Muhammed Göktepe
2017400162

Assignment 2

Describe the steps you have performed for data preprocessing.

In this project we are assigned to implement a document retrieval system for simple boolean queries using the positional inverted indexing scheme. Reuters news is our dataset, and we have some stopwords to discard. First, I stored stopwords in a list for not to include these words while indexing texts. I have a dictionary of dictionary named index_dict. I store newId's of each word and their corresponding positions in documents in here. Then I started to traverse .sgm files in reuters21578 directory. I used BeautifulSoup library to ease .sgm file processing. Then, I store all articles in a file in news_list list. I iterate over each news and look for their titles and body's. If the article contains either of these, I combine the title and body and send it to the normalize function. In normalize function I made punctuation removal, case folding and finally stopwords removal. After these steps I added these words to dictionary and create a dictionary for each word. For these words I added their news ids and their positions in the document. And printed this dictionary to file.

Describe the data structures (hash, b-tree, linked list etc.) that you used for representing the inverted index (i.e., the dictionary and the postings lists).

I used dictionary of dictionary to store words, news ids and word positions in documents. Since .sgm files read in ascending order of news ids it stores new ids in ascending order. Also, positions are in ascending order in document. I used AND operator like algorithm while implementing phrase queries. In text free queries I used Lec03 as my reference.

Provide a screenshot of running the indexing module of your system.

Since I am using BeautifulSoup library it must be installed on working environment.

```
PS C:\Users\muhammed\Desktop\Ders\7.Dönem\Cmpe493\Assignment1> pip install bs4
WARNING: Ignoring invalid distribution -ip (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution - (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution -ip (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution - (c:\python310\lib\site-packages)
Requirement already satisfied: bs4 in c:\python310\lib\site-packages (0.0.1)
Requirement already satisfied: beautifulsoup4 in c:\python310\lib\site-packages (from bs4) (4.10.0)
Requirement already satisfied: soupsieve>1.2 in c:\python310\lib\site-packages (from beautifulsoup4->bs4) (2.3.1)
WARNING: Ignoring invalid distribution -ip (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution - (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution -ip (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution - (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution -ip (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution - (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution -ip (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution - (c:\python310\lib\site-packages)
WARNING: You are using pip version 21.2.3; however, version 21.3.1 is available.
You should consider upgrading via the 'C:\Python310\python.exe -m pip install --upgrade pip' command.
PS C:\Users\muhammed\Desktop\Ders\7.Dönem\Cmpe493\Assignment1> python .\preprocessing.py
```

Muhammed Göktepe
2017400162

You can install it by typing "pip install bs4". Then you can run my preprocessing code by typing "python preprocessing.py".

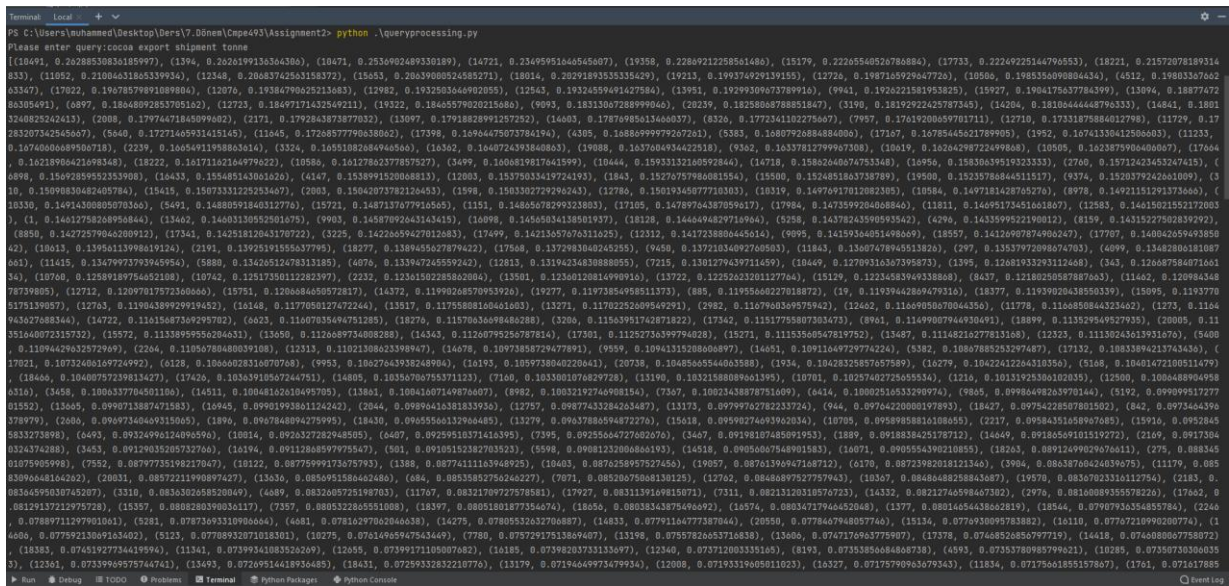Provide four screenshots of running your system for each of the four types of queries.

"old crop cocoa"

```
PS C:\Users\muhammed\Desktop\Ders\7.Dönem\Cmpe493\Assignment2> python .\queryprocessing.py
Please enter query:"old crop cocoa"
['1']
PS C:\Users\muhammed\Desktop\Ders\7.Dönem\Cmpe493\Assignment2>
```

"sugar price"

```
PS C:\Users\muhammed\Desktop\Ders\7.Dönem\Cmpe493\Assignment2> python .\queryprocessing.py
Please enter query:"sugar price"
['9470', '15804']
PS C:\Users\muhammed\Desktop\Ders\7.Dönem\Cmpe493\Assignment2>
```

"leverage position"

```
PS C:\Users\muhammed\Desktop\Ders\7.Dönem\Cmpe493\Assignment2> python .\queryprocessing.py
Please enter query:"leverage position"
['6001']
PS C:\Users\muhammed\Desktop\Ders\7.Dönem\Cmpe493\Assignment2>
```

"United States District"

```
PS C:\Users\muhammed\Desktop\Ders\7.Dönem\Cmpe493\Assignment2> python .\queryprocessing.py
Please enter query:"United States District"
['6002', '15355', '18573', '19008']
PS C:\Users\muhammed\Desktop\Ders\7.Dönem\Cmpe493\Assignment2>
```
▶ Run   ≔ TODO   ❶ Problems   ⊡ Terminal   ≋ Python Packages   ⬡ Python Console

Muhammed Göktepe
2017400162

cocoa export shipment tonne

Oil price

United States

Muhammed Göktepe
2017400162

payments market