

INFORMATION RETRIEVAL - CSE 508

ASSIGNMENT 1

SUBMITTED BY:	Abhinav garg-	MT22002
	Riya garg -	MT22058
	Saksham Nautiyal-	MT22061

ASSUMPTIONS:

1. All files have text based content that we need to preprocess.

Q1. DATA PREPROCESSING:

LIBRARIES USED:

- Nltk
- Beautiful Soup
- punkt

(i)First we print 5 files before extraction of data between tags.

To extract contents between contents between the <TITLE>...</TITLE> and <TEXT>...</TEXT> tags we have used the **Beautiful Soup** library which provides simple methods for navigating, searching, and modifying a parse tree in HTML, XML files.

Then we concatenate the 2 strings using blank space and rewrite them into the files.

```

from bs4 import BeautifulSoup
for f1 in list1:
    f=open(f'{path}/{f1}',mode='r')
    data = f.read()
    f.close()
    title1=""
    text=""
    bsText = BeautifulSoup(data, 'html.parser')
    #title1 = BeautifulSoupText.findAll('TITLE').text
    #text = BeautifulSoupText.findAll('TEXT').text
    # print(BeautifulSoupText)
    for tag in bsText.findAll('title'):
        title1 = bsText.find(tag.name).text
        #print(title1)
    for tag in bsText.findAll('text'):
        text = bsText.find(tag.name).text
    f=open(f'{path}/{f1}',mode='w')
    x = title1 + " " + text
    f.write(x);
    f.close();

```

(ii)Preprocessing:

1. **Lowercase the text:** For this we have used the String lower() method. It converts all uppercase characters in a string into lowercase characters and returns it.
2. **Perform tokenization:** NLTK contains a module called tokenize with a word_tokenize() method that will help us split a text into tokens.
3. **Remove stopwords :** for this stopwords library from nltk.corpus
4. **Remove punctuations:** For this we use regular expressions and string function.
5. **Remove blank space tokens** using string function.

Output after preprocessing:

```

File 1

experimental investigation aerodynamics wing slipstream experimental study wing propeller slipstream made order determine spanwise distribution lift
increase due slipstream different angles attack wing different free stream slipstream velocity ratios results intended part evaluation basis different
theoretical treatments problem comparative span loading curves together supporting evidence showed substantial part lift increment produced slipstream
due destalling boundarylayercontrol effect integrated remaining lift increment subtracting destalling lift found agree well potential flow theory
empirical evaluation destalling effects made specific configuration experiment

File 2

simple shear flow past flat plate incompressible fluid small viscosity study highspeed viscous flow past twodimensional body usually necessary
consider curved shock wave emitting nose leading edge body consequently exists inviscid rotational flow region shock wave boundary layer situation
arises instance study hypersonic viscous flow past flat plate situation somewhat different prandtl s classical boundarylayer problem prandtls original
problem inviscid free stream outside boundary layer irrotational hypersonic boundarylayer problem inviscid free stream must considered rotational
possible effects vorticity recently discussed ferri libby present paper simple shear flow past flat plate fluid small viscosity investigated shown
problem treated boundarylayer approximation novel feature free stream constant vorticity discussion restricted twodimensional incompressible steady
flow

```

Ques 2: Boolean Queries:

Assumption: Stop words are irrelevant for creating a positional or inverted index.

UNIGRAM INVERTED INDEX:

To store the document ids for a given term, we created a dictionary we call postings = {}. On each distinct word, we have created a for loop. If a word is already in our dictionary postings, we append the document id to the list if it is. We shall create a new value pair if the term has just appeared.

Key in postings is current word and value is list which contains the document ids in which the word is present.

After creating all the posting lists for all the terms, we store them using the "unigram postings" Python pickle.

We have created four functions, OR_merge(), AND_merge(), AND_NOT_merge(), and OR_NOT(), which when given two parameters—the document ids lists of two terms—perform the necessary logic.

After accepting the input in the required format, we first load the pickle we earlier saved in order to retrieve/load the posting lists of the terms.

Next we call the function process_query(), which checks the operators list and calls the appropriate function based on its findings.

For the first two terms of the query, we first execute the process_query() method, after which we save the result and the number of comparisons.

Then, after doing the process with the term posts and results for the remaining terms, we have the final result merged list and the minimal number of comparisons.

In order to retrieve the names of the documents, we also maintain a doc map dictionary using Python pickle, which maps document names to document ids.

No. of comparisons required for query 1 : 1

QUES 3:

BIGRAM INVERTED INDEX:

To find the bigrams first we create a tuple containing two adjacent words. Then we are creating a dictionary as in ques 2 `bi_postings={}`. In this dictionary tuple of adjacent words is a key and value is a list which contains the document ids in which the tuple is present.

On each distinct tuple, we have created a for loop. If both words in tuple are already in our dictionary postings, we append the document id to the list if it is. We shall create a new value pair if the tuple has just appeared.

After creating all the posting lists for all the terms, we store them using the "bigram_postings" Python pickle.

Positional index:

To store the document ids and locations of a given word, we created a dictionary named `posInd = {}`.

The current word serves as the dictionary's key, and the value is a list, where the first element of the list denotes frequency and the second element contains the dictionary.

The dictionary's keys are the document ids in which the word appears, and its value is a list of the locations in the specified document where the current word can be found.

After creating all the posting lists for all the terms, we store them using the "positional_postings" Python pickle.

To do a search in the bigram inverted index, we turn the query into a collection of tuples, or sets of two words, and then we process it by invoking the `bigram_query()` method, which performs an AND merge boolean query on every tuple in the collection before printing the results.

The process `positional_query()` method is called by the query and gets the tokens processed. It then checks the positional postings dictionary to see if the tokens are in adjacent positions and stores them in the list if they are.

For each query the output will be in the following format:

- i. Number of documents retrieved using bigram inverted index
- ii. Names of documents retrieved using bigram inverted index

- iii. Number of documents retrieved using positional index
- iv. Names of documents retrieved using positional index

```
Query 1 : flow past flat
Number of documents retrieved for query 1 using bigram inverted index: 7
Names of documents retrieved for query 1 using bigram inverted index:
cranfield0002,cranfield0003,cranfield0388,cranfield0389,cranfield1186,cranfield0308,cranfield0663,
Number of documents retrieved for query 1 using positional index: 7
Number of documents retrieved for query 1 using positional index:
cranfield0002,cranfield0003,cranfield0308,cranfield0388,cranfield0389,cranfield0663,cranfield1186,
```

Comparison

When a query is processed using the positional index we check for adjacent positions whereas using the bigram inverted index we check for combinations of words together which might lead to different documents where the words occur together but the total query does not occur together.

```
Query 1 : design highspeed flight
Number of documents retrieved for query 1 using bigram inverted index: 1
Names of documents retrieved for query 1 using bigram inverted index: cranfield0012,
Number of documents retrieved for query 1 using positional index: 0
Number of documents retrieved for query 1 using positional index:
```

Here in document number 12 the words design highspeed flight are not together but (design,highspeed) and (highspeed,flight) are present in different locations so using bigram inverted index it is printing that the query is present and using positional index it is shows not present.