

Information Retrieval and Web Search

Project Report

1. System Design

- a. Background:* This project is written entirely in python and python libraries are used. For ease of execution, the entire code is also provided within the python notebook to execute the entire code at once.
- b. Input Provided:* A folder called “ft911” that has all the files containing the documents and the information is provided along with a stopwords list file that needs to be used to while performing tokenization on the word content.
- c. File Reading and Text Processing:*
 - i) The first step of the process is to extract the relevant information from the given files. To do this, TextParser’s fetchDocs() method is called upon which takes in the given ft911 folder and reads all the documents and loads the name of the documents and the information to memory.
 - ii) The next step consists of creating dictionaries, where the file names are stored in a file dictionary along with a unique identifier. For the word dictionary, words are first tokenized (by taking only non-alphanumeric words and removing stop words) and then stemmed using the NLTK library’s PorterStemmer and the unique words are stored in Word dictionary along with a unique identifier.
- d. Indexing:* There are two indexes that are generated in this step. Forst a forward_index, which has FileIds with all the stemmed words that it contains along with their frequency and an inverted_index that has all the words and the ft911_file in which it occurs along with the frequency in that file.
- e. Output Generation:*
 - i) parser_output.txt has the word dictionary and file dictionary that are ordered by their IDs.

- ii) Two index files, forward_index.txt and inverted_index.txt are generated by the code to store the forward and inverted indexes.

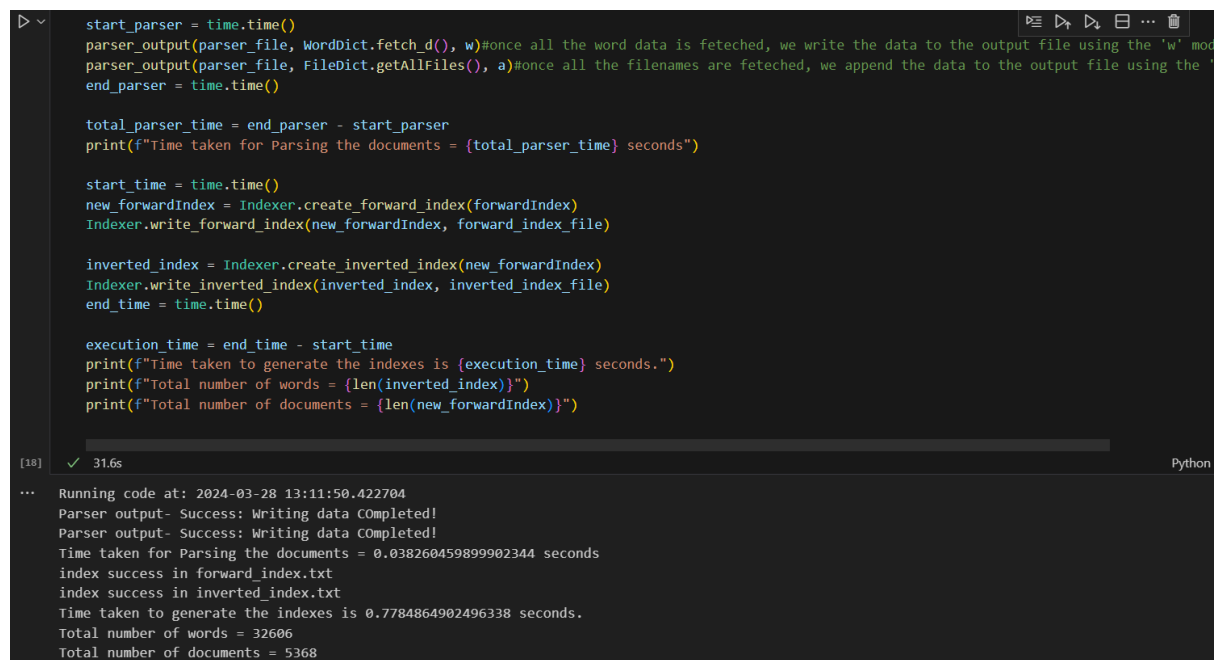
2. Execution Time Analysis and Index Size:

Time taken to Parse and create IDs for all documents is <1second.

Time taken to generate Forward and Inverted Indexes is <1second.

Total Number of Documents = 5368

Total Number of Words = 32606



```
start_parser = time.time()
parser_output(parser_file, WordDict.fetch_d(), w)#once all the word data is feteched, we write the data to the output file using the 'w' mod
parser_output(parser_file, FileDict.getAllFiles(), a)#once all the filenames are feteched, we append the data to the output file using the '
end_parser = time.time()

total_parser_time = end_parser - start_parser
print(f"Time taken for Parsing the documents = {total_parser_time} seconds")

start_time = time.time()
new_forwardIndex = Indexer.create_forward_index(forwardIndex)
Indexer.write_forward_index(new_forwardIndex, forward_index_file)

inverted_index = Indexer.create_inverted_index(new_forwardIndex)
Indexer.write_inverted_index(inverted_index, inverted_index_file)
end_time = time.time()

execution_time = end_time - start_time
print(f"Time taken to generate the indexes is {execution_time} seconds.")
print(f"Total number of words = {len(inverted_index)}")
print(f"Total number of documents = {len(new_forwardIndex)}")
```

[18] ✓ 31.6s Python

... Running code at: 2024-03-28 13:11:50.422704
Parser output- Success: Writing data Completed!
Parser output- Success: Writing data Completed!
Time taken for Parsing the documents = 0.038260459899902344 seconds
index success in forward_index.txt
index success in inverted_index.txt
Time taken to generate the indexes is 0.7784864902496338 seconds.
Total number of words = 32606
Total number of documents = 5368

3. Search Engine Performance:

To measure the performance of the search engine, first the query from a given input document needs to be parsed. This is done by calling the query_parser() method in the main.py file. This method will separate the titles, descriptions and narratives from the given input document and add it to a dictionary which will be used to calculate the cosine similarity.

To calculate the cosine score, we pass the content we obtained above in the following settings to the getCosineScore() method in main.py:

- i. title
- ii. title + description
- iii. title + narrative

The following formula is used to calculate the cosine score for the given query:

Cosine Similarity: When comparing two vectors—in this case, the query and a document—cosine similarity calculates the cosine of their angle.

Its formula is,

$$\frac{V(d1).V(d2)}{|V(d1)| |V(d2)|}$$

Here, d1 and d2 are two documents and V(d1) and V(d2) are its vector representatives. In the `getCosineScore()` method, we initialize two dictionaries, one to store the scores the other to store the length of the documents. We iterate through each term in the given query and get its forward index which is used to calculate the idf of that term. Using a set of documents, this algorithm determines the cosine score for a particular query. A combination of dictionaries and lists are used in this search engine engineering. It computes IDF, extracts pertinent documents from the inverted index, and repeats this process for each query phrase. It then normalises the scores after adding them up for every document based on TF-IDF values. In the end, the normalised scores are given back as the cosine similarity between each document and the query.

The outputs for the vector space models are provided in the three text files `vsm_title.txt`, `vsm_description.txt` and `vsm_narrative.txt`.

The relevant documents retrieved for each of the three settings are in the format

Setting1 = { topic number1 : number of relevant documents.... }

topic = {352: 1112, 353: 124, 354: 445, 359: 872}

topic + description = {352: 1885, 353: 1876, 354: 2375, 359: 1485}

topic + narrative = {352: 3787, 353: 3394, 354: 2123, 359: 2722}