

Information Retrieval Phase 3 Report

PROGRAM STRUCTURE

1. Open and read through all the html files
2. Using BeautifulSoup library to parse through all the .html files and extract only the text from it
3. Tokenise the extracted words
4. Handles cases such as special characters, numerical inside a word, and also decimal numbers
5. Converts all the words to lower case
6. Reads the text file stopwords.txt and stores the given specific stop words in a list
7. The extracted words are now passed into a counter to get the word and its frequency value
8. Checks the given conditions that is
 - a. Words of length 1
 - b. Words with frequency 1
 - c. Words that match the stop words
9. Updates the counter filtering out the words from those three conditions
10. Computes the term frequency and stores the value in a global dictionary with file name as key and another dictionary as value, containing word and tf
11. Computes the value that gives the occurrence count of any given word in number of files
12. Computes the inverse document frequency and stores the value in a global dictionary with file name as key and another dictionary containing word and idf as value
13. With these values the tf-idf is calculated for each word in the global dictionary
14. With the tf-idf values we create a Document Term Matrix
15. With the help of TDM and tf-idf dictionaries we create the posting file.

EXPLANATION

The program is built on top of the previous submitted logic with few minor changes and with better efficiency.

The developed program successfully parses through all of the html files with the help of beautiful soup library and tokenises all of the words. The main basis of tokenisation is spacing. The inbuilt function '.split' is used to split the words extracted.

The program handles punctuations and special characters using regular expression. We give a set of predefined symbols that is filtered out of the extracted and tokenised words. Then all the tokenised words are converted to lower case.

The program then reads the given set of stop words from a text files and filters out the words based on the three given conditions. Then the counter is updated accordingly with the new values. The globally declared dictionary has the file name as the key and a counter of word and its frequency as its value.

The **compute_tf** function in the program then takes this dictionary as the parameter and computes the term frequency. It first computes the word count of every file and stores

the file name and its word count in a dictionary. Then with the word count of that word and total count of words in that file term frequency is calculated. This is done for each and every word in the dictionary and the dictionary is accordingly updated with the new set of values.

Then the **doc_containing** function takes a word as its parameter and calculates the count of documents that word exists in and returns the count.

The **compute_idf** takes the dictionary as an argument and computes the inverse document frequency of the words. It first checks the number of documents to be checked and stores that length in the variable num_doc. For every word in the dictionary it calculates idf with the formula, log of number of documents divided by the word count(calculated using the doc_containing function). With these values it updates another dictionary with the key as file name and another dictionary as its value with the word and its idf.

The function **compute_tf_idf** takes the dictionary as its argument and then multiplies the values of tf and idf to compute the tf-idf of every word. It then writes the obtained value to a text file.

The function **calculate_tdm** forms a document term matrix from the tdf-dictionary. And populates an another dictionary with words as its key and a dictionary with files names as its keys and tf-idf values of that word as its value.

The function **dictionary_file** creates a text file with word, number of files that has that word and the first occurrence.

The function **posting_file** creates a file text file referencing from this tdm dictionary and dictionary file to map a words first occurrence and its tf-idf value.

The program takes about 22 seconds to complete its execution

RESULTS

The result of the program is that it creates a directoryfile.txt file with word, number of files that has that word and the first occurrence and postingfile.txt file with first occurrence of a word and it's tf-idf value.

Document.txt:

```
nemzet
42
1
eeeeeee
81
43
zzzzzzz
120
124
```

Posting.txt:

```
340,0.004510300729911663
356,0.014107503135320288
376,0.006072566314439305
360,0.002437422007673139
337,0.006339371621658819
336,0.0070778807064320725
361,0.005137962859423426
```

FORMULAS IMPLEMENTED

1. `compute_tf()` : number of times a word appears in a document dividing by the total number of words.
2. `computer_idf()` : log of total number of documents to the number of documents containing word.
3. `compute_tf_idf()` : product of tf and idf.

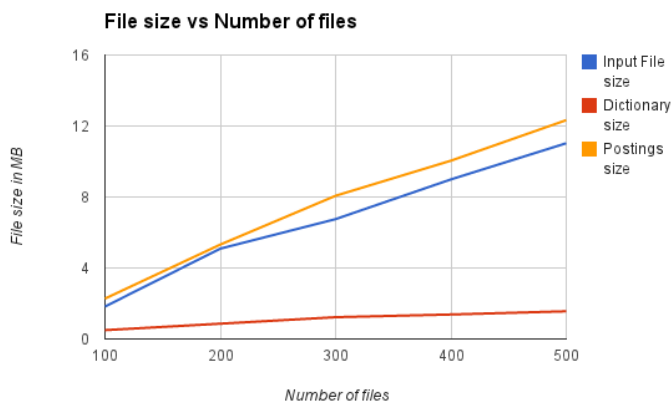
EXECUTION

The program was coded in PyCharm. It can be run on any IDE that supports Python.

Imports,

1. Import os
2. Import glob
3. Import math
4. from collections import Counter
5. from bs4 import BeautifulSoup
6. Import time

GRAPHS



⌚ Time elapsed for 503 files in seconds

