

1. Use the web crawler you built in Project 1 that crawled a limited space, looking for text, html files and php files. You may need to modify how you saved the information from the documents that you traversed to support the query engine. You may assume a maximum of 60 documents will need to be indexed (the number of files fetched will be higher than this). Describe in detail what you changed to support the second half of the project. [15/15 points]

Answer:

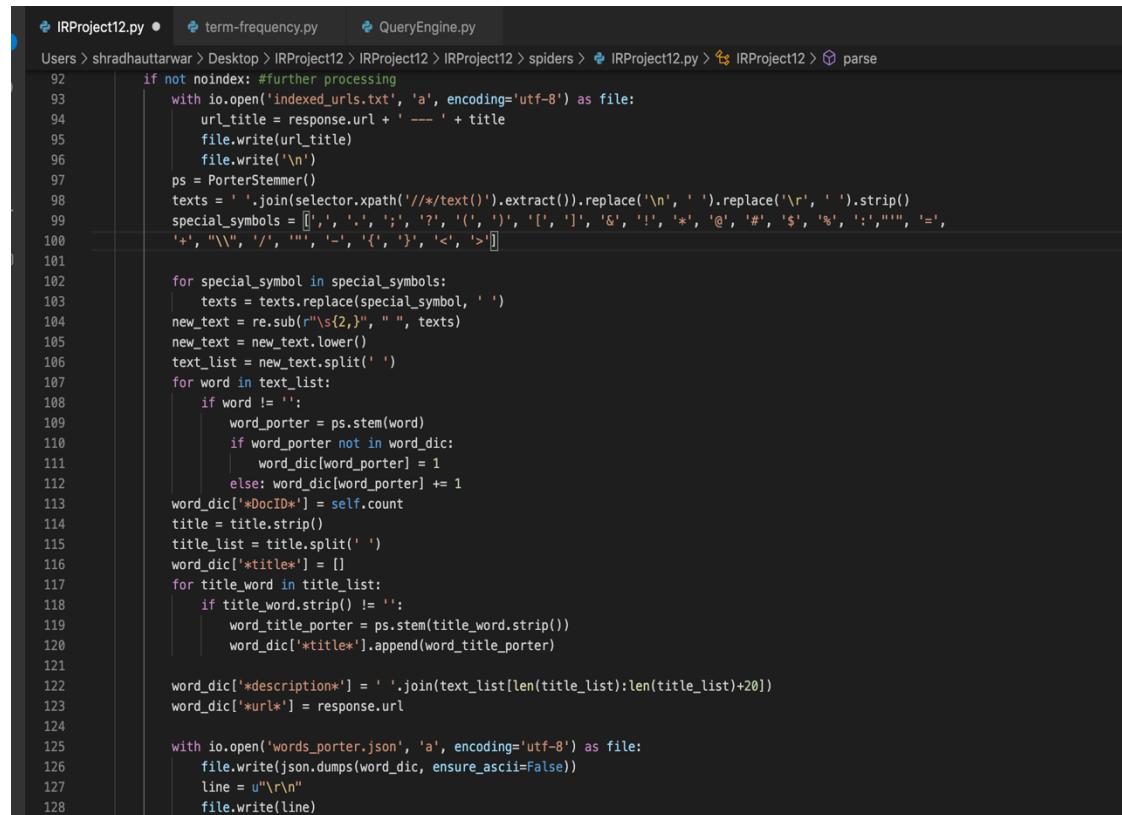
Using the web crawler in Project 1, I indexed 40 pages this time. The URLs of these pages are following:

```
● ● ● indexed_urls.txt
https://s2.smu.edu/~fmoore/ --- Freeman Moore - SMU Spring 2020
https://s2.smu.edu/~fmoore/schedule.htm --- SMU CS 5337/7337 Spring 2020 Schedule
https://s2.smu.edu/~fmoore/textfiles/cow3.txt ---
https://s2.smu.edu/~fmoore/textfiles/cow2.txt ---
https://s2.smu.edu/~fmoore/textfiles/cow1.txt ---
https://s2.smu.edu/~fmoore/misc/noindex.html --- Don't index this page
https://s2.smu.edu/~fmoore/misc/count_letters_duplicate.txt ---
https://s2.smu.edu/~fmoore/misc/seo_email.txt ---
https://s2.smu.edu/~fmoore/textfiles/index.html --- SMU CSE 5/7337 Spring 2018 Textfiles
https://s2.smu.edu/~fmoore/textfiles/extratextfiles/index.php --- SMU CS 5/7337 Spring 2020 text files"
https://s2.smu.edu/~fmoore/misc/vb_byte.c.txt ---
https://s2.smu.edu/~fmoore/misc/levenshtein.html --- Levenshtein Distance demo
https://s2.smu.edu/~fmoore/textfiles/cow4.txt ---
https://s2.smu.edu/~fmoore/textfiles/baseball3.txt ---
https://s2.smu.edu/~fmoore/textfiles/baseball5.txt ---
https://s2.smu.edu/~fmoore/textfiles/baseball4.txt ---
https://s2.smu.edu/~fmoore/textfiles/football4.txt ---
https://s2.smu.edu/~fmoore/textfiles/football3.txt ---
https://s2.smu.edu/~fmoore/textfiles/football2.txt ---
https://s2.smu.edu/~fmoore/textfiles/football1.txt ---
https://s2.smu.edu/~fmoore/textfiles/football5.txt ---
https://s2.smu.edu/~fmoore/textfiles/basketball3.txt ---
https://s2.smu.edu/~fmoore/textfiles/basketball2.txt ---
https://s2.smu.edu/~fmoore/textfiles/baseball2.txt ---
https://s2.smu.edu/~fmoore/textfiles/baseball1.txt ---
https://s2.smu.edu/~fmoore/textfiles/basketball1.txt ---
https://s2.smu.edu/~fmoore/textfiles/basketball5.txt ---
https://s2.smu.edu/~fmoore/textfiles/golf5.txt ---
https://s2.smu.edu/~fmoore/textfiles/basketball4.txt ---
https://s2.smu.edu/~fmoore/textfiles/extratextfiles/building2.txt ---
https://s2.smu.edu/~fmoore/textfiles/extratextfiles/magictext.html --- This is the magic file
https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockingbird5.html --- Mockingbird part 5
https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockingbird4.html --- Mockingbird part 4
https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockingbird3.html --- Mockingbird part 3
https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockingbird2.html --- Mockingbird part 2
https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockingbird1.html --- Mockingbird novel part 1
https://s2.smu.edu/~fmoore/textfiles/golf3.txt ---
https://s2.smu.edu/~fmoore/textfiles/golf2.txt ---
https://s2.smu.edu/~fmoore/textfiles/golf1.txt ---
https://s2.smu.edu/~fmoore/textfiles/golf4.txt ---
https://s2.smu.edu/~fmoore/ --- Freeman Moore - SMU Spring 2020
```

In the second half of the project, We, can modified and add some new data structures and try to crawl indexed documents.

1.1 Web Crawling phase:

Web crawling part for each indexed page, used xpath to extract its title, url, and the 1st of the 20 words of this page and all texts of this page. Therefore, as we can see for every word, we can put it into a dictionary, the key of the dictionary is the word and value of the dictionary if there is a frequency of that word in this document. Then we can add a key called “*title*” and the value of this key “*title*” is the title of this page. Whereas we have also added a key called “*url*” and the value of this key “*url*” is the url of this page. However, we have added a key called “*description*” and the value of this key “*description*” of this page. Alongside assigned a unique document ID for every page, so also added an additional key called “*sequence*” and the value of the key “*sequence*” is the document ID. Then write dictionary in to a intermediate json file “word_porter.json”. Each dictionary in this file represents an indexed page.



```
IRProject12.py • term-frequency.py QueryEngine.py
Users > shradhauttarwar > Desktop > IRProject12 > IRProject12 > spiders > IRProject12.py > IRProject12 > parse
92     if not noindex: #further processing
93         with io.open('indexed_urls.txt', 'a', encoding='utf-8') as file:
94             url_title = response.url + ' --- ' + title
95             file.write(url_title)
96             file.write('\n')
97             ps = PorterStemmer()
98             texts = ' '.join(selector.xpath('//*/text()').extract()).replace('\n', ' ').replace('\r', ' ').strip()
99             special_symbols = [',', '.', ';', '?', '(', ')', '[', ']', '&', '!', '*', '@', '#', '$', '%', ':', "'", '=',
100                '+', '\\\\\", \/\', \"\", \_, \{', '\}', '<', '>']
101
102             for special_symbol in special_symbols:
103                 texts = texts.replace(special_symbol, ' ')
104             new_text = re.sub(r"\s{2,}", " ", texts)
105             new_text = new_text.lower()
106             text_list = new_text.split(' ')
107             for word in text_list:
108                 if word != '':
109                     word_porter = ps.stem(word)
110                     if word_porter not in word_dic:
111                         word_dic[word_porter] = 1
112                     else:
113                         word_dic[word_porter] += 1
114             word_dic['*DocID*'] = self.count
115             title = title.strip()
116             title_list = title.split(' ')
117             word_dic['*title*'] = []
118             for title_word in title_list:
119                 if title_word.strip() != '':
120                     word_title_porter = ps.stem(title_word.strip())
121                     word_dic['*title*'].append(word_title_porter)
122
123             word_dic['*description*'] = ' '.join(text_list[len(title_list):len(title_list)+20])
124             word_dic['*url*'] = response.url
125
126             with io.open('words_porter.json', 'a', encoding='utf-8') as file:
127                 file.write(json.dumps(word_dic, ensure_ascii=False))
128                 line = u"\r\n"
129                 file.write(line)
```

Whereas this file is large so it cannot be shown entirely using screenshot, therefore only put small part of the json file. Below is the result from the “word_porter.Json” file.

```

1 {"freeman": 4, "moor": 4, "smu": 3, "spring": 3, "2020": 3, "u": 1, "phd": 1, "email": 1, "fmoor": 2, "lyle": 2, "ed": 1, "u": 1, "fall": 1, "2019": 1, "cs": 1, "533
2 {"smu": 2, "cs": 2, "5337": 2, "337": 2, "spring": 3, "2020": 2, "schedule": 3, "preliminari": 1, "hi": 1, "page": 2, "is": 1, "maintain": 1, "as": 1, "the": 3, "la
3 {"what": 1, "cow": 2, "is": 1, "a": 1, "brown": 1, "*DocID": 3, "title*": [], "*description*": "cow is a brown cow", "*url*": "https://s2.smu.edu/~fmoore/textfile
4 {"how": 1, "now": 1, "brown": 1, "cow": 1, "*DocID": 4, "title*": [], "*description*": "now brown cow", "*url*": "https://s2.smu.edu/~fmoore/textfiles/cow2.txt",
5 {"how": 1, "brown": 2, "is": 1, "the": 1, "cow": 1, "*DocID": 5, "title*": [], "*description*": "brown is the brown cow", "*url*": "https://s2.smu.edu/~fmoore/tex
6 {"smu": 1, "cse": 1, "5": 5, "337": 1, "spring": 1, "2018": 1, "textfil": 2, "for": 1, "cluster": 1, "golf": 5, "i": 4, "2": 4, "3": 4, "4": 4, "basketbal": 5, "bas
7 {"brown": 1, "cow": 1, "*DocID": 7, "title*": [], "*description*": "cow", "*url*": "https://s2.smu.edu/~fmoore/textfiles/cow4.txt", "*sequences*": 7
8 {"levenshtein": 2, "distan": 5, "demo": 1, "function": 3, "levenshteindist": 2, "a": 6, "b": 6, "var": 10, "cost": 3, "m": 7, "length": 4, "n": 8, "make": 1, "sure"
9 {"includ": 1, "int": 7, "main": 1, "unsign": 3, "*testdata": 3, "12": 1, "222": 1, "654": 1, "36000": 1, "100000": 1, "const": 1, "num_testdata": 1, "sizeof": 2, "2020
10 {"includ": 3, "int": 4, "main": 1, "argc": 1, "char": 2, "argv": 3, "letter": 5, "26": 3, "file": 2, "fin": 5, "i": 10, "onechar": 4, "for": 2, "0": 3, "fopen": 1, "
11 {"smu": 1, "cs": 1, "5": 2, "337": 1, "spring": 1, "2020": 2, "text": 2, "file": 4, "addit": 1, "to": 1, "support": 1, "queri": 1, "implement": 1, "part": 5, "1": 1
12 {"origin": 1, "messag": 1, "from": 2, "atul": 1, "kumar": 1, "malto": 1, "atul939": 1, "outlook": 2, "com": 1, "sent": 2, "thursday": 1, "march": 1, "5": 2, "2020
13 {"after": 1, "24": 1, "win": 1, "season": 4, "nav": 2, "coach": 1, "care": 1, "not": 2, "to": 1, "promis": 1, "too": 1, "muchi": 1, "for": 3, "next": 1, "s": 1, "bask
14 {"terenc": 1, "mann": 1, "ray": 6, "peopl": 4, "will": 5, "come": 5, "they": 16, "11": 10, "to": 3, "iowa": 1, "for": 4, "reason": 1, "can": 1, "t": 2, "even": 2, "f
15 {"the": 3, "golden": 2, "state": 2, "warrior": 1, "basketbal": 2, "team": 1, "had": 1, "someth": 1, "to": 3, "prove": 1, "all": 1, "right": 1, "they": 1, "plan": 1,
16 {"the": 5, "oddsmak": 1, "at": 1, "bettds": 1, "sportsbook": 1, "have": 1, "ohant": 2, "fever": 1, "as": 4, "muchi": 1, "all": 1, "those": 1, "basebal": 2, "scout": 1
17 {"season": 6, "ticket": 2, "sale": 1, "for": 2, "the": 7, "ranger": 3, "are": 1, "off": 1, "after": 1, "team": 2, "s": 1, "78": 1, "84": 1, "record": 1, "last": 3, "
18 {"basketbal": 2, "is": 6, "a": 11, "limit": 1, "contact": 1, "sport": 2, "play": 3, "on": 6, "rectangular": 1, "court": 3, "while": 2, "most": 2, "offem": 2, "as": 1
19 {"it": 4, "s": 6, "open": 3, "day": 6, "and": 4, "for": 4, "one": 3, "our": 1, "team": 1, "is": 1, "go": 3, "to": 8, "win": 2, "the": 7, "whole": 1, "bleep": 1, "thi
20 {"the": 13, "texa": 1, "ranger": 2, "have": 1, "to": 2, "be": 2, "4": 4, "their": 1, "most": 1, "present": 1, "with": 2, "cole": 1, "hamel": 4, "pitch": 2, "and": 5
21 {"uta": 4, "fire": 3, "a": 8, "good": 1, "man": 1, "and": 7, "now": 1, "they": 2, "look": 3, "for": 2, "the": 10, "ark": 1, "of": 2, "coven": 1, "without": 1
22 {"under": 1, "normal": 1, "circumst": 1, "the": 15, "news": 1, "that": 2, "dirk": 3, "nowitzki": 1, "is": 8, "done": 1, "for": 7, "year": 1, "would": 1, "be": 1, "de
23 {"the": 5, "resour": 1, "poulter": 3, "will": 2, "have": 1, "that": 1, "chanc": 1, "sunday": 2, "after": 1, "shoot": 1, "a": 3, "4": 1, "under": 1, "par": 1, "67": 1
24 {"kelli": 1, "kraft": 5, "got": 1, "a": 7, "bird": 1, "just": 1, "not": 2, "the": 8, "one": 1, "he": 1, "want": 1, "s": 2, "ball": 3, "hit": 1, "bird": 2, "im": 3
25 {"bryson": 1, "dechambeau": 3, "shot": 1, "a": 3, "career": 1, "best": 1, "hi": 1, "under": 1, "par": 2, "friday": 1, "to": 2, "take": 1, "the": 7, "lead": 2, "into": 1
26 {"jay": 1, "baa": 2, "nearli": 1, "shot": 1, "hi": 1, "age": 1, "friday": 1, "to": 1, "take": 1, "the": 5, "lead": 5, "into": 1, "a": 3, "36": 1, "hole": 2, "saturda
27 {"the": 7, "champagn": 1, "pop": 1, "echo": 1, "you": 2, "hear": 1, "are": 1, "from": 1, "cb": 1, "and": 6, "pga": 1, "tour": 1, "offici": 1, "celebr": 1, "becaus": 1
28 {"the": 7, "dalla": 2, "cowboy": 4, "overhaul": 1, "their": 1, "group": 1, "of": 4, "assis": 1, "coach": 1, "thi": 1, "off": 1, "season": 1, "richard": 1, "is": 2,
29 [{"buildingon": 1, "buildingtwo": 1, "buildingthre": 1, "*DocID": 30, "title*": [], "*description*": "buildingtwo buildingthree", "*url*": "https://s2.smu.edu/~fmo

```

1.2 Created a dictionary for “words_porter.json”:

“Word_porter.json” automatically read file when the query engine is running. So, create a dictionary called words_documents= {} to save the file. Therefore, the key if this dictionary is the unique document ID which is the value of “*Sequence*”. The value of this dictionary is the dictionary of each document we saved in file “words_porter.json”. This document used to find the title of url, and the 20 words of the document based on its unique document Id.

```

14
15     words_documents = {}
16     document_number = 0
17
18     for line in io.open("words_porter.json", 'r', encoding='utf-8'):
19         new_line = line.replace('\n', '').replace('\r', '').strip()
20         if new_line != '':
21             document_number += 1
22             word_dic = json.loads(new_line)
23             words_documents[document_number] = word_dic

```

1.3 Created a dictionary for “words_frequency_matrix”:

The word frequency matrix will run the query engine with the automatically for that we have created a dictionary called words_matrix= {} to save that. The key of this dictionary in each unique word. The value of this dictionary is a list that contains word frequency in each document, so the length of this list is the number of documents.

```
Users > shradhauttarwar > Desktop > IRProject12 > IRProject12 > QueryEngine.py > ...
24
25     words_matrix = {}
26     for line in io.open("Words_Frequency_Matrix_porter_query.json", 'r', encoding='utf-8'):
27         new_line = line.replace('\n', '').replace('\r', '').strip()
28         if new_line != '':
29             word_dic_temp = json.loads(new_line)
30             for key in word_dic_temp:
31                 words_matrix[key] = word_dic_temp[key]
32
```

1.4 Created a dictionary to save term frequency of each term in each document:

Purpose of creating this dictionary to make calculating cosine similarity conveniently. The key of this dictionary is the unique document ID, and the value of this dictionary is term frequency of each term in this document. Therefore, if we see the example of 40 documents and 1000 words, then this dictionary has 40 keys. The value of each key is list that contains term frequency of these 1000 words in this document, so the length of this is 1000.

```
Users > shradhauttarwar > Desktop > IRProject12 > IRProject12 > QueryEngine.py
32
33     documents_list = []
34     for word in words_matrix.keys():
35         for idx, frequency in enumerate(words_matrix[word]):
36             if idx+1 not in documents_list:
37                 documents_list[idx+1] = []
38                 documents_list[idx+1].append(frequency)
```

1.5 Created a dictionary to save the document frequency of each word:

The document frequency dictionary also makes calculating cosine similarly conveniently. During creating this dictionary, we assign a unique word ID to each word. So, the key of this dictionary is the unique word ID and the value of this dictionary is the document frequency of this word.

```

    IRProject12.py ● term-frequency.py QueryEngine.py ×
Users > shradhauttarwar > Desktop > IRProject12 > IRProject12 > QueryEngine.py > ...
39
40     idf_dict = {}
41     word_number = 0
42     for word in words_matrix.keys():
43         word_number += 1
44         idf_dict[word_number] = document_number - words_matrix[word].count(0) + 1

```

Modification done is the next question.

2. The user will be able to enter multiple queries, consisting of one or more query words separated by space. The single word query “stop” will cause your program to stop. [15/15 points]
 - a) Describe what happens if a user enters a word that is not in the dictionary.
 - b) You may assume the user only enters “words” per your definition.
 - c) Do not recrawl the website for each query, use your collected data.
 - d) CS7337: You will need to implement thesaurus expansion on each term using the list provided at the end.

Answer:

The query engine will run inside a while loop. Therefore, in this while loop, the query engine is waiting for user to enter a query. When the user enters, “stop”, the query engine stops running.

```

Users > shradhauttarwar > Desktop > IRProject12 > IRProject12 > QueryEngine.py
55     while True:
56         if not is_expanded:
57             query = input('\nInput query you wish: ')
58
59             if query == 'stop':
60                 print('\nClosing the query engine...')
61                 break
62             query_list = query.split(' ')

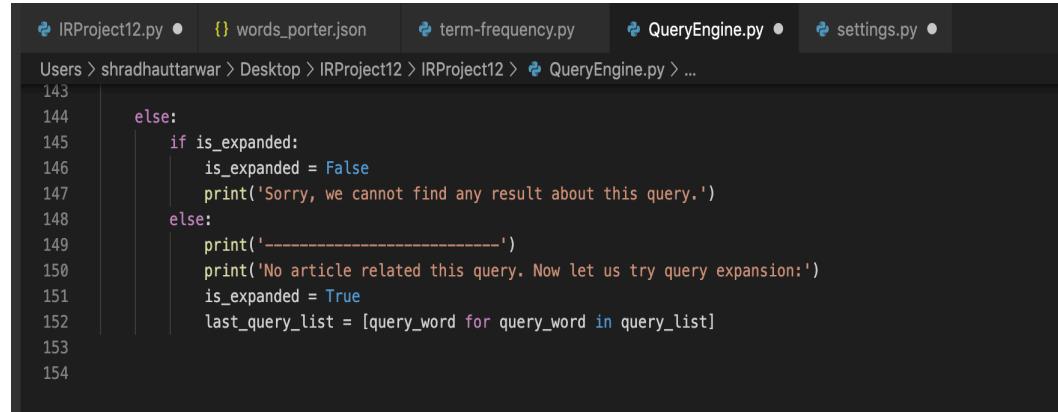
```

```

Input query you wish: stop

Closing the query engine...
(IRProject12) Shradhas-Air:IRProject12 shradhauttarwar$
```

However, we can see the query engine is waiting for user to enter query as long as user doesn't enter “stop”. When user enters “stop”, the query engine stops running.



The screenshot shows a code editor with several tabs at the top: IRProject12.py, words_porter.json, term-frequency.py, QueryEngine.py (which is the active tab), and settings.py. The code in the editor is a Python script for a query engine. It includes logic for handling expanded queries and printing results. The code is as follows:

```
143
144     else:
145         if is_expanded:
146             is_expanded = False
147             print('Sorry, we cannot find any result about this query.')
148         else:
149             print('-----')
150             print('No article related this query. Now let us try query expansion:')
151             is_expanded = True
152             last_query_list = [query_word for query_word in query_list]
153
154
```

Let's see some examples:

```
Input query you wish: smu sjhdjdahfuiia
Score: 0.22278402875996756
URL: https://s2.smu.edu/~fmoore/
Title: freeman moor - smu spring 2020
Description: 2020 freeman l moore phd email fmoore lyle smu ed u fall 2019 cs 5330 7330 keep looking at the
-----
Score: 0.17362363285124685
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/index.php
Title: smu CS 5/7337 spring 2020 text files"
Description: files additional text files to support query implementation part 1 part 2 part 3 part 4 part 5 hocuspocus 3
-----
Score: 0.17337378897914474
URL: https://s2.smu.edu/~fmoore/textfiles/index.html
Title: smu cse 5/7337 spring 2018 textfil
Description: textfiles textfiles for clustering golf 1 golf 2 golf 3 golf 4 golf 5 basketball 1 basketball 2 basketball 3
-----
Score: 0.15615318273648168
URL: https://s2.smu.edu/~fmoore/schedule.htm
Title: smu CS 5337/7337 spring 2020 schedul
Description: schedule smu cs 5337 7337 spring 2020 preliminary schedule this page is maintained as the latest schedule of content and
```

In the above example for query “*smu sjhdjdahfuiia*”, word “*smu*” is in the dictionary but word “*sjhdjdahfuiia*” is not in the dictionary. So, we will only use word “*smu*” to calculate the similar and returns the result.

```

Input query you wish: nowork
-----
No article related this query. Now let us try query expansion:
Sorry, we cannot find any result about this query.

Input query you wish: no work
-----
Score: 0.040508971607146824
URL: https://s2.smu.edu/~fmoore/textfiles/baseball4.txt
Title: No title for this article.
Description: oddsmakers at betdsi sportsbook have ohtani fever as much as all those baseball scouts and writers tripping over themselves as

Score: 0.038654316820208194
URL: https://s2.smu.edu/~fmoore/misc/seo_email.txt
Title: No title for this article.
Description: original message from atul kumar mailto atul1930 outlook com sent thursday march 5 2020 12 58 am subject first page

Score: 0.03464360745416091
URL: https://s2.smu.edu/~fmoore/schedule.htm
Title: smu CS 5337/7337 spring 2020 schedul
Description: schedule smu cs 5337 7337 spring 2020 preliminary schedule this page is maintained as the latest schedule of content and

Score: 0.0334946388512385414
URL: https://s2.smu.edu/~fmoore/textfiles/baseball3.txt
Title: No title for this article.
Description: ticket sales for the rangers are off after the team s 78 84 record last season according to a source

Score: 0.025695692829510352
URL: https://s2.smu.edu/~fmoore/textfiles/basketball2.txt
Title: No title for this article.
Description: fired a good man and now they are looking for the ark of the covenant without indiana jones - they

No article related this query. Now you can reenter Query:
-----
Score: 0.040508971607146824
URL: https://s2.smu.edu/~fmoore/textfiles/baseball4.txt
Title: No title for this article.
Description: oddsmakers at betdsi sportsbook have ohtani fever as much as all those baseball scouts and writers tripping over themselves as

Score: 0.038654316820208194
URL: https://s2.smu.edu/~fmoore/misc/seo_email.txt
Title: No title for this article.
Description: original message from atul kumar mailto atul1930 outlook com sent thursday march 5 2020 12 58 am subject first page

Score: 0.03464360745416091
URL: https://s2.smu.edu/~fmoore/schedule.htm
Title: smu CS 5337/7337 spring 2020 schedul
Description: schedule smu cs 5337 7337 spring 2020 preliminary schedule this page is maintained as the latest schedule of content and

Score: 0.0334946388512385414
URL: https://s2.smu.edu/~fmoore/textfiles/baseball3.txt
Title: No title for this article.
Description: ticket sales for the rangers are off after the team s 78 84 record last season according to a source

Score: 0.025695692829510352
URL: https://s2.smu.edu/~fmoore/textfiles/basketball2.txt
Title: No title for this article.
Description: fired a good man and now they are looking for the ark of the covenant without indiana jones - they

Input query you wish: ■

```

For query “*nowork*”, word “*nowork*” is not in the dictionary, then we try to expand this word to its alternatives word “*no*” and “*work*” use expanded query “*nowork no work*” find articles.

```

Input query you wish: hoehaiuhbsfanck
-----
No article related this query. Now let us try query expansion:
Sorry, we cannot find any result about this query.

Input query you wish: ■

```

In the above query we try to find the word “*hoehaiuhbsfanck*” which we could not find the word so we print the statement that the word is “Sorry, we cannot find any result about this query”.

The user enters “words” per definition we crawled the text from the documents, we stem each word. So, after we get the query, we can split the query into `query_list` and then we stem each query word using Porter stemmer. Therefore, after stemming the query word, we used stemmed query to calculate the cosine similarity.

```

shradhauttarwar > Desktop > IRProject12 > IRProject12 > QueryEngine.py
    ps = PorterStemmer()
    query_list_porter = []
    for original_word in query_list:
        word_porter = ps.stem(original_word)
        query_list_porter.append(word_porter) |

```

Query document will make query as a document and calculated term frequency for each term in this query document and add this frequency list into dictionary document_list= {} (A dictionary that contain term frequency for each term in each document) for calculating cosine similarity. The key for query in the dictionary is assigned as 0.

```

Users > shradhauttarwar > Desktop > IRProject12 > IRProject12 > QueryEngine.py > ...
73
74     for matrix_word in words_matrix.keys():
75         words_matrix[matrix_word].append(0)
76
77     word_flag = False
78     for porter_word in query_list_porter:
79         for matrix_word in words_matrix.keys():
80             if matrix_word == porter_word:
81                 word_flag = True
82                 words_matrix[matrix_word][-1] += 1
83
84     if word_flag:
85         documents_list[0] = []
86         for matrix_word in words_matrix.keys():
87             documents_list[0].append(words_matrix[matrix_word][-1])
88

```

We have implemented thesaurus expansion on each term and then searching term documents. The words_expansion has in place regardless of its being found in document or not to be. Therefore, the words are alternative in the dictionary words expansion as shown in below:

```

words_expansion = {'beautiful': ['nice', 'fancy'], 'chapter': ['chpt'],
'chpt': ['chapter'], 'responsible': ['owner', 'accountable'], 'freemanmoore':
['freeman', 'moore'], 'dept': ['department'], 'photo': ['photograph', 'image',
'picture'], 'brown': ['beige', 'tan', 'auburn'], 'tues': ['Tuesday'], 'sole':
['owner', 'single', 'shoe', 'boot'], 'homework': ['hmwk', 'home', 'work'],
'novel': ['book', 'unique'], 'computer': ['cse'], 'story': ['novel', 'book'],
'hocuspocus': ['magic', 'abracadabra'], 'thisworks': ['this', 'work'], }

```

```

else:
    query_list = [query_word for query_word in last_query_list]
    for query_word in last_query_list:
        if words_expansion.get(query_word):
            for expanded_word in words_expansion[query_word]:
                query_list.append(expanded_word)

```

3. Implement the cosine similarity of the query against each document. [40/40 points]

- a) What document/query weighting scheme did you implement (suggestion: lnc.ltc)?

In weighting scheme used lnc.ltc to calculate cosine similarity. For document, using Inc scheme and for query, using ltc scheme. The term frequency of each term in each document and query is stored at dictionary document_list. The key of the dictionary is the unique document ID and special key 0 represents query. Then we use lnc scheme to process document and use ltc to process query. After that, we do product between query and calculate the cosine similarity. Then we store the results into the dictionary similarity_score_dict= {}. The key of this dictionary is unique document ID and the value of this dictionary is the cosine similarity between query and this document.

```

Users > shradhauttarwar > Desktop > IRProject12 > IRProject12 > QueryEngine.py > ...
83
84     if word_flag:
85         documents_list[0] = []
86         for matrix_word in words_matrix.keys():
87             documents_list[0].append(words_matrix[matrix_word][-1])
88
89         similarity_scores_dict = {}
90         for i in range(1, document_number+1):
91             query_score = []
92             query_sum = 0
93             for idx, tf in enumerate(documents_list[0]):
94                 if tf > 0:
95                     tf_temp = (1 + math.log(tf)) * math.log(document_number/(df_dict[idx+1]+1))
96                 else:
97                     tf_temp = 0
98                     query_sum += tf_temp * tf_temp
99                     query_score += [tf_temp]
100                     query_normalized = math.sqrt(query_sum)
101                     query_score_processed = np.array([score/(query_normalized+1) for score in query_score])
102
103                     document_score = []
104                     document_sum = 0
105                     for tf in documents_list[i]:
106                         if tf > 0:
107                             tf_temp = (1 + math.log(tf))
108                         else:
109                             tf_temp = 0
110                             document_sum += tf_temp * tf_temp
111                             document_score += [tf_temp]
112                             document_normalized = (math.sqrt(document_sum) )
113                             document_score_processed = np.array([score/(document_normalized+1) for score in document_score])

```

- b) If any of the query words appear in the <title> of a document, add 0.1 to its score.**

Adding score for title we have to remember can access the title of each dictionary words_documents= {} using unique document ID. For example, words_documents[1]['*title*'] represents the title of document 1. By doing this way we determined the query word appears in the title of each documents. If it does, then we add 0.1 to the cosine similarity score between query and this document.

```
Users > shradhauttarwar > Desktop > IRProject12 > IRProject12 > QueryEngine.py > ...
116
117     query_in_title = False
118     for porter_word in query_list_porter:
119         if porter_word in words_documents[i]['*title*']:
120             query_in_title = True
121     if query_in_title:
122         similarity_score += 0.1
123     similarity_scores_dict[i] = similarity_score
124 similarity_scores_dict_new = sorted(similarity_scores_dict.items(), key=lambda d: d[1], reverse=True)
```

- c) Display the resulting score, document URL, and document title in descending numerical order for the top K=5 results.**

```
similarity_scores_dict_new = sorted(similarity_scores_dict.items(), key=lambda d: d[1], reverse=True)
result_count = 0
for number, score in similarity_scores_dict_new[:5]:
    if score > 0:
        result_count += 1
        print('-----')
        print('Score: ', score)
        print('URL: ', words_documents[number]['*url*'])
        if not words_documents[number]['*title*']:
            title = 'No title for this article.'
        else:
            title = ' '.join(words_documents[number]['*title*'])
        print('Title: ', title)
        print('Description: ', words_documents[number]['*description*'])
    if not is_expanded:
        print('-----')
        print('No article related this query. Now you can reenter Query:')
        is_expanded = True
last_query_list = [query_word for query_word in query_list]
```

Displaying the top 5 results after we get cosine similarity scores in dictionary similarity_score_dict= {}, we then sort this dictionary in descending order based on its cosine similarity score and print the score, title, url and the first 20 words of the top 5 documents. Remember we can access the title, url and the first 20 words of each document in dictionary words_documents= {} using unique document ID(word_documents[1]['*title*'],words_documents[1]['*url*'], words_documents[1]['*description*']).

Result shows below:

```

Input query you wish: three year story
-----
Score: 0.0969805355852488
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird4.html
Title: mockingbird part 4
Description: mockingbird part 4 the story takes place during three years 1933-35 of the great depression in the fictional tired old
-----
Score: 0.09447253232678378
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird1.html
Title: mockingbird novel part 1
Description: mockingbird part 1 to kill a mockingbird is primarily a novel about growing up under extraordinary circumstances in the 1930s
-----
Score: 0.06320967274651282
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird2.html
Title: mockingbird part 2
Description: mockingbird part 2 the only neighbor who puzzles them is the mysterious arthur radley nicknamed boo who never comes outside
-----
Score: 0.05239093230701331
URL: https://s2.smu.edu/~fmoore/textfiles/basketball5.txt
Title: No title for this article.
Description: is a limited contact sport played on a rectangular court while most often played as a team sport with five
-----
Score: 0.04968828394775876
URL: https://s2.smu.edu/~fmoore/textfiles/basketball1.txt
Title: No title for this article.
Description: normal circumstances the news that dirk nowitzki is done for the year would be devastating but thursday s announcement was

```

The query is “*three year story*”, and we print the score, url, title and top 5 documents based on their cosine similarity score in descending order.

d) CS7337: Display the first 20 words of the document (this can be the stemmed version).

In order to save the first 20 words of the documents, we saved the first 20 words of the document into file ‘words_porter.json’ during web crawling phase. This file will then be read into dictionary words_documents= {} during query phase. We can access to the 20 words of the documents using unique document ID such as words_documents [1] [*description*]. Then we print the description when we show the results given the query.

```

IRProject12.py ● { } words_porter.json ● term-frequency.py ● QueryEngine.py ● settings.py ●
Users > shradhauitarwar > Desktop > IRProject12 > IRProject12 > QueryEngine.py > ...
124     similarity_scores_dict_new = sorted(similarity_scores_dict.items(), key=lambda d: d[1], reverse=True)
125     result_count = 0
126     for number, score in similarity_scores_dict_new[:5]:
127         if score > 0:
128             result_count += 1
129             print('-----')
130             print('Score: ', score)
131             print('URL: ', words_documents[number]['*url*'])
132             if not words_documents[number]['*title*']:
133                 title = 'No title for this article.'
134             else:
135                 title = ' '.join(words_documents[number]['*title*'])
136                 print('Title: ', title)
137                 print('Description: ', words_documents[number]['*description*'])
138             if not is_expanded:
139                 print('-----')
140                 print('No article related this query. Now you can reenter Query:')
141                 is_expanded = True |
142                 last_query_list = [query_word for query_word in query_list]
143             else:
144                 is_expanded = False
145                 last_query = []

```

For example, of 1st 5 documents and 20 words:

```

Input query you wish: three year story
-----
Score: 0.09698605355852488
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird4.html
Title: mockingbird part 4
Description: mockingbird part 4 the story takes place during three years 1933-35 of the great depression in the fictional tired old

-----
Score: 0.09447253232678378
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird1.html
Title: mockingbird novel part 1
Description: mockingbird part 1 to kill a mockingbird is primarily a novel about growing up under extraordinary circumstances in the 1930s

-----
Score: 0.06328967274651282
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird2.html
Title: mockingbird part 2
Description: mockingbird part 2 the only neighbor who puzzles them is the mysterious arthur radley nicknamed boo who never comes outside

-----
Score: 0.05239093230701331
URL: https://s2.smu.edu/~fmoore/textfiles/basketball5.txt
Title: No title for this article.
Description: is a limited contact sport played on a rectangular court while most often played as a team sport with five

-----
Score: 0.04968823394775876
URL: https://s2.smu.edu/~fmoore/textfiles/basketball1.txt
Title: No title for this article.
Description: normal circumstances the news that dirk nowitzki is done for the year would be devastating but thursday s announcement was

```

As we can see the query “*three year story*”, we can see it returns 5 results and there is a description (first 20 words) of each results.

4. Explain why you believe the results are correct. [15/15 points]

You must show the results of these queries (tentative):

a) **moore smu:**

```

|Input query you wish: moore smu
-----
Score: 0.31503030217522243
URL: https://s2.smu.edu/~fmoore/
Title: freeman moor - smu spring 2020
Description: 2020 freeman l moore phd email fmoore lyle smu ed u fall 2019 cs 5330 7330 keep looking at the

Score: 0.1487599845963782
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/index.php
Title: smu CS 5/7337 spring 2020 text files"
Description: files additional text files to support query implementation part 1 part 2 part 3 part 4 part 5 hocuspocus 3

Score: 0.14859451091675568
URL: https://s2.smu.edu/~fmoore/textfiles/index.html
Title: smu cse 5/7337 spring 2018 textfil
Description: textfiles textfiles for clustering golf 1 golf 2 golf 3 golf 4 golf 5 basketball 1 basketball 2 basketball 3

Score: 0.13718953030748318
URL: https://s2.smu.edu/~fmoore/schedule.htm
Title: smu CS 5337/7337 spring 2020 schedul
Description: schedule smu cs 5337 7337 spring 2020 preliminary schedule this page is maintained as the latest schedule of content and

No article related this query. Now you can reenter Query:

Score: 0.31503030217522243
URL: https://s2.smu.edu/~fmoore/
Title: freeman moor - smu spring 2020
Description: 2020 freeman l moore phd email fmoore lyle smu ed u fall 2019 cs 5330 7330 keep looking at the

Score: 0.1487599845963782
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/index.php
Title: smu CS 5/7337 spring 2020 text files"
Description: files additional text files to support query implementation part 1 part 2 part 3 part 4 part 5 hocuspocus 3

Score: 0.14859451091675568
URL: https://s2.smu.edu/~fmoore/textfiles/index.html
Title: smu cse 5/7337 spring 2018 textfil
Description: textfiles textfiles for clustering golf 1 golf 2 golf 3 golf 4 golf 5 basketball 1 basketball 2 basketball 3

Score: 0.13718953030748318
URL: https://s2.smu.edu/~fmoore/schedule.htm
Title: smu CS 5337/7337 spring 2020 schedul
Description: schedule smu cs 5337 7337 spring 2020 preliminary schedule this page is maintained as the latest schedule of content and

Input query you wish: |

```

Query “*moore smu*” returns 4 results and after 4 results it shows another result which is expanded query. The expanded query returns 4 results too, but which is similar result. Therefore, there are only 4 documents that are relevant to the query “*moore smu*”.

For 1st document:

<https://s2.smu.edu/~fmoore/>, moore and smu appears in both the title and the content of this documents.

For 2nd document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/index.php>

appears both in the title and the content of this document.

For 3rd and 4th documents:

<https://s2.smu.edu/~fmoore/textfiles/index.html>

<https://s2.smu.edu/~fmoore/schedule.htm>

smu appears in the title of the document

b) Bob Ewell where Scout

Result for “Bob Ewell Scout”:

```

Input query you wish: Bob Ewell where Scout
-----
Score: 0.172887913202098
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird5.html
Title: mockingbird part 5
Description: mockingbird part 5 atticus does not want jem and scout to be present at tom robinson s trial no seat

-----
Score: 0.064930001283103
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird4.html
Title: mockingbird part 4
Description: mockingbird part 4 the story takes place during three years 1933-35 of the great depression in the fictional tired old

-----
Score: 0.06458299614527417
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird1.html
Title: mockingbird novel part 1
Description: mockingbird part 1 to kill a mockingbird is primarily a novel about growing up under extraordinary circumstances in the 1930s

-----
Score: 0.051609001028916655
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird3.html
Title: mockingbird part 3
Description: mockingbird part 3 suddenly scout and jem have to tolerate a barrage of racial slurs and insults because of atticus

-----
Score: 0.04352345183907334
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird2.html
Title: mockingbird part 2
Description: mockingbird part 2 the only neighbor who puzzles them is the mysterious arthur radley nicknamed boo who never comes outside

```

For this Query 5 results are returned:

For 1st document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird5.html>

Bob appears 3 times whereas Ewell and Scout appears 6 times in the content and where appears once in the content.

For 2nd document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird4.html>

where appears once and Scout appears 2 times.

For 3rd document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird1.html>

Ewell and Scout appears in the content

For 4th document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird3.html>

After visiting this link, we can see Scout appears 4 times in the content.

For 5th document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird2.html>

After visiting this link, we can see Scout appears 3 times in the content.

c) three year story

```

Input query you wish: three year story
-----
Score: 0.09698605355852488
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird4.html
Title: mockingbird part 4
Description: mockingbird part 4 the story takes place during three years 1933-35 of the great depression in the fictional tired old
-----
Score: 0.09447253232678378
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird1.html
Title: mockingbird novel part 1
Description: mockingbird part 1 to kill a mockingbird is primarily a novel about growing up under extraordinary circumstances in the 1930s
-----
Score: 0.06320967274651282
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird2.html
Title: mockingbird part 2
Description: mockingbird part 2 the only neighbor who puzzles them is the mysterious arthur radley nicknamed boo who never comes outside
-----
Score: 0.05239893230701331
URL: https://s2.smu.edu/~fmoore/textfiles/basketball5.txt
Title: No title for this article.
Description: is a limited contact sport played on a rectangular court while most often played as a team sport with five
-----
Score: 0.0496882839477876
URL: https://s2.smu.edu/~fmoore/textfiles/basketball1.txt
Title: No title for this article.
Description: normal circumstances the news that dirk nowitzki is done for the year would be devastating but thursday s announcement was

```

Result shows 5 document for this query.

For 1st document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird4.html>

Query word Appears in this content.

For 2nd document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird1.html>

Query word Appears in this content.

For 3rd document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird2.html>

Query word Appears in this content.

For 4th document:

<https://s2.smu.edu/~fmoore/textfiles/basketball5.txt>

Query word Appears in this content.

For 5th document:

<https://s2.smu.edu/~fmoore/textfiles/basketball1.txt>

Query word Appears in this content.

- d) Atticus to defend Maycomb

```

Input query you wish: Atticus to defend Maycomb
-----
Score: 0.17256190070635355
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird4.html
Title: mockingbird part 4
Description: mockingbird part 4 the story takes place during three years 1933-35 of the great depression in the fictional tired old

-----
Score: 0.13205000158378427
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird1.html
Title: mockingbird novel part 1
Description: mockingbird part 1 to kill a mockingbird is primarily a novel about growing up under extraordinary circumstances in the 1930s

-----
Score: 0.07255591397566555
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird5.html
Title: mockingbird part 5
Description: mockingbird part 5 atticus does not want jem and scout to be present at tom robinson s trial no seat

-----
Score: 0.04432221289342825
URL: https://s2.smu.edu/~fmoore/textfiles/baseball2.txt
Title: No title for this article.
Description: s opening day and for one day our team is going to win the whole bleeping thing it s opening

-----
Score: 0.041069353485050086
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird3.html
Title: mockingbird part 3
Description: mockingbird part 3 suddenly scout and jem have to tolerate a barrage of racial slurs and insults because of atticus

```

Result shows 5 document for the query:

For 1st document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird4.html>

Query word Appears in this content.

For 2nd document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird1.html>

Query word Appears in this content.

For 3rd document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird5.html>

Query word Appears in this content.

For 4th document:

<https://s2.smu.edu/~fmoore/textfiles/baseball2.txt>

Query word Appears in this content.

For 5th document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird3.html>

Query word Appears in this content.

e) hocuspocus thisworks

```

Input query you wish: hocuspocus thiswork
-----
Score: 0.08189520556775087
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/index.php
Title: smu CS 5/7337 spring 2020 text files"
Description: files additional text files to support query implementation part 1 part 2 part 3 part 4 part 5 hocuspocus 3

-----
No article related this query. Now you can reenter Query:

```

Query “hocuspocus thisworks” returns 1 result which is less then 5 results. So, we try to expand the query “hocuspocus thisworks”

```
-----  
Score: 0.27756189592705854  
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/magictext.html  
Title: thi is the magic file  
Description: magic shows up here and in the title brown      beige    tan    auburn  
-----  
Score: 0.06682754269175717  
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/index.php  
Title: smu CS 5/7337 spring 2020 text files"  
Description: files additional text files to support query implementation part 1 part 2 part 3 part 4 part 5 hocuspocus 3  
-----  
Score: 0.02942547213772617  
URL: https://s2.smu.edu/~fmoore/textfiles/baseball5.txt  
Title: No title for this article.  
Description: mann ray people will come ray they ll come to iowa for reasons they can t even fathom they ll  
  
Input query you wish: ■
```

The expanded query “hocuspocus thisworks magic abracadabra this work” returns 3 results this time.

For 1st document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/magictext.html>

Query word Appears in this content.

For 2nd document:

<https://s2.smu.edu/~fmoore/textfiles/extratextfiles/index.php>

Query word Appears in this content.

For 3rd document:

<https://s2.smu.edu/~fmoore/textfiles/baseball5.txt>

Query word Appears in this content.

5. Instructor testing with additional queries. [15/15 points]

Answer: You can also check, just follow the “*Readme file*”.

1. **beautiful:** This query result shows below also shows the expanded result.

```
Input query you wish: beautiful
-----
Score: 0.05266360012571011
URL: https://s2.smu.edu/~fmoore/textfiles/baseball1.txt
Title: No title for this article.
Description: texas rangers have to be at their most presentable with cole hamels pitching and if thursday is the best version
-----
No article related this query. Now you can reenter Query:
-----
Score: 0.05176960701656745
URL: https://s2.smu.edu/~fmoore/textfiles/golf5.txt
Title: No title for this article.
Description: resurgent poult will have that chance sunday after shooting a 4 under par 67 to take a one shot lead
-----
Score: 0.04020464212013578
URL: https://s2.smu.edu/~fmoore/textfiles/baseball1.txt
Title: No title for this article.
Description: texas rangers have to be at their most presentable with cole hamels pitching and if thursday is the best version
-----
Input query you wish: ■
```

URL: <https://s2.smu.edu/~fmoore/textfiles/baseball1.txt>

This URL content word beautiful.

2. **chapter:** There is not result for chapter therefore the expanded query ‘*chpt*’ document is shown below,

```
Input query you wish: chapter
-----
No article related this query. Now let us try query expansion:
-----
Score: 0.1404520810435768
URL: https://s2.smu.edu/~fmoore/schedule.htm
Title: smu CS 5337/7337 spring 2020 schedul
Description: schedule smu cs 5337 7337 spring 2020 preliminary schedule this page is maintained as the latest schedule of content and
```

URL: <https://s2.smu.edu/~fmoore/schedule.htm>

3. **chpt:** This query result shows below only 1 document have shown result.

URL : <https://s2.smu.edu/~fmoore/schedule.htm>

```
Input query you wish: chpt
-----
Score: 0.1342486096984444
URL: https://s2.smu.edu/~fmoore/schedule.htm
Title: smu CS 5337/7337 spring 2020 schedul
Description: schedule smu cs 5337 7337 spring 2020 preliminary schedule this pa
ge is maintained as the latest schedule of content and
-----
No article related this query. Now you can reenter Query:
-----
Score: 0.1342486096984444
URL: https://s2.smu.edu/~fmoore/schedule.htm
Title: smu CS 5337/7337 spring 2020 schedul
Description: schedule smu cs 5337 7337 spring 2020 preliminary schedule this pa
ge is maintained as the latest schedule of content and
-----
Input query you wish: ■
```

Extended Url content the same result so only one document is relevant.

4. **responsible:** This query result shows below also shows the expanded result.

```
Input query you wish: responsible
-----
Score: 0.0684465946854374
URL: https://s2.smu.edu/~fmoore/
Title: freeman moor - smu spring 2020
Description: 2020 freeman l moore phd email fmoore lyle smu ed u fall 2019 cs 5330 7330 keep looking at the
-----
Score: 0.05778459443945048
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird5.html
Title: mockingbird part 5
Description: mockingbird part 5 atticus does not want jem and scout to be present at tom robinson s trial no seat
-----
Score: 0.04079005947397189
URL: https://s2.smu.edu/~fmoore/misc/seo_email.txt
Title: No title for this article.
Description: original message from atul kumar mailto atul1930 outlook com sent thursday march 5 2020 12 58 am subject first page
-----
No article related this query. Now you can reenter Query:
-----
Score: 0.04344653371691458
URL: https://s2.smu.edu/~fmoore/textfiles/baseball2.txt
Title: No title for this article.
Description: s opening day and for one day our team is going to win the whole bleeping thing it s opening
-----
Score: 0.041767611115076105
URL: https://s2.smu.edu/~fmoore/
Title: freeman moor - smu spring 2020
Description: 2020 freeman l moore phd email fmoore lyle smu ed u fall 2019 cs 5330 7330 keep looking at the
-----
Score: 0.03992821235917886
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird5.html
Title: mockingbird part 5
Description: mockingbird part 5 atticus does not want jem and scout to be present at tom robinson s trial no seat
-----
Score: 0.028185265858824873
URL: https://s2.smu.edu/~fmoore/misc/seo_email.txt
Title: No title for this article.
Description: original message from atul kumar mailto atul1930 outlook com sent thursday march 5 2020 12 58 am subject first page
-----
Input query you wish: ■
```

5. **freemanmoore:** There is not result for query “freemanmoore” therefore the expanded query ‘*freeman moore*’ document is shown below:

```
Input query you wish: freemanmoore
-----
No article related this query. Now let us try query expansion:
-----
Score: 0.345059952060554
URL: https://s2.smu.edu/~fmoore/
Title: freeman moor - smu spring 2020
Description: 2020 freeman l moore phd email fmoore lyle smu ed u fall 2019 cs 5330 7330 keep looking at the
```

6. dept: There is no result shown, not even expanded query.

```
Input query you wish: dept
-----
No article related this query. Now let us try query expansion:
Sorry, we cannot find any result about this query.
```

7. photo: There is no result shown, not even expanded query.

```
Input query you wish: photo
-----
No article related this query. Now let us try query expansion:
Sorry, we cannot find any result about this query.
```

8. **brown:** This query result shows below, there are 5 documents resulted related to brown.

```
Input query you wish: brown
-----
Score: 0.3048160515790603
URL: https://s2.smu.edu/~fmoore/textfiles/cow1.txt
Title: No title for this article.
Description: brown is the brown cow
-----
Score: 0.2699788495141498
URL: https://s2.smu.edu/~fmoore/textfiles/cow4.txt
Title: No title for this article.
Description: cow
-----
Score: 0.21726220001698174
URL: https://s2.smu.edu/~fmoore/textfiles/cow2.txt
Title: No title for this article.
Description: now brown cow
-----
Score: 0.18002927038998093
URL: https://s2.smu.edu/~fmoore/textfiles/cow3.txt
Title: No title for this article.
Description: cow is a brown cow
-----
Score: 0.03213994433135377
URL: https://s2.smu.edu/~fmoore/schedule.htm
Title: smu CS 5337/7337 spring 2020 schedule
Description: schedule smu cs 5337 7337 spring 2020 preliminary schedule this page is maintained as the latest schedule of content and
-----
No article related this query. Now you can reenter Query:
-----
Score: 0.3048160515790603
URL: https://s2.smu.edu/~fmoore/textfiles/cow1.txt
Title: No title for this article.
Description: brown is the brown cow
-----
Score: 0.2699788495141498
URL: https://s2.smu.edu/~fmoore/textfiles/cow4.txt
Title: No title for this article.
Description: cow
-----
Score: 0.21726220001698174
URL: https://s2.smu.edu/~fmoore/textfiles/cow2.txt
Title: No title for this article.
Description: now brown cow
-----
Score: 0.18002927038998093
URL: https://s2.smu.edu/~fmoore/textfiles/cow3.txt
Title: No title for this article.
Description: cow is a brown cow
-----
Score: 0.03213994433135377
URL: https://s2.smu.edu/~fmoore/schedule.htm
Title: smu CS 5337/7337 spring 2020 schedule
Description: schedule smu cs 5337 7337 spring 2020 preliminary schedule this page is maintained as the latest schedule of content and
-----
Input query you wish: ■
```

9. **tues:** There is not result for query “*tues*” therefore the expanded query ‘*Tuesday*’ document is shown below:

```
Input query you wish: tues
-----
No article related this query. Now let us try query expansion:
-----
Score: 0.05935465695968675
URL: https://s2.smu.edu/~fmoore/textfiles/football4.txt
Title: No title for this article.
Description: cowboys vice president stephen jones commented on the release of
eceiver dez bryant for the first time tuesday jones said
-----
Input query you wish: ■
```

10. **sole:** This query result shows below also shows the expanded result.

```

Input query you wish: sole
-----
Score: 0.06508050313047531
URL: https://s2.smu.edu/~fmoore/
Title: freeman moor - smu spring 2020
Description: 2020 freeman l moore phd email fmoore lyle smu ed u fall 2019 cs 5330 7330 keep looking at the
-----
No article related this query. Now you can reenter Query:
-----
Score: 0.0496840054522204
URL: https://s2.smu.edu/~fmoore/
Title: freeman moor - smu spring 2020
Description: 2020 freeman l moore phd email fmoore lyle smu ed u fall 2019 cs 5330 7330 keep looking at the
-----
Score: 0.03962132638954409
URL: https://s2.smu.edu/~fmoore/textfiles/baseball2.txt
Title: No title for this article.
Description: s opening day and for one day our team is going to win the whole bleeping thing it s opening
Input query you wish: ■

```

- 11. homework:** There is not result for query “*homework*” therefore the expanded query ‘*hmwk home work*’ document is shown below:

```

Input query you wish: homework
-----
No article related this query. Now let us try query expansion:
-----
Score: 0.10537658283054217
URL: https://s2.smu.edu/~fmoore/schedule.htm
Title: smu CS 5337/7337 spring 2020 schedule
Description: schedule smu cs 5337 7337 spring 2020 preliminary schedule this page is maintained as the latest schedule of content and
-----
Score: 0.03445883261389325
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird5.html
Title: mockingbird part 5
Description: mockingbird part 5 atticus does not want jem and scout to be present at tom robinson s trial no seat
-----
Score: 0.031143142489181604
URL: https://s2.smu.edu/~fmoore/misc/seo_email.txt
Title: No title for this article.
Description: original message from atul kumar mailto atul1930 outlook com sent thursday march 5 2020 12 58 am subject first page
-----
Score: 0.02353328868919592
URL: https://s2.smu.edu/~fmoore/textfiles/baseball1.txt
Title: No title for this article.
Description: texas rangers have to be at their most presentable with cole hamels pitching and if thursday is the best version
-----
Score: 0.022722136742997865
URL: https://s2.smu.edu/~fmoore/textfiles/basketball1.txt
Title: No title for this article.
Description: normal circumstances the news that dirk nowitzki is done for the year would be devastating but thursday s announcement was
Input query you wish: ■

```

- 12. novel:** This query result shows below and also shows the expanded result.

```

Input query you wish: novel
-----
Score: 0.1948099458782001
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird1.html
Title: mockingbird novel part 1
Description: mockingbird part 1 to kill a mockingbird is primarily a novel about growing up under extraordinary circumstances in the 1930s
-----
No article related this query. Now you can reenter Query:
-----
Score: 0.1948099458782001
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird1.html
Title: mockingbird novel part 1
Description: mockingbird part 1 to kill a mockingbird is primarily a novel about growing up under extraordinary circumstances in the 1930s

```

- 13. computer:** There is not result for query “*computer*” therefore the expanded query ‘*cse*’ document is shown below:

```

Input query you wish: computer
-----
No article related this query. Now let us try query expansion:
-----
Score: 0.18161728298182483
URL: https://s2.smu.edu/~fmoore/textfiles/index.html
Title: smu cse 5/7337 spring 2018 textfil
Description: textfiles textfiles for clustering golf 1 golf 2 golf 3 golf 4 golf 5 basketball 1 basketball 2 basketball 3

Input query you wish: ■

```

14. story: This query result shows below and also shows the expanded result.

```

Input query you wish: story
-----
Score: 0.05034007174288458
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird1.html
Title: mockingbird novel part 1
Description: mockingbird part 1 to kill a mockingbird is primarily a novel about growing up under extraordinary circumstances in the 1930s

-----
Score: 0.04670248852855358
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird2.html
Title: mockingbird part 2
Description: mockingbird part 2 the only neighbor who puzzles them is the mysterious arthur radley nicknamed boo who never comes outside

-----
Score: 0.03412274563928915
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird4.html
Title: mockingbird part 4
Description: mockingbird part 4 the story takes place during three years 1933-35 of the great depression in the fictional tired old

-----
Score: 0.033033545254433264
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird5.html
Title: mockingbird part 5
Description: mockingbird part 5 atticus does not want jem and scout to be present at tom robinson s trial no seat

-----
No article related this query. Now you can reenter Query:
-----
Score: 0.2149687382467323
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird1.html
Title: mockingbird novel part 1
Description: mockingbird part 1 to kill a mockingbird is primarily a novel about growing up under extraordinary circumstances in the 1930s

-----
Score: 0.030930457151792632
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird2.html
Title: mockingbird part 2
Description: mockingbird part 2 the only neighbor who puzzles them is the mysterious arthur radley nicknamed boo who never comes outside

-----
Score: 0.02259905532126555
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird4.html
Title: mockingbird part 4
Description: mockingbird part 4 the story takes place during three years 1933-35 of the great depression in the fictional tired old

-----
Score: 0.02187769192297675
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/mockbird5.html
Title: mockingbird part 5
Description: mockingbird part 5 atticus does not want jem and scout to be present at tom robinson s trial no seat

Input query you wish: ■

```

15. hocuspocus: This query result shows below and also shows the expanded result.

```
Input query you wish: hocuspocus
-----
Score: 0.08189520556775087
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/index.php
Title: smu CS 5/7337 spring 2020 text files"
Description: files additional text files to support query implementation part 1 part 2 part 3 part 4 part 5 hocuspocus 3
-----
No article related this query. Now you can reenter Query:
-----
Score: 0.27756189592705854
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/magictext.html
Title: thi is the magic file
Description: magic shows up here and in the title brown      beige    tan    auburn
-----
Score: 0.06602754269175717
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/index.php
Title: smu CS 5/7337 spring 2020 text files"
Description: files additional text files to support query implementation part 1 part 2 part 3 part 4 part 5 hocuspocus 3
-----
Score: 0.02942547213772617
URL: https://s2.smu.edu/~fmoore/textfiles/baseball5.txt
Title: No title for this article.
Description: mann ray people will come ray they ll come to iowa for reasons they can t even fathom they ll
Input query you wish: █
```

16. thisworks: There is not result for query “*thiswork*” therefore the expanded query ‘*this work*’ document is shown below:

```
Input query you wish: thiswork
-----
No article related this query. Now let us try query expansion:
Sorry, we cannot find any result about this query.
```

```
Input query you wish: this work
-----
Score: 0.14669291132566536
URL: https://s2.smu.edu/~fmoore/textfiles/extratextfiles/magictext.html
Title: thi is the magic file
Description: magic shows up here and in the title brown      beige    tan    auburn
-----
Score: 0.06536835341303449
URL: https://s2.smu.edu/~fmoore/misc/seo_email.txt
Title: No title for this article.
Description: original message from atul kumar mailto atul1930 outlook com sent thursday march 5 2020 12 58 am subject first page
-----
Score: 0.044771061983734387
URL: https://s2.smu.edu/~fmoore/textfiles/basketball3.txt
Title: No title for this article.
Description: 24 win season mavs coach careful not to promise too much for next season s basketball team our fans deserve
-----
Score: 0.037300292874449144
URL: https://s2.smu.edu/~fmoore/textfiles/football1.txt
Title: No title for this article.
Description: the first time since dez bryant was released last friday cowboy head coach jason garrett and team captions spoke publicly
-----
Score: 0.034156228153657014
URL: https://s2.smu.edu/~fmoore/
Title: freeman moor - smu spring 2020
Description: 2020 freeman l moore phd email fmoore lyle smu ed u fall 2019 cs 5330 7330 keep looking at the
```

Source Code:

IRProject12 for crawl:

```
#ShradhaUttarwar

import scrapy
from scrapy.http import Request
from scrapy.selector import Selector
import re
import json
import io
import hashlib
from nltk.stem import PorterStemmer

N = int(input("Enter value of N=40. 40 document to be indexed :"))
class IRProject12(scrapy.Spider):
    name = 'IRProject12'
    start_urls = ['https://s2.smu.edu/~fmoore/']
    dup_dic = {}
    count = 0

    def parse(self, response): #parse function is defined
        selector = Selector(response)
        md = hashlib.md5()
        md.update(response.body)
        hashed_data = md.hexdigest()
        outgoing_urls = []
        ingoing_urls = []
        full_url = []
        broken_urls = []
        nontxt_urls = []
        in_going_pdf_xlsx_urls = []
        duplicated_urls = []
        meta_noindex = []
        URL_Frontier = []
        graphic_urls = []

        if hashed_data in self.dup_dic:
            duplicated_urls.append(response.url)
        else:
            self.dup_dic[hashed_data] = 1

        if response.status == 404:
            broken_urls.append(response.url)
```

```

if response.status == 200:
    if self.count < int(N): #Limiting the number of pages to N
        self.count += 1
        word_dic = {}
        title = selector.xpath('//title') #extracting titles
        title = ''.join(title.xpath('string(.)').extract())
        title = title.replace('\n', '').replace('\r', '').strip()
        title = re.sub(r"\s{2,}", " ", title)
        urls = selector.xpath('//a/@href | //img/@src').extract()
        for url_part in urls: #filtering out URLs outside the domain of
            http and fmoore
                if re.match(r'^http', url_part):
                    url = url_part
                else:
                    if not re.search(r'~fmoore/$', response.url):
                        response_list = response.url.split('/')
                        del response_list[-1]
                        url = '/'.join(response_list) + '/' + url_part
                    else: url = response.url + url_part

                if re.match(r'^https://s2.smu.edu/~fmoore/', url) or
re.match(r'^http://s2.smu.edu/~fmoore/', url):
                    ingoing_urls.append(url)
                    # appending non-text URLs
                    if not (re.search(r'txt$', url_part) or re.search(r'htm$', url_part) or re.search(r'html$', url_part) or re.search(r'php$', url_part)):
                        nontxt_urls.append(url)
                        #filter out non-txt URLs (pdf, xlsx, jpg, jpeg, png, gif,
                        pptx, dontgohere)
                            if re.search(r'gif$', url_part) or re.search(r'jpg$', url_part) or re.search(r'jpeg$', url_part) or re.search(r'png$', url_part) or
re.search(r'pptx$', url_part):
                                graphic_urls.append(url)
                            elif re.search(r'pdf$', url) or re.search(r'xlsx$', url):
                                in_going_pdf_xlsx_urls.append(url)
                            elif 'dontgohere/' in url:
                                broken_urls.append(url)
                            elif '@lyle' in url:
                                broken_urls.append(url)
                            else: URL_Frontier.append(url)
                    else:
                        outgoing_urls.append(url)
                        full_url.append(url)

noindex = False
meta_tag = selector.xpath('//head/meta/@content').extract()
if meta_tag:

```

```

        if re.match(r'^noindex', ''.join(meta_tag).lower()):
            self.count -= 1
            meta_noindex.append(response.url)
            noindex = True
            exit

        if not noindex: #further processing
            with io.open('indexed_urls.txt', 'a', encoding='utf-8') as
file:
                url_title = response.url + ' --- ' + title
                file.write(url_title)
                file.write('\n')
                ps = PorterStemmer()
                texts = '
'.join(selector.xpath('//*/text()').extract()).replace('\n', ' ').replace('\r', ' ')
.strip()

                special_symbols = [',', '.', ';', '?', '(', ')', '[', ']',
'&', '!', '*', '@', '#', '$', '%', ':', "'", '=', '+', "\\", '/', '"', '-', '{', '}', '<', '>']

                for special_symbol in special_symbols:
                    texts = texts.replace(special_symbol, ' ')
                new_text = re.sub(r"\s{2,}", " ", texts)
                new_text = new_text.lower()
                text_list = new_text.split(' ')
                for word in text_list:
                    if word != '':
                        word_porter = ps.stem(word)
                        if word_porter not in word_dic:
                            word_dic[word_porter] = 1
                        else: word_dic[word_porter] += 1
                word_dic['*DocID*'] = self.count
                title = title.strip()
                title_list = title.split(' ')
                word_dic['*title*'] = []
                for title_word in title_list:
                    if title_word.strip() != '':
                        word_title_porter = ps.stem(title_word.strip())
                        word_dic['*title*'].append(word_title_porter)

                word_dic['*description*'] =
''.join(text_list[len(title_list):len(title_list)+20])
                word_dic['*url*'] = response.url

            with io.open('words_porter.json', 'a', encoding='utf-8') as
file:
                file.write(json.dumps(word_dic, ensure_ascii=False))
                line = u"\r\n"

```

```

        file.write(line)

    else:
        print("Too high value of N")
        exit

    if URL_Frontier:
        for in_url in URL_Frontier:
            yield Request(in_url, callback=self.parse, dont_filter=True)

    if outgoing_urls:
        for outgoing_url in outgoing_urls:
            with io.open('outgoing_urls.txt', 'a', encoding='utf-8') as file:
                file.write(outgoing_url)
                file.write('\n')

    if broken_urls:
        for broken_url in broken_urls:
            with io.open('broken_urls.txt', 'a', encoding='utf-8') as file1:
                file1.write(broken_url)
                file1.write('\n')

    if nontxt_urls:
        for nontxt_url in nontxt_urls:
            with io.open('nontext_urls.txt', 'a', encoding='utf-8') as file2:
                file2.write(nontxt_url)
                file2.write('\n')

    if graphic_urls:
        for graphic_url in graphic_urls:
            with io.open('graphic_urls.txt', 'a', encoding='utf-8') as file2:
                file2.write(graphic_url)
                file2.write('\n')

    if in_going_pdf_xlsx_urls:
        for in_going_pdf_xlsx_url in in_going_pdf_xlsx_urls:
            with io.open('in_going_pdf_xlsx_urls.txt', 'a', encoding='utf-8') as
file3:
                file3.write(in_going_pdf_xlsx_url)
                file3.write('\n')

    if duplicated_urls:
        for duplicated_url in duplicated_urls:
            with io.open('duplicated_urls.txt', 'a', encoding='utf-8') as file4:
                file4.write(duplicated_url)
                file4.write('\n')

    if meta_noindex:

```

```

for meta_noindex_url in meta_noindex:
    with io.open('meta_noindex.txt', 'a', encoding='utf-8') as file5:
        file5.write(meta_noindex_url)
        file5.write('\n')

```

Term frequency matrix:

```

#ShradhaUttarwar

import io
import json
word_dic_full = {}
for line in io.open("words_porter.json", 'r', encoding='utf-8'):
    new_line = line.replace('\n', '').replace('\r', '').strip()
    if new_line != '':
        word_dic = json.loads(new_line)
        for key in word_dic.keys():
            if key != '*DocID*' and key != '*description*' and key != '*url*' and key != '*title*':
                if key not in word_dic_full:
                    word_dic_full[key] = [0 for _ in range(40)]
                    word_dic_full[key][word_dic['*DocID*']-1] += word_dic[key]
                else:
                    word_dic_full[key][word_dic['*DocID*']-1] += word_dic[key]
with io.open('Words_Frequency_Matrix_porter_query.json', 'a', encoding='utf-8') as file:
    for item in word_dic_full.keys():
        word_dic_temp = {}
        word_dic_temp[item] = word_dic_full[item]
        file.write(json.dumps(word_dic_temp, ensure_ascii=False))
        line = u"\r\n"
        file.write(line)

```

Query Engine:

```

# -*- coding: utf-8 -*-
#####
Created on Sat May  2 16:49:18 2020

@author: Shradha
#####

import json
import io
import math
import numpy as np

```

```

from nltk.stem import PorterStemmer

number_punctuations = [',', '.', ':', ';', '?', '(', ')', '[', ']', '&', '!', '*', '@', '#', '$', '%', ': ', '0', '1', '2', '3', '4', '5', '6', '7', '8', '9']

words_documents = {}
document_number = 0

for line in io.open("words_porter.json", 'r', encoding='utf-8'):
    new_line = line.replace('\n', '').replace('\r', '').strip()
    if new_line != '':
        document_number += 1
        word_dic = json.loads(new_line)
        words_documents[document_number] = word_dic

words_matrix = {}
for line in io.open("Words_Frequency_Matrix_porter_query.json", 'r', encoding='utf-8'):
    new_line = line.replace('\n', '').replace('\r', '').strip()
    if new_line != '':
        word_dic_temp = json.loads(new_line)
        for key in word_dic_temp:
            words_matrix[key] = word_dic_temp[key]

documents_list = {}
for word in words_matrix.keys():
    for idx, frequency in enumerate(words_matrix[word]):
        if idx+1 not in documents_list:
            documents_list[idx+1] = []
        documents_list[idx+1].append(frequency)

idf_dict = {}
word_number = 0
for word in words_matrix.keys():
    word_number += 1
    idf_dict[word_number] = document_number - words_matrix[word].count(0) + 1

words_expansion = {'beautiful': ['nice', 'fancy'], 'chapter': ['chpt'],
'chpt': ['chapter'], 'responsible': ['owner', 'accountable'], 'freemanmoore':
['freeman', 'moore'], 'dept': ['department'], 'photo': ['photograph', 'image'],
'picture': ['brown': ['beige', 'tan', 'auburn'], 'tues': ['Tuesday'], 'sole':
['owner', 'single', 'shoe', 'boot'], 'homework': ['hmwk', 'home', 'work'],
'novel': ['book', 'unique'], 'computer': ['cse'], 'story': ['novel', 'book'],
'hocuspocus': ['magic', 'abracadabra'], 'thisworks': ['this', 'work'], } }

is_expanded = False
last_query_list = []
while True:
    if not is_expanded:

```

```

query = input('\nInput query you wish: ')

if query == 'stop':
    print('\nClosing the query engine...')
    break
query_list = query.split(' ')
else:
    query_list = [query_word for query_word in last_query_list]
    for query_word in last_query_list:
        if words_expansion.get(query_word):
            for expanded_word in words_expansion[query_word]:
                query_list.append(expanded_word)

ps = PorterStemmer()
query_list_porter = []
for original_word in query_list:
    word_porter = ps.stem(original_word)
    query_list_porter.append(word_porter)

for matrix_word in words_matrix.keys():
    words_matrix[matrix_word].append(0)

word_flag = False
for porter_word in query_list_porter:
    for matrix_word in words_matrix.keys():
        if matrix_word == porter_word:
            word_flag = True
            words_matrix[matrix_word][-1] += 1

if word_flag:
    documents_list[0] = []
    for matrix_word in words_matrix.keys():
        documents_list[0].append(words_matrix[matrix_word][-1])

similarity_scores_dict = {}
for i in range(1, document_number+1):
    query_score = []
    query_sum = 0
    for idx, tf in enumerate(documents_list[0]):
        if tf > 0:
            tf_temp = (1 + math.log(tf)) *
math.log(document_number/(idf_dict[idx+1]+1))
        else:
            tf_temp = 0
            query_sum += tf_temp * tf_temp
            query_score += [tf_temp]
    query_normalized = math.sqrt(query_sum)

```

```

query_score_processed = np.array([score/(query_normalized+1) for score in query_score])

document_score = []
document_sum = 0
for tf in documents_list[i]:
    if tf > 0:
        tf_temp = (1 + math.log(tf))
    else:
        tf_temp = 0
    document_sum += tf_temp * tf_temp
    document_score += [tf_temp]
document_normalized = (math.sqrt(document_sum) )
document_score_processed = np.array([score/(document_normalized+1) for score in document_score])

similarity_score = np.dot(query_score_processed, document_score_processed)

query_in_title = False
for porter_word in query_list_porter:
    if porter_word in words_documents[i]['*title*']:
        query_in_title = True
if query_in_title:
    similarity_score += 0.1
similarity_scores_dict[i] = similarity_score
similarity_scores_dict_new = sorted(similarity_scores_dict.items(), key=lambda d: d[1], reverse=True)
result_count = 0
for number, score in similarity_scores_dict_new[:5]:
    if score > 0:
        result_count += 1
        print('-----')
        print('Score: ', score)
        print('URL: ', words_documents[number]['*url*'])
        if not words_documents[number]['*title*']:
            title = 'No title for this article.'
        else:
            title = ' '.join(words_documents[number]['*title*'])
        print('Title: ', title)
        print('Description: ', words_documents[number]['*description*'])
if not is_expanded:
    print('-----')
    print('No article related this query. Now you can reenter Query:')
    is_expanded = True
    last_query_list = [query_word for query_word in query_list]
else:
    if is_expanded:
        is_expanded = False

```

```
last_query = []

else:
    if is_expanded:
        is_expanded = False
        print('Sorry, we cannot find any result about this query.')
    else:
        print('-----')
        print('No article related this query. Now let us try query expansion:')
        is_expanded = True
        last_query_list = [query_word for query_word in query_list]
```

-----End-----