

Information Retrieval Homework 1

Pradeep Agrawal (IQ57213)

Language used: Python

Following is the step-by-step explanation of how I did this:

1. I downloaded the tarfile and then unzipped it. By doing this, I got 503 HTML files. I put all the files in a directory named 'files'.
2. Then I wrote a python code in which I used `os.scandir` to iterate through all the HTML files. In each iteration, I was checking if the current iteration holds a file.
3. For each valid file, I was calling a method 'extract_tokens' by passing the filepath. This method was responsible for extracting the words out of the HTML file.
4. For this word extraction, I am using BeautifulSoup (a python library). This library is providing me all the text present in the file. Once I got the text, I broke that into words by using `split()` method. After getting the words, I am using a Regular Expression to check if the words have any special characters or numbers. If they had such things, then I removed those characters. And then I took only those words which were greater than 2 in length and lowercase all of them.
5. Once I got those tokens, I passed those tokens to a method called 'create_token_files'. This method also requires output directory path. So, this method will check if the output directory exists. If not, will create that directory and then write these tokens into a text file separated by newline. So, this will result in having 503 txt files of tokens.

```
def extract_tokens(file_path):
    with open(file_path, 'r', encoding='utf-8', errors='ignore') as fp:
        soup = BeautifulSoup(fp, features="html.parser")
        for script in soup(["script", "style"]):
            script.extract()
        text = soup.get_text()
        words = text.split()
        words = [re.sub(r'[0-9]+|\W+|_', '', word) for word in words]
        words = [word.lower() for word in words if len(word) > 2]
        return words

def create_token_files(output_dir, filename, words):
    if not os.path.exists(output_dir):
        os.makedirs(output_dir)
        print("The new directory is created!")
    filepath = output_dir + '/' + filename.split('.')[0] + '.txt'
    with open(filepath, 'w') as f:
        for word in words:
            f.write(word)
            f.write('\n')
```

6. I wrote a method 'create_frequency_files' to create frequency files. So, when I got the token files, I called this method by passing `output_dir` as parameter. This method iterates over these text files and read words in them. Here I am using python dictionary to hold frequency of each word. So, if the word does not exist in the dictionary, then I add that word as key with 1 as its value. And if it exists then I increase its count by 1. At the end of this loop, I get frequencies of every word respectively in this dictionary. Using this, I created 2 files. One having words sorted by tokens alphabetically and second sorted by their respective frequency.

Code Execution Guide:

We just need to run the python file with 2 command line arguments (input files path, output files path). Keep in mind that we have to have the input-directory in the same path where we have the `tokenize_search.py` file.

`python tokenize_search.py <input-directory> <output-directory>`

Ex. `python tokenize_search.py html-files-dir text-files-dir`

In above command 'tokenize_search.py' is the python file name. 'html-files-dir' is the directory which holds all the HTML files and then 'text-files-dir' is the directory name where we will put all the text token files.

Following are screenshots of the output:

```
((search_engine@3.9.10) → search_engine python3 tokenize_search.py files text_files_dir
All the times are in milliseconds
Parsed Files 100 and Time taken 1123.539
Parsed Files 200 and Time taken 881.482
Parsed Files 300 and Time taken 819.741
Parsed Files 400 and Time taken 685.798
Parsed Files 500 and Time taken 759.315
Total elapsed time: 4487.913
(search_engine@3.9.10) → search_engine █
```

Above shown the time taken in parsing those HTML files and converting them into token files. And finally, creating 2 frequency files.

Code Execution Flow:

When we run the command, it will start execution and will call the 'main' method that will run the entire script in the required sequence as shown in the screenshot.

```
def main(input_dir, output_dir):
    file_count = 0
    start = time.process_time_ns()
    for filename in os.listdir(input_dir):
        if filename.is_file():
            file_count += 1
            tokens = extract_tokens(filename.path)
            create_token_files(output_dir, filename.name, tokens)
    if file_count%100 == 0:
        end = time.process_time_ns()
        print(f'Parsed Files {file_count} and Time taken {(end - start)/1000000}')
        start = time.process_time_ns()
    create_frequency_files(output_dir)
```

Incorrectly Tokenized Words:

After this entire process, I investigated both the generated frequency files and noticed that there are some random words that does not make any sense are there. Those words are some random collection of characters.

Ex. abaloldalizzasnak, abortusztanacsadas, acelparreorganizacios, acsaladtamogatasok etc

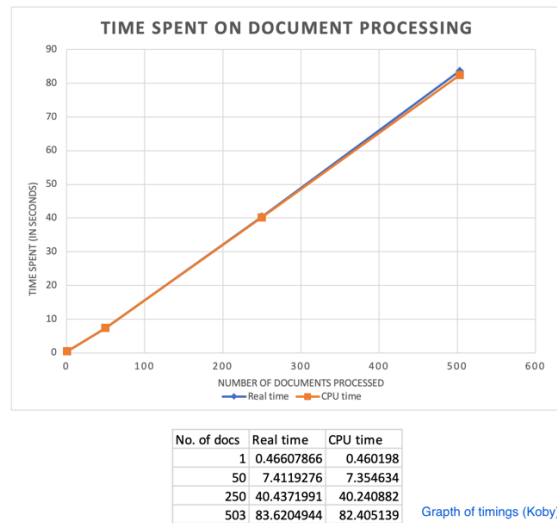
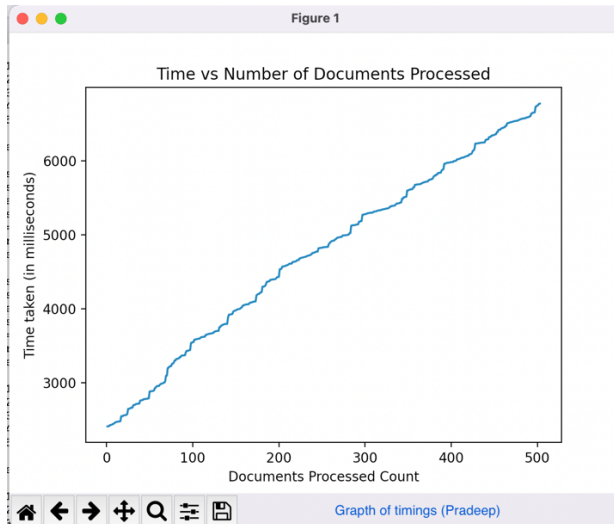
It seems like there were some special characters in between these words inside somewhere which got removed while cleaning the data and I end up getting such words.

Results Comparison:

- I compared my results with Vinayak's and Koby's works. And I noticed that Vinayak is using Jsoup parser to parse the HTML files while I used BeautifulSoup.
- Vinayak is removing all the special characters and converting it into lowercase and then splitting them by whitespace. While I am removing special characters, numbers and symbols. In addition to that I am keeping only the tokens which have character length of greater than 2.

- Koby's frequency file have tokens like 'a', 's', '79yr', '€' etc. While my frequency file contains only words which does not have any integers in them.
- Koby's code is taking approx. 83 seconds in the execution of code. While my code takes ~6 seconds in worst case.

Graphs:



Token files sorted by alphabets:

Word	Frequency	
1	aaa	21
2	baas	2
3	aaathomas	1
4	aacutevosok	1
5	aeliberalis	1
6	aemassago	1
7	aemegvalosult	1
8	aardvark	1
9	aaron	8
10	aarons	1
11	aarp	2
12	aas	4
13	aasad	1
14	ababa	2
15	abacha	1
16	abacsiskun	1
17	abalczo	1
18	abalkanon	1
19	abaloldalzasnak	1
20	aban	1
21	abandon	5
22	abandoned	1
23	abank	1
24	abanknak	1
25	abankoknak	2
26	abankrendszer	1
27	abankuralom	1
28	abankuralom	1
29	abankuralom	1
30	abankuralom	1
31	abankuralom	1
32	abankuralom	1
33	abankuralom	1
34	abankuralom	1
35	abankuralom	1
36	abankuralom	1
37	abankuralom	1
38	abankuralom	1
39	abankuralom	1
40	abankuralom	1
41	abankuralom	1
42	abankuralom	1
43	abankuralom	1
44	abankuralom	1
45	abankuralom	1
46	abankuralom	1
47	abankuralom	1
48	abankuralom	1
49	abankuralom	1
50	abankuralom	1
51	abankuralom	1
52	abankuralom	1
53	abankuralom	1
54	abankuralom	1
55	abankuralom	1
56	abankuralom	1
57	abankuralom	1
58	abankuralom	1
59	abankuralom	1
60	abankuralom	1
61	abankuralom	1
62	abankuralom	1
63	abankuralom	1
64	abankuralom	1
65	abankuralom	1
66	abankuralom	1
67	abankuralom	1
68	abankuralom	1
69	abankuralom	1
70	abankuralom	1
71	abankuralom	1
72	abankuralom	1
73	abankuralom	1
74	abankuralom	1
75	abankuralom	1
76	abankuralom	1
77	abankuralom	1
78	abankuralom	1
79	abankuralom	1
80	abankuralom	1
81	abankuralom	1
82	abankuralom	1
83	abankuralom	1
84	abankuralom	1
85	abankuralom	1
86	abankuralom	1
87	abankuralom	1
88	abankuralom	1
89	abankuralom	1
90	abankuralom	1
91	abankuralom	1
92	abankuralom	1
93	abankuralom	1
94	abankuralom	1
95	abankuralom	1
96	abankuralom	1
97	abankuralom	1
98	abankuralom	1
99	abankuralom	1
100	abankuralom	1

Token files sorted by frequencies:

frequencies_sorted_by_frequency.txt			Sorted_by_Frequency.txt		freqfreq.txt	
1	the	33420	1	Word	Frequency	0.00 : 62496
2	and	17210	2	=====		the : 33107
3	pgs	6327	3	the	32978	a : 27611
4	for	4678	4	a	26471	of : 21042
5	that	3529	5	of	20896	and : 17135
6	hogy	3264	6	and	17029	to : 11181
7	with	2920	7	to	11226	in : 10044
8	nem	2534	8	in	9998	az : 6918
9	are	1933	9	s	7393	is : 5605
10	from	1820	10	az	6347	for : 4621
11	this	1798	11	pgs	6327	on : 4123
12	which	1235	12	yr	6327	's : 3871
13	not	1218	13	fn	5795	es : 3735
14	meg	1192	14	is	5469	that : 3515
15	his	1190	15	for	4585	as : 3431
16	was	1176	16	on	4025	de : 3399
17	magyar	1143	17	as	3461	hogy : 3388
18	their	1101	18	that	3438	by : 3263
19	have	1081	19	de	3316	with : 2892
20	examines	1077	20	es	3264	an : 2861
21	will	1077	21	by	3251	nem : 2575
22	analysis	1065	22	hogy	3092	or : 2528
23	has	1047	23	an	2897	be : 2109
24	its	1036	24	with	2871	are : 1915
25	new	965	25	or	2527	from : 1796
26	egy	958	26	nem	2394	this : 1774
27	all	864	27	be	2117	0s : 1766
28	law	820	28	are	1905	at : 1715
29	other	807	29	this	1781	0fn : 1435
30	can	751	30	from	1729	it : 1408
31	may	750	31	at	1699	not : 1308

Pradeep

Vinayak

Koby