, 19: 0.0, 20: 0.1134188166651773, 21: 0.0, 22: 0.0, 23: 0.0, 24: 0.0, 25: 0.0, 26: 0.19887862731566365, 27: 0.0, 28: 0.0, 29: 0.0, 30: 0.0, 31: 0.0, 32: 0.0, 33: 0.0}, 'Conservative': {0: 0.0, 1: 0.0, 2: 0.0, 3: 0.0, 4: 0.24756743683743648, 5: 0.0, 6: 0.0, 7: 0.0, 8: 0.0, 9: 0.0, 10: 0.0, 11: 0.0, 12: 0.0, 13: 0.0, 14: 0.0, 15: 0.0, 16: 0.0, 17: 0.0, 18: 0.0, 19: 0.0, 20: 0.0, 21: 0.0, 22: 0.0, 23: 0.0, 24: 0.0, 25: 0.0, 26: 0.0, 27: 0.0, 28: 0.0, 29: 0.0, 30: 0.0, 31: 0.0, 32: 0.0, 33: 0.0}, 'NDP': {0: 0.0, 1: 0.0, 2: 0.0, 3: 0.0, 4: 0.0, 5: 0.0, 6: 0.0, 7: 0.0, 8: 0.0, 9: 0.0, 10: 0.0, 11: 0.0, 12: 0.0, 13: 0.0, 14: 0.0, 15: 0.0, 16: 0.0, 17: 0.0, 18: 0.0, 19: 0.0, 20: 0.0, 21: 0.0, 22: 0.0, 23: 0.0, 24: 0.0, 25: 0.0, 26: 0.0, 27: 0.0, 28: 0.0, 29: 0.0, 30: 0.0, 31: 0.0, 32: 0.0, 33: 0.0}, 'Top-K Document Numbers': [{(26, 2.49983749742545), (4, 1.7212754875831898), (14, 0.9566724914123441), (27, 0.913452697788263), (23, 0.848896432618853), (28, 0.64373833579667538), (3, 0.8173759943876615), (24, 0.4898669385671763), (22, 0.42294633478422), (19, 0.84198798175737), (20, 0.38937862618469), (21, 0.38663217587264265), (18, 0.3286122988768715), (11, 0.24392116854168897), (31, 0.188817179586974384), (25, 0.15991784261629167), (19, 0.1281147954772842), (15, 0.108649985253875487), (33, 0.09656922467879252), (1, 0.88986171931331843), (29, 0.88134755313641994), (32, 0.82265749453244485), (2, 0.828989262654181315), (12, 0.8888371435491715143), (13, 0.8879897407215924), (38, 0.884929897992621423), (7, 0.88498541769481682), (6, 0.0827919110722715973), (17, 0.8824384214194947), (5, 0.0), (8, 0.0), (9, 0.0), (16, 0.0)}]
```

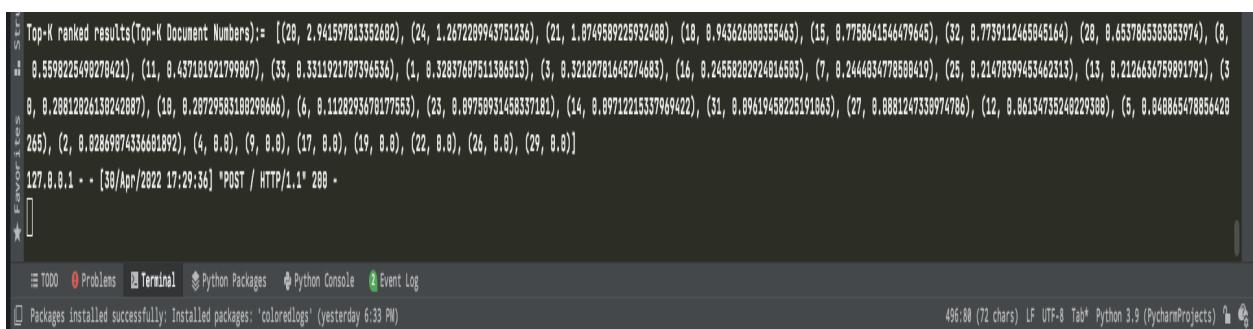
127.0.0.1 - - [30/Apr/2022 17:25:08] "POST / HTTP/1.1" 208 -

```
(base) harshpatel@Harshs-MacBook-Pro ~ % curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"Progressive Affiliation Conservative NDP\"}" "127.0.0.1:5000/"
[{"id": 26, "score": 2.499830749742545}, {"id": 4, "score": 1.7212754875031098}, {"id": 14, "score": 0.9566724914123441}, {"id": 27, "score": 0.9134526970788263}, {"id": 23, "score": 0.848896432618853}, {"id": 28, "score": 0.6437303357966738}, {"id": 3, "score": 0.5173759943076615}, {"id": 24, "score": 0.48019660305671763}, {"id": 22, "score": 0.4229263334704422}, {"id": 19, "score": 0.41987298175737}, {"id": 20, "score": 0.3898372601016489}, {"id": 21, "score": 0.38663217587264265}, {"id": 18, "score": 0.32061229008706715}, {"id": 11, "score": 0.24392116854160087}, {"id": 31, "score": 0.18801729586074384}
```

6th Query: -

```
curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"Despite losing percentage point terms vote share\"}" "127.0.0.1:5000/"
```

6th Query Result would be: - (In response in pycharm terminal and in terminal)



```
Top-K ranked results(Top-K Document Numbers)= [(28, 2.841597813352682), (24, 1.2672289943751236), (21, 1.8749589225932468), (18, 0.943626688335463), (15, 0.7758841546479646), (32, 0.7739112465845164), (28, 0.653785383853974), (8, 0.5598225409278421), (11, 0.437101921799867), (33, 0.3311921787396536), (1, 0.32837687511386513), (3, 0.32182781645274683), (16, 0.24558282924016593), (7, 0.2444034778589419), (25, 0.21478399453462313), (13, 0.2126636759891791), (3, 0.288122926138242887), (18, 0.28729583188298666), (6, 0.1128293678177553), (23, 0.09758931458337181), (14, 0.09712215337969422), (31, 0.09619458225191863), (27, 0.0881247338974786), (12, 0.06134735246229398), (5, 0.040845478856428265), (2, 0.02869874336631892), (4, 0.0), (9, 0.0), (17, 0.0), (19, 0.0), (22, 0.0), (26, 0.0), (29, 0.0)]
```

127.0.0.1 - [30/Apr/2022 17:29:36] "POST / HTTP/1.1" 200

PyCharm Terminal window showing the command and its execution.

```
(base) harshpatel@Harshs-MacBook-Pro ~ % curl -X POST -H "Content-type: application/json" -d "{\"query1\" : \"Despite losing percentage point terms vote share\"}" "127.0.0.1:5000/"
[{"id": 28, "score": 2.941597813352602}, {"id": 24, "score": 1.2672209943751236}, {"id": 21, "score": 1.8749589225932408}, {"id": 18, "score": 0.9436260008355463}, {"id": 15, "score": 0.7758641546479645}, {"id": 32, "score": 0.7739112465045164}, {"id": 28, "score": 0.6537865303853974}, {"id": 8, "score": 0.5598225490270421}, {"id": 21, "score": 0.437101921799867}, {"id": 33, "score": 0.3311921787396536}, {"id": 1, "score": 0.32837607511386513}, {"id": 3, "score": 0.32182781645274683}, {"id": 16, "score": 0.24558202924016583}, {"id": 7, "score": 0.2444034778500419}, {"id": 25, "score": 0.21470399453462313}, {"id": 1, "score": 0.1470399453462313}
```

Success: - The system is almost working fine except following flaws.

1. It will remove stop words from query even if present there. Top-K results has almost all accuracy. I got almost all the result relating big corpus too.

Cautions/Flaws [Test needed to be taken care for future, while checking for large corpus]: -

My system is not fault-proof. I have some bug to work on, pops on when I tested large corpus on last minute.

1. The major thing I faced today is, if I search one query for top-k results in flask server, and get result, and if I search another one, **I got Error like AttributeError: 'dict' object has no attribute 'extend', but if I think it's my system design fault. If you try multiple queries JSON, please bear to restart flask server again to get the query result, and check if no server address already in use.**
2. **If you run another URL, in case, please try to change,**

```
query_to_check = "Oshawa Vancouver"
```

Change it to, words in new wordlist from new corpus from new URL.

```
query_to_check = "{words that are in new token list/You can search it from any words from new webpage that saved}"
```

7 Data Sources: -

- ⇒ For Indexing, Processing -> I use primarily Dictionary data structure for Inverted Index for each, like, tokenized words dictionary, DF for each word, IDF score, TD-IDF score, query vector and cosine similarity. I almost use dictionary for indexing and query processing in Flask.
- ⇒ For corpus data, I used user defined argument to give user-input url. I used url as, https://en.wikipedia.org/wiki/George_Minaker
- ⇒ -a define user-defined argument for url. Example is given below.

```
scrapy crawl MainSpider -a url=https://en.wikipedia.org/wiki/George_Minaker
```

- ⇒ I use BeautifulSoup package for html parser for locally saved HTML pages from tokenization of words from like 'UTF-8'. The import package is given below.

```
from bs4 import BeautifulSoup
```

- ⇒ BeautifulSoup Link: - <https://www.crummy.com/software/BeautifulSoup/bs4/doc/>

8 Test Cases: -

- ⇒ Please check Test case with 6 different queries in conclusion part, on page 16.

9 Source Code: -

- ⇒ Project directory structure is shown in Page 4.
- ⇒ The source code is given below.

FOR SCRAPY BASED CRAWLER: -

- ⇒ **Python File: cs429crawler/spiders/mainspider.py**

```
import scrapy
import os
import requests
import logging

class MainSpider(scrapy.Spider):
    name = "MainSpider"
    logging.getLogger("urllib3").setLevel(logging.WARNING)

    custom_settings = {
        'DEPTH_LIMIT': '2',
    }
```

```

def __init__(self, url='', **kwargs):
    self.links = []
    self.allowed_domains = ['en.wikipedia.org', 'www.iit.edu',
'www.phdessay.com']
    self.start_urls = [url] # py36
    self.base_url = 'en.wikipedia.org'
    super(MainSpider, self).__init__(**kwargs) # python3

save_path = 'htmlFiles'

def parse(self, response, dest_path = save_path, save_path=save_path):
    init = 0
    for link in response.xpath('//div/p/a'):
        # if init == 2:
        #     break
        yield {
            "link": self.base_url + link.xpath('.//@href').get()
        }
    print(self.base_url + link.xpath('.//@href').get())
    filePath = self.base_url + link.xpath('.//@href').get()
    words = ['tel', 'phone', 'mob', 'mailto', 'Special']
    filePath_upd = filePath.rsplit('/', 1)[-1] + '.html'
    for word in words:
        if word in filePath:
            filePath = self.base_url
    print(requests.get('https://{}'.format(filePath)))
    if requests.get('https://{}'.format(filePath)).status_code ==
200:
        for f in os.listdir(save_path):
            if init == 0:
                if not f.endswith(".html"):
                    continue
                os.remove(os.path.join(save_path, f))
            if os.path.join('htmlFiles', filePath_upd):
                with open(os.path.join('htmlFiles', filePath_upd), 'w') :
                    pass
                with open(os.path.join('htmlFiles', filePath_upd), "wb") as
file:
                    # response = get(self.base_url +
link.xpath('.//@href').get())
                    response = requests.get('https://{}'.format(filePath))
                    file.write(response.content)
                # if requests.get('http://{}'.format(filePath)).status_code ==
200:
                    # for f in os.listdir(save_path):
                    #     if init == 0:
                    #         if not f.endswith(".html"):
                    #             continue
                    #         os.remove(os.path.join(save_path, f))
                    #     if os.path.join('htmlFiles', filePath_upd):
                    #         with open(os.path.join('htmlFiles', filePath_upd),
'w') :
                    #             pass
                    #         with open(os.path.join('htmlFiles', filePath_upd), "wb") as
file:
                        #             # response = get(self.base_url +
link.xpath('.//@href').get())

```

```

#           response = requests.get('http://{}'.format(filePath))
#           file.write(response.content)
else:
    pass
init += 1

```

⇒ **Python File: cs429crawler/middlewares.py**

```

# Define here the models for your spider middleware
#
# See documentation in:
# https://docs.scrapy.org/en/latest/topics/spider-middleware.html

from scrapy import signals
from scrapy.http import Request
import tldextract
import logging

logger = logging.getLogger(__name__)
# useful for handling different item types with a single interface
from itemadapter import is_item, ItemAdapter


class Cs429CrawlerSpiderMiddleware:
    # Not all methods need to be defined. If a method is not defined,
    # scrapy acts as if the spider middleware does not modify the
    # passed objects.

    def __init__(self, maxdepth, stats, verbose_stats=False, prio=1):
        self.maxdepth = maxdepth
        self.stats = stats
        self.verbose_stats = verbose_stats
        self.prio = prio

    @classmethod
    def from_crawler(cls, crawler):
        # This method is used by Scrapy to create your spiders.

        # s = cls()
        # crawler.signals.connect(s.spider_opened,
signal=signals.spider_opened)
        # return s

        settings = crawler.settings
        maxdepth = settings.getint('DEPTH_LIMIT')
        verbose = settings.getbool('DEPTH_STATS_VERBOSE')
        prio = settings.getint('DEPTH_PRIORITY')
        return cls(maxdepth, crawler.stats, verbose, prio)

    def process_spider_input(self, response, spider):
        # Called for each response that goes through the spider
        # middleware and into the spider.

```

```

# Should return None or raise an exception.
return None

def process_spider_output(self, response, result, spider):
    def _filter(request):
        if isinstance(request, Request):
            depth = response.meta['depth'] + 1
            request.meta['depth'] = depth
            if self.prio:
                request.priority -= depth * self.prio
            if self.maxdepth and depth > self.maxdepth:
                logging.debug(
                    "Ignoring link (depth > %(maxdepth)d): %(requrl)s",
                    {'maxdepth': self.maxdepth, 'requrl': request.url},
                    extra={'spider': spider})
        )
        return False
    else:
        if self.verbose_stats:
            self.stats.inc_value(f'request_depth_count/{depth}', spider=spider)
            self.stats.max_value('request_depth_max', depth, spider=spider)
    return True

# base case (depth=0)
if 'depth' not in response.meta:
    response.meta['depth'] = 0
    if self.verbose_stats:
        self.stats.inc_value('request_depth_count/0', spider=spider)

    return (r for r in result or () if _filter(r))

def process_spider_exception(self, response, exception, spider):
    # Called when a spider or process_spider_input() method
    # (from other spider middleware) raises an exception.

    # Should return either None or an iterable of Request or item
    objects.
    pass

def process_start_requests(self, start_requests, spider):
    # Called with the start requests of the spider, and works
    # similarly to the process_spider_output() method, except
    # that it doesn't have a response associated.

    # Must return only requests (not items).
    for r in start_requests:
        yield r

def spider_opened(self, spider):
    spider.logger.info('Spider opened: %s' % spider.name)

class Cs429CrawlerDownloaderMiddleware:
    # Not all methods need to be defined. If a method is not defined,
    # scrapy acts as if the downloader middleware does not modify the

```

```

# passed objects.

@classmethod
def from_crawler(cls, crawler):
    # This method is used by Scrapy to create your spiders.
    s = cls()
    crawler.signals.connect(s.spider_opened,
signal=signals.spider_opened)
    return s

def process_request(self, request, spider):
    # Called for each request that goes through the downloader
    # middleware.

    # Must either:
    # - return None: continue processing this request
    # - or return a Response object
    # - or return a Request object
    # - or raise IgnoreRequest: process_exception() methods of
    #   installed downloader middleware will be called
    return None

def process_response(self, request, response, spider):
    # Called with the response returned from the downloader.

    # Must either:
    # - return a Response object
    # - return a Request object
    # - or raise IgnoreRequest
    return response

def process_exception(self, request, exception, spider):
    # Called when a download handler or a process_request()
    # (from other downloader middleware) raises an exception.

    # Must either:
    # - return None: continue processing this exception
    # - return a Response object: stops process_exception() chain
    # - return a Request object: stops process_exception() chain
    pass

def spider_opened(self, spider):
    spider.logger.info('Spider opened: %s' % spider.name)

```

⇒ [Python File: cs429crawler/settings.py](#)

```

# Scrapy settings for cs429crawler project
#
# For simplicity, this file contains only settings considered important or
# commonly used. You can find more settings consulting the documentation:
#
#     https://docs.scrapy.org/en/latest/topics/settings.html
#     https://docs.scrapy.org/en/latest/topics/downloader-middleware.html
#     https://docs.scrapy.org/en/latest/topics/spider-middleware.html

```

```

BOT_NAME = 'cs429crawler'

SPIDER_MODULES = ['cs429crawler.spiders']
NEWSPIDER_MODULE = 'cs429crawler.spiders'

DEPTH_LIMIT = 2
SCHEDULER_DISK_QUEUE = 'scrapy.squeues.PickleFifoDiskQueue'
SCHEDULER_MEMORY_QUEUE = 'scrapy.squeues.FifoMemoryQueue'

SPIDER_MIDDLEWARES = {'cs429crawler.middlewares.Cs429CrawlerSpiderMiddleware': 0}

# FEED_FORMAT="csv"
# FEED_URI="iitEdu.csv"

# Crawl responsibly by identifying yourself (and your website) on the user-agent
#USER_AGENT = 'cs429crawler (+http://www.yourdomain.com)'

# Obey robots.txt rules
ROBOTSTXT_OBEY = True

# Configure maximum concurrent requests performed by Scrapy (default: 16)
#CONCURRENT_REQUESTS = 32

# Configure a delay for requests for the same website (default: 0)
# See https://docs.scrapy.org/en/latest/topics/settings.html#download-delay
# See also autothrottle settings and docs
#DOWNLOAD_DELAY = 3
# The download delay setting will honor only one of:
#CONCURRENT_REQUESTS_PER_DOMAIN = 16
#CONCURRENT_REQUESTS_PER_IP = 16

# Disable cookies (enabled by default)
#COOKIES_ENABLED = False

# Disable Telnet Console (enabled by default)
#TELNETCONSOLE_ENABLED = False

# Override the default request headers:
#DEFAULT_REQUEST_HEADERS = {
#    'Accept': 'text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8',
#    'Accept-Language': 'en',
#}

# Enable or disable spider middlewares
# See https://docs.scrapy.org/en/latest/topics/spider-middleware.html
#SPIDER_MIDDLEWARES = {
#    'cs429crawler.middlewares.Cs429CrawlerSpiderMiddleware': 543,
#}

# Enable or disable downloader middlewares
# See https://docs.scrapy.org/en/latest/topics/downloader-middleware.html
#DOWNLOADER_MIDDLEWARES = {
#    'cs429crawler.middlewares.Cs429CrawlerDownloaderMiddleware': 543,
#}

```

```

# Enable or disable extensions
# See https://docs.scrapy.org/en/latest/topics/extensions.html
#EXTENSIONS = {
#     'scrapy.extensions.telnet.TelnetConsole': None,
# }

# Configure item pipelines
# See https://docs.scrapy.org/en/latest/topics/item-pipeline.html
#ITEM_PIPELINES = {
#     'cs429crawler.pipelines.Cs429CrawlerPipeline': 300,
# }

# Enable and configure the AutoThrottle extension (disabled by default)
# See https://docs.scrapy.org/en/latest/topics/autothrottle.html
#AUTOTHROTTLE_ENABLED = True
# The initial download delay
#AUTOTHROTTLE_START_DELAY = 5
# The maximum download delay to be set in case of high latencies
#AUTOTHROTTLE_MAX_DELAY = 60
# The average number of requests Scrapy should be sending in parallel to
# each remote server
#AUTOTHROTTLE_TARGET_CONCURRENCY = 1.0
# Enable showing throttling stats for every response received:
#AUTOTHROTTLE_DEBUG = False

# Enable and configure HTTP caching (disabled by default)
# See https://docs.scrapy.org/en/latest/topics/downloader-
middleware.html#httpcache-middleware-settings
#HTTPCACHE_ENABLED = True
#HTTPCACHE_EXPIRATION_SECS = 0
#HTTPCACHE_DIR = 'httpcache'
#HTTPCACHE_IGNORE_HTTP_CODES = []
#HTTPCACHE_STORAGE = 'scrapy.extensions.httpcache.FilesystemCacheStorage'

```

⇒ **Python File: cs429crawler/items.py**

```

# Define here the models for your scraped items
#
# See documentation in:
# https://docs.scrapy.org/en/latest/topics/items.html

import scrapy


class Cs429CrawlerItem(scrapy.Item):
    # define the fields for your item here like:
    # name = scrapy.Field()
    pass

```

FOR SCIKIT-LEARN BASED INDEXER: -

⇒ **Python File: indexer/main.py**

```
⇒ from bs4 import BeautifulSoup
import re
import os
import sys
import numpy as np
import pickle

html_parsing_list = list()
df_list = list()
directory = os.fsencode('htmlFiles')
ROOT_PATH = os.path.dirname(os.path.abspath('htmlFiles'))

def main():
    init = 0
    start = 0
    dict = {}
    idf_dict = {}
    tf_idf_dict = {}
    tf_idf_dict_opt_ext = {}
    query_vector_dict = {}
    cosine_similarity_dict_parent = {}
    cosine_similarity_dict_parent_2 = {}
    cosine_similarity_dict_parent_final = {}
    cosine_similarity_dict_parent_term_list = {}
    # word = "Oklahoma"

    # Tokenized-words and DF-list
    for file in os.listdir(directory):
        # if init == 1:
        #     break
        with open('iitEdu.csv', 'w') as f:
            f.truncate()
        filename = os.fsdecode(file)
        if filename.endswith(".html"):
            filePath = ROOT_PATH + "/htmlFiles/" + filename
            # print(filePath)
            HtmlFile = open(filePath, 'r', encoding='utf-8')
            source_code = HtmlFile.read()
            soup = BeautifulSoup(source_code, "html.parser")
            with open('iitEdu.csv', 'w') as f:
                for line in soup.text:
                    f.write(line)
            # upd_clear_str = re.sub("[^\w]", " ", str(soup.text))
            # wordList = re.split('\s+', upd_clear_str)
            # print("length of wordList:= ", len(wordList))

            wordList = re.sub("[^\w]", " ", str(soup.text)).split()
            stops = ["", "-", "i", "me", "my", "myself", "we", "our",
                     "ours", "ourselves", "you", "your", "yours", "yourself", "yourselves",
                     "he", "him", "his", "himself", "she", "her", "hers", "herself", "it",
                     "its", "itself", "they", "them", "their", "theirs", "themselves", "what",
                     "which", "who", "whom", "this", "that", "these", "those", "am", "is",
                     "are", "was", "were", "be", "been", "being", "have", "has", "had",
                     "having", "do", "does", "did", "doing", "a", "an", "the", "and", "but",
                     "if", "or", "because", "as", "until", "while", "of", "at", "by", "for",
                     "with", "about", "against", "between", "into", "through", "during",
                     "before", "after", "above", "below", "to", "from", "up", "down", "in"]
```

```

"out", "on", "off", "over", "under", "again", "further", "then", "once",
"here", "there", "when", "where", "why", "how", "all", "any", "both",
"each", "few", "more", "most", "other", "some", "such", "no", "nor",
"not", "only", "own", "same", "so", "than", "too", "very", "s", "t",
"can", "will", "just", "don", "should", "now"]
    tokenized_word = [i for i in wordList if str(i).lower() not in
stops]
    # print(tokenized_word)
    # print("length of tokenized_word:= ", len(tokenized_word))
    html_parsing_list.append(tokenized_word)

    # DF - list implementation
    DF = {}
    upd_DF = {}
    upd_upd_DF = {}
    for w in tokenized_word:
        if w in DF:
            DF[w].append(w)
        else:
            DF[w] = [w]
    for i in DF:
        upd_DF[i] = len(DF.get(i))
    for i in DF:
        upd_upd_DF[i] = upd_DF.get(i) / len(tokenized_word)
    df_list.append(upd_upd_DF)

    # Create a dictionary for all the words.
    check = str(soup.text)
    # with open('iitEdu_1.csv', 'w') as f:
    #     for line in check:
    #         f.write(line)
    # print("html_parsing_list[init]:= ", html_parsing_list[init])
    # if word in check:
    #     print(">> Oklahoma word is present!")
    for item in html_parsing_list[init]:
        if item in check:
            if item not in dict:
                dict[item] = []
            if item in dict and (init + 1) not in dict.get(item):
                dict[item].append(init + 1)

    # Convert Index pickling of dictionary index pages pickle
files.
    fileName = "save_" + str(init) + ".pkl"
    if not os.path.isfile(os.path.join('picklefiles', fileName)):
        with open(os.path.join('picklefiles', fileName), 'wb') as
f:
            pickle.dump(dict, f, protocol=pickle.HIGHEST_PROTOCOL)
    else:
        os.remove(os.path.join('picklefiles', fileName))
        with open(os.path.join('picklefiles', fileName), 'wb') as
f:
            pickle.dump(dict, f, protocol=pickle.HIGHEST_PROTOCOL)

    init += 1

```

```

# if word in dict:
#     print("The word is in the list!")
# print("dict score of term 'Binary':= ", dict["Binary"])
# print("dict score of term 'Binary':= ", dict["Oshawa"])

# print("df_list:= ", df_list[18])

# Pickle file - html_parsing_list list
if not os.path.isfile('picklefiles/html_parsing_list.pkl'):
    with open('picklefiles/html_parsing_list.pkl', 'wb') as f:
        pickle.dump(html_parsing_list, f,
protocol=pickle.HIGHEST_PROTOCOL)
else:
    os.remove('picklefiles/html_parsing_list.pkl')
    with open('picklefiles/html_parsing_list.pkl', 'wb') as f:
        pickle.dump(html_parsing_list, f,
protocol=pickle.HIGHEST_PROTOCOL)

# IDF_score logic
for key, values in dict.items():
    idf_score = np.log(len(html_parsing_list) / len(values))
    if key not in idf_dict:
        idf_dict[key] = list()
    for value in values:
        idf_dict[key].append([value, idf_score])
# print("idf_dict score of term 'Binary':= ", idf_dict["Binary"])
# print("idf_dict score of term 'Binary':= ", idf_dict["Oshawa"])

# Pickle file - IDF
if not os.path.isfile('picklefiles/idf.pkl'):
    with open('picklefiles/idf.pkl', 'wb') as f:
        pickle.dump(idf_dict, f, protocol=pickle.HIGHEST_PROTOCOL)
else:
    os.remove('picklefiles/idf.pkl')
    with open('picklefiles/idf.pkl', 'wb') as f:
        pickle.dump(idf_dict, f, protocol=pickle.HIGHEST_PROTOCOL)

# Find TF-IDF score
init_token = 0
for token_item in range(len(html_parsing_list)):
    # if init_token == 1:
    #     break
    for item in df_list[token_item]:
        if item not in tf_idf_dict:
            tf_idf_dict[item] = {}
        if item in tf_idf_dict:
            tf_idf_dict_opt = {}
            doc_ids = [docs[0] for docs in idf_dict[item]]
            doc_ids_idf_score = [docs[1] for docs in idf_dict[item]]
            # print("item:= ", item)
            # print("doc_ids:= ", doc_ids)
            # print("doc_ids_idf_score:= ", doc_ids_idf_score)
            count = 0
            for number in doc_ids:
                if number not in tf_idf_dict_opt:
                    tf_idf_dict_opt[number] = []
                if number not in tf_idf_dict_opt_ext:

```

```

        tf_idf_dict_opt_ext[number] = []
        if number in tf_idf_dict_opt:
            # print("doc_ids_idf_score[count]:= ", doc_ids_idf_score[count])
            # print("df_list[number - 1][item]:= ", df_list[number - 1][item])

        tf_idf_dict_opt[number].append(doc_ids_idf_score[count] * df_list[number - 1][item])

        tf_idf_dict_opt_ext[number].append(doc_ids_idf_score[count] * df_list[number - 1][item])
            count = count + 1
            tf_idf_dict[item] = tf_idf_dict_opt
            init_token += 1
            #elementsArticles
            # print("tf_idf_dict_opt_ext:= ", tf_idf_dict_opt_ext)
            # print("tf_idf_dict score of term 'Binary':= ", tf_idf_dict["Binary"])
            tf_idf_dict["Binary"])
            # print("tf_idf_dict score of term 'Oshawa':= ", tf_idf_dict["Oshawa"])
            # print("tf_idf_dict score of term 'Vancouver':= ", tf_idf_dict["Vancouver"])
            # print("tf_idf_dict_opt_ext score:= ", len(tf_idf_dict_opt_ext))

        # Pickle file - TF/IDF
        if not os.path.isfile('picklefiles/td_idf.pkl'):
            with open('picklefiles/td_idf.pkl', 'wb') as f:
                pickle.dump(tf_idf_dict, f, protocol=pickle.HIGHEST_PROTOCOL)
        else:
            os.remove('picklefiles/td_idf.pkl')
            with open('picklefiles/td_idf.pkl', 'wb') as f:
                pickle.dump(tf_idf_dict, f, protocol=pickle.HIGHEST_PROTOCOL)

        # Pickle file - TF_IDF_OPT
        if not os.path.isfile('picklefiles/tf_idf_opt.pkl'):
            with open('picklefiles/tf_idf_opt.pkl', 'wb') as f:
                pickle.dump(tf_idf_dict_opt_ext, f,
protocol=pickle.HIGHEST_PROTOCOL)
        else:
            os.remove('picklefiles/tf_idf_opt.pkl')
            with open('picklefiles/tf_idf_opt.pkl', 'wb') as f:
                pickle.dump(tf_idf_dict_opt_ext, f,
protocol=pickle.HIGHEST_PROTOCOL)

        # The vector should be a dict with query terms as keys and IDF scores as values.
        query_to_check = "Oshawa Vancouver"
        stops = ["", "-", "i", "me", "my", "myself", "we", "our", "ours",
"ourselves", "you", "your", "yours", "yourself",
"yourselves", "he", "him", "his", "himself", "she", "her",
"hers", "herself", "it", "its", "itself",
"they", "them", "their", "theirs", "themselves", "what",
"which", "who", "whom", "this", "that", "these",
"those", "am", "is", "are", "was", "were", "be", "been",
"being", "have", "has", "had", "having", "do",
"does", "did", "doing", "a", "an", "the", "and", "but", "if",

```

```

"or", "because", "as", "until", "while",
        "of", "at", "by", "for", "with", "about", "against",
"between", "into", "through", "during", "before",
        "after", "above", "below", "to", "from", "up", "down", "in",
"out", "on", "off", "over", "under", "again",
        "further", "then", "once", "here", "there", "when", "where",
"why", "how", "all", "any", "both", "each",
        "few", "more", "most", "other", "some", "such", "no", "nor",
"not", "only", "own", "same", "so", "than",
        "too", "very", "s", "t", "can", "will", "just", "don",
"should", "now"]
    data = re.sub('[^\\w|_| ', ' ', query_to_check) # only keep numbers and
letters and spaces
    # data = data.lower()
    data = re.sub(r'[^\\x00-\\x7f]', r'', data) # remove non ascii texts
    tokens = [data_upt for data_upt in data.split(' ') if data_upt] #
split the words and remove empty words
    tokens = [word for word in tokens if word.lower() not in stops]
    for token in tokens:
        doc_ids_idf_score = [idf_dict[token][0][1]]
        query_vector_dict[token] = doc_ids_idf_score
    print("sample Query vector score:= ", query_vector_dict)

    # Cosine similarity
    # The result of the dot product between the query vector and the
    documents (per term) need to be converted into a final sorted list of
    pairs containing the document identifier and score
    # The sorting is based on highest score to lowest.
    terms = [docs for docs in query_vector_dict]
    terms_idf_val = [query_vector_dict.get(docs)[0] for docs in
query_vector_dict]
    for term_num in range(len(terms)):
        doc_length = []
        cosine_similarity_dict_clone_res = {}
        dict_keys = list(tf_idf_dict[terms[term_num]].keys())
        cosine_similarity_dict = {}
        wt = terms_idf_val[term_num]
        for idx in range(len(html_parsing_list)):
            if idx not in dict_keys:
                cosine_similarity_dict[idx] = 0
            else:
                cosine_similarity_dict[idx] = terms_idf_val[0] *
tf_idf_dict[terms[term_num]].get(idx)[0]
                doc_length.append(tf_idf_dict_opt_ext.get(idx)[0])
        doc_square = [i ** 2 for i in doc_length]
        doc_length_sqrt_val = pow(sum(doc_square), 0.5)
        for key, values in cosine_similarity_dict.items():
            cosine_similarity_dict_clone_res[key] =
cosine_similarity_dict.get(key) / doc_length_sqrt_val
        for key, values in cosine_similarity_dict_clone_res.items():
            if key not in cosine_similarity_dict_parent and key != 0:
                cosine_similarity_dict_parent[key] = []
            if key in cosine_similarity_dict_parent:
                cosine_similarity_dict_parent[key].append(values)
        for key, values in cosine_similarity_dict_parent.items():
            if key not in cosine_similarity_dict_parent:
                cosine_similarity_dict_2[key] = 0

```

```

        if key in cosine_similarity_dict_parent:
            cosine_similarity_dict_parent_2[key] =
sum(cosine_similarity_dict_parent.get(key))
            cosine_similarity_dict_parent_term_list[terms[term_num]] =
cosine_similarity_dict_clone_res
            print("cosine_similarity_dict_parent_term_list:= ",
cosine_similarity_dict_parent_term_list)
            cosine_similarity_dict_parent_final =
sorted(cosine_similarity_dict_parent_2.items(), key=lambda t: t[1],
reverse=True)

        # "document identifier, score" for query and document pages for
combine query vector and document vector.
        print("Sorted cosine_similarity_dict_parent_final:= ",
cosine_similarity_dict_parent_final)

if __name__ == '__main__':
    main()

```

FOR FLASK BASED PROCESSOR: -

⇒ **Python File: processor/config.py**

```

import os

LOG_LEVEL = 'DEBUG'
LOG_FORMAT = '%(asctime)s %(levelname)s : %(message)s'
APPLICATION_ROOT = '/cs429crawler'
FLASK_RUN_MODE = os.environ.get('MODE') or 'PROD'
WEB_PORT = os.environ.get('PORT') or 5000

```

⇒ **Python File: processor/logger.py**

```

import coloredlogs
import logging

from config import LOG_LEVEL, LOG_FORMAT

# Create a logger object.
logger = logging.getLogger(__name__)

# inject color log
coloredlogs.install(level=LOG_LEVEL, fmt=LOG_FORMAT)

```

⇒ **Python File: processor/app.py**

```

from flask import Flask, request, jsonify, make_response
import pickle
import os
import re

```

```

from pathlib import Path

app = Flask(__name__)

# TF-IDF dict objects <- To use pickle files
main_res_dict_tf_idf = {}
tf_idf_dict = {}
dict_new_tf_idf = {}

# IDF dict objects <- To use pickle files
idf_dict = {}
dict_new_idf = {}

# TF_IDF_OPT dict objects <- To use pickle files
main_res_dict_tf_idf_opt = {}
tf_idf_dict_opt_ext = {}
dict_new_tf_idf_opt = {}

# Query vector objects <- To use pickle files
query_vector_dict = {}
dict_new_inverted_index = {}
main_res_dict_inverted_index = {}

# Cosine similarity objects -> To get similarity score.
cosine_similarity_dict_parent = {}
cosine_similarity_dict_parent_2 = {}
cosine_similarity_dict_parent_final = {}
cosine_similarity_dict_parent_term_list = {}

# html_parsing_list object <- To use pickle files
html_parsing_list = list()

@app.route('/', methods=['POST'])
def process_json():
    global cosine_similarity_dict_parent_final
    content_type = request.headers.get('Content-Type')
    if (content_type == 'application/json'):
        query_list = list()
        json = request.json
        # query_1 = json['query1']
        # print("query_1:= ", query_1)
        for idx in range(len(json)):
            init = idx + 1
            queryName = "query" + str(init)
            query_list.append(json[queryName])
        print("Query list(from Json format):= ", query_list)

        # Using block-sorting algorithm, merged/combined together for the
        final one inverted index.
        myDir = Path('picklefiles/')
        fileNames = [file.name for file in myDir.iterdir() if
        file.name.startswith('save')]
        for file in fileNames:
            with open(os.path.join(myDir, file), 'rb') as fileName:
                dict_new_inverted_index = pickle.load(fileName)
                for k, v in dict_new_inverted_index.items():
                    try:

```

```

        main_res_dict_inverted_index[k].extend(v)
    except KeyError:
        main_res_dict_inverted_index[k] = v
    res_main_res_dict = {key: list(set(value)) for key, value in
main_res_dict_inverted_index.items()}
    print("Invrted index length: = ", len(res_main_res_dict))

    # TF-IDF file
    if os.path.isfile('picklefiles/html_parsing_list_pk1.pkl'):
        with open('picklefiles/html_parsing_list_pk1.pkl', 'rb') as
fileName:
            html_parsing_list = pickle.load(fileName)

    # TF-IDF file to tf_idf_dict
    if os.path.isfile('picklefiles/td_idf.pkl'):
        with open('picklefiles/td_idf.pkl', 'rb') as fileName:
            dict_new_tf_idf = pickle.load(fileName)
            for k, v in dict_new_tf_idf.items():
                try:
                    tf_idf_dict[k].extend(v)
                except KeyError:
                    tf_idf_dict[k] = v
    # tf_idf_dict = {key: list(set(value)) for key, value in
main_res_dict_tf_idf.items()}

    # IDF file to idf_dict
    if os.path.isfile('picklefiles/idf.pkl'):
        with open('picklefiles/idf.pkl', 'rb') as fileName:
            dict_new_idf = pickle.load(fileName)
            for k, v in dict_new_idf.items():
                try:
                    idf_dict[k].extend(v)
                except KeyError:
                    idf_dict[k] = v
    #
    # # TF_IDF_OPT dict to tf_idf_dict_opt_ext
    if os.path.isfile('picklefiles/tf_idf_opt.pkl'):
        with open('picklefiles/tf_idf_opt.pkl', 'rb') as fileName:
            dict_new_tf_idf_opt = pickle.load(fileName)
            for k, v in dict_new_tf_idf_opt.items():
                try:
                    tf_idf_dict_opt_ext[k].extend(v)
                except KeyError:
                    tf_idf_dict_opt_ext[k] = v
    #
    # # Query vector and similarity
    for idx in range(len(query_list)):
        stops = ["-", "i", "me", "my", "myself", "we", "our", "ours",
"ourselves", "you", "your", "yours",
                     "yourself", "yourselves", "he", "him", "his", "himself",
"she", "her", "hers", "herself", "it", "its",
                     "itself", "they", "them", "their", "theirs",
"themselves", "what", "which", "who", "whom", "this", "that",
                     "these", "those", "am", "is", "are", "was", "were",
"be", "been", "being", "have", "has", "had", "having",
                     "do", "does", "did", "doing", "a", "an", "the", "and",
"but", "if", "or", "because", "as", "until",

```

```

        "while", "of", "at", "by", "for", "with", "about",
"against", "between", "into", "through", "during",
            "before", "after", "above", "below", "to", "from", "up",
"down", "in", "out", "on", "off", "over", "under",
            "again", "further", "then", "once", "here", "there",
"when", "where", "why", "how", "all", "any", "both",
            "each", "few", "more", "most", "other", "some", "such",
"no", "nor", "not", "only", "own", "same", "so",
            "than", "too", "very", "s", "t", "can", "will", "just",
"don", "should", "now"]
    query_to_check = query_list[idx]
    data = re.sub('[^\w|\_]', ' ', query_to_check) # only keep
numbers and letters and spaces
    # data = data.lower()
    data = re.sub(r'[\x00-\x7f]', r'', data) # remove non ascii
texts
    tokens = [data_upt for data_upt in data.split(' ') if data_upt]
# split the words and remove empty words
    tokens = [word for word in tokens if word.lower() not in stops]
    for token in tokens:
        doc_ids_idf_score = [idf_dict[token][0][1]]
        query_vector_dict[token] = doc_ids_idf_score
    #
    # Cosine similarity
    # The result of the dot product between the query vector and the
documents (per term) need to be converted into a final sorted list of pairs
containing the document identifier and score
    # The sorting is based on highest score to lowest.
    terms = [docs for docs in query_vector_dict]
    terms_idf_val = [query_vector_dict.get(docs)[0] for docs in
query_vector_dict]
    for term_num in range(len(terms)):
        doc_length = []
        cosine_similarity_dict_clone_res = {}
        dict_keys = list(tf_idf_dict[terms[term_num]].keys())
        cosine_similarity_dict = {}
        for idx in range(len(html_parsing_list)):
            if idx not in dict_keys:
                cosine_similarity_dict[idx] = 0
            else:
                cosine_similarity_dict[idx] = terms_idf_val[0] *
tf_idf_dict[terms[term_num]].get(idx)[0]
                doc_length.append(tf_idf_dict_opt_ext.get(idx)[0])
        doc_square = [i ** 2 for i in doc_length]
        doc_length_sqrt_val = pow(sum(doc_square), 0.5)
        for key, values in cosine_similarity_dict.items():
            cosine_similarity_dict_clone_res[key] =
cosine_similarity_dict.get(key) / doc_length_sqrt_val
        for key, values in cosine_similarity_dict_clone_res.items():
            if key not in cosine_similarity_dict_parent and key != 0:
                cosine_similarity_dict_parent[key] = []
            if key in cosine_similarity_dict_parent:
                cosine_similarity_dict_parent[key].append(values)
        for key, values in cosine_similarity_dict_parent.items():
            if key not in cosine_similarity_dict_parent:
                cosine_similarity_dict_parent_2[key] = 0
            if key in cosine_similarity_dict_parent:

```

```

                cosine_similarity_dict_parent_2[key] =
sum(cosine_similarity_dict_parent.get(key))
                cosine_similarity_dict_parent_term_list[terms[term_num]] =
cosine_similarity_dict_clone_res
                print("\n\n\nTop-K ranked results(Query word -> similarity
result):= ",cosine_similarity_dict_parent_term_list)
                cosine_similarity_dict_parent_final =
sorted(cosine_similarity_dict_parent_2.items(), key=lambda t: t[1],
reverse=True)
                print("Top-K ranked results(Top-K Document Numbers):= ",
cosine_similarity_dict_parent_final)
# return json
jsonify("The result is given below (HTML page saved and their K-top
results using similarity")
return jsonify(cosine_similarity_dict_parent_final), 200
else:
    return 'Content-Type not supported!'

if __name__ == '__main__':
    app.run(debug=True)

```

Packages installed in my machine: - (In the case if you need)

appdirs	1.4.4	1.4.4
attrs	21.4.0	21.4.0
automat	20.2.0	20.2.0
bcrypt	3.2.0	3.2.0
beautifulsoup4	4.11.1	4.11.1
blas	1.0	1.0
bottleneck	1.3.4	1.3.4
brotlipy	0.7.0	0.7.0
ca-certificates	2022.3.29	2022.3.29
certifi	2021.10.8	2021.10.8
cffi	1.15.0	1.15.0
charset-normalizer	2.0.4	2.0.4
click	8.1.2	8.0.4
colorama	0.4.4	0.4.4
coloredlogs	15.0.1	15.0.1
constantly	15.1.0	15.1.0
cryptography	36.0.0	36.0.0
cssselect	1.1.0	1.1.0
dataclasses	0.8	0.8
dm-tree	0.1.7	0.1.5
filelock	3.6.0	3.6.0

flask	2.0.3	2.0.3
greenlet	1.1.1	1.1.1
humanfriendly	10.0	10.0
hyperlink	21.0.0	21.0.0
icu	58.2	68.1
idna	3.3	3.3
incremental	21.3.0	21.3.0
intel-openmp	2021.4.0	2022.0.0
itemadapter	0.3.0	0.3.0
itemloaders	1.0.4	1.0.4
itsdangerous	2.0.1	2.0.1
jinja2	3.0.3	3.0.3
jmespath	0.10.0	0.10.0
joblib	1.1.0	1.1.0
libcxx	12.0.0	12.0.0
libffi	3.3	3.4.2
libiconv	1.16	1.16
libxml2	2.9.12	2.9.12
libxslt	1.1.34	1.1.34
lxml	4.8.0	4.8.0
markupsafe	2.0.1	2.0.1
mkl	2021.4.0	2022.0.0
mkl-service	2.4.0	2.4.0
mkl_fft	1.3.1	1.3.1
mkl_random	1.2.2	1.2.2
ncurses	6.3	6.3
nltk	3.7	3.7
numexpr	2.8.1	2.8.1
numpy	1.21.5	1.21.5
numpy-base	1.21.5	1.21.5
openssl	1.1.1n	1.1.1n
packaging	21.3	21.3
pandas	1.4.2	1.4.2
parsel	1.6.0	1.6.0
pillow	9.1.0	9.0.1
pip	21.2.4	21.2.4
protego	0.1.16	0.1.16
pyasn1	0.4.8	0.4.8

pyasn1-modules	0.2.8	0.2.8
pycparser	2.21	2.21
pydispatcher	2.0.5	2.0.5
pyhamcrest	2.0.2	2.0.2
pyopenssl	22.0.0	22.0.0
pyparsing	3.0.4	3.0.4
pysocks	1.7.1	1.7.1
python	3.9.12	3.10.4
python-dateutil	2.8.2	2.8.2
pytz	2021.3	2021.3
queuelib	1.5.0	1.5.0
readline	8.1.2	8.1.2
regex	2022.3.15	2022.3.15
requests	2.27.1	2.27.1
requests-file	1.5.1	1.5.1
scrapy	2.6.1	2.6.1
service_identity	18.1.0	18.1.0
setuptools	61.2.0	61.2.0
six	1.16.0	1.16.0
soupsieve	2.3.1	2.3.1
sqlalchemy	1.4.32	1.4.32
sqlite	3.38.2	3.38.2
svgwrite	1.4.2	
tk	8.6.11	8.6.11
tldextract	3.2.0	3.2.0
tqdm	4.64.0	4.64.0
tree	0.2.4	
twisted	22.2.0	22.2.0
typing_extensions	4.1.1	4.1.1
tzdata	2022a	2022a
urllib3	1.26.9	1.26.9
w3lib	1.21.0	1.21.0
werkzeug	2.0.3	2.0.3
wheel	0.37.1	0.37.1
xz	5.2.5	5.2.5
zlib	1.2.12	1.2.12
zope	1.0	1.0
zope.interface	5.4.0	5.4.0

10 Bibliography:-

1. Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze, 2008 Introduction to Information Retrieval, Cambridge University Press.
2. Donald Metzler, Trevor Strohman, and W. Bruce Croft, 2010, Search Engines Information Retrieval in Practice, Pearson Education, Inc
3. https://en.wikipedia.org/wiki/Information_retrieval