

CMPE 493 INFORMATION RETRIEVAL

ASSIGNMENT 2: A Simple Search System for Phrase and Free Text Queries

In this assignment we built a simple information retrieval system for phrase and free text queries with tf-idf and cosine similarity approach.

Index_building module is working with:

```
python3 index_build.py ./reuters21578/ ./stopwords.txt
```

As the index is built, Query_processing module is working with:

```
python3 query_process.py
```

Assignment consists of 3 main steps:

- 1- Data Preprocessing
- 2- Index Building
- 3- Query Processing, Phrase Search and Free Text Search

1- Data Preprocessing

Reuters data was read via an HTML file processing library, BeautifulSoup. As the documents are read, they are trimmed from punctuations according to the string.punctuation list and cleaned from stop words specified by the user. Latin-1 encoding was used for the documents. Multiprocessing which creates a process for each document was used to speed up the data reading at data preprocessing step.

2- Index Building

I built two different indexes that will be used for phrase and free text queries. For phrase queries I used a nested dictionary structure in which words are hashed to dictionary of documents which are hashed to positions of the corresponding word in the particular document. ($\{\text{word1}:\{\text{doc1}:[\text{pos1},\text{pos2}],\text{doc2}[\text{pos1},\text{pos2}]\}\}$) For free text queries I built three different dictionaries, one nested structure for tf_idf values for words in documents and another dictionary structure for idf values of all the words and one for the document norms.

I preferred dictionary (hashmap structure in python) for indexes because it uses a hash function to link key and values and it provides an $O(1)$ time complexity to store and retrieve a particular object. I preferred a basic array structure for positions list as it has $O(1)$ complexity for adding new elements. It has $O(n)$ complexity for searching items but as the positions lists of words are not large, I thought it won't create a problem during query search. Also, during the index building step, I didn't use any sorting function, which made the process a lot faster.

```
Last login: Tue Dec  7 16:41:04 on ttys001
[maral@192 ~ % cd /Users/maral/Documents/Classes/INFO_RETRIEVAL/Assignment2_tf_idf_cosine_search_engine
maral@192 Assignment2_tf_idf_cosine_search_engine % python3 index_build.py ./reuters21578/ ./stopwords.txt
[- Index is created in 7.914770 seconds -
maral@192 Assignment2_tf_idf_cosine_search_engine % █
```

Figure 1 (Index Build step takes 7.914770 second for the given reuters data.)

3- Query Processing

At the last step of the assignment, query processing, a phrase query search engine and a free text query search engine is created. Query input is taken from the user and preprocessed, cleaned from stop words and punctuations before it was sent to the engines. Phrase query search engine performs a conjunction operation in the documents and positions level. Free text query search engine uses the tf-idf values of the query terms and tf-idf values of the same terms in the document. Vectors for the query and the documents is created consisting of the tf-idf values. Finally, vectors are compared in terms of their cosine similarity and the list of documents sorted by their cosine similarity is displayed along with their cosine similarity score. For cosine similarity calculation, I computed the tf-idf values for queries and took the tf-idf values of

Maral Dicle Maral

03.12.2021

documents' query words from the `tf_idf_index` file. Also, during cosine similarity calculation, because the norm of document vectors are not the same with the query vectors, when I took only the query terms into account, this created a problem with the cosine similarity calculation. Therefore I used the actual norms of documents from `norm_index`.

```
maral@192 Assignment2_tf_idf_cosine_search_engine % python3 query_process.py
Index is successfully loaded!
Please enter a query (Use "w1 w2...wn" for phrase queries, w1 w2...wn for free text queries.). Enter E to exit.
Please enter a query('E' for exit):"James Baker"
"James Baker" is found at:
[52, 175, 190, 348, 458, 854, 965, 1357, 1392, 2052, 2078, 2190, 2452, 2620, 2632, 2633, 2648, 2701, 3421, 4008, 4038, 4139, 4764, 5176, 5186, 5290, 5294, 5759, 5869, 5944, 5964, 59
78, 6016, 6028, 6032, 6061, 6066, 6285, 6337, 6631, 6896, 7135, 7343, 7560, 7485, 7493, 7651, 7669, 7695, 7706, 7721, 7811, 8072, 8077, 8097, 8118, 8189, 8195, 8202, 8293, 8252, 830
9, 8318, 8315, 8483, 8548, 8621, 8664, 8714, 8963, 9022, 9055, 9134, 9220, 9498, 9501, 9515, 9524, 9533, 9535, 9603, 9628, 9656, 9657, 9675, 9689, 9690, 9693, 9694, 9697, 9764, 9788
, 9795, 9797, 9798, 9813, 9816, 9871, 9946, 10733, 10780, 10799, 11209, 11281, 11297, 11314, 11322, 11418, 11430, 11444, 11446, 11450, 11455, 11460, 11551, 11734, 11753, 11998, 1202
7, 12047, 12070, 12092, 12120, 12145, 12333, 12650, 12653, 12792, 12806, 12848, 13046, 13144, 13278, 13347, 13631, 13657, 13982, 14043, 14654, 14708, 14710, 14727, 14734, 14748, 147
57, 14780, 15048, 15453, 15470, 15549, 15656, 15957, 16006, 16051, 16075, 16092, 16200, 16540, 16780, 16852, 17126, 17138, 17200, 17218, 17231, 17268, 17296, 17320, 17321, 17325, 17
363, 17368, 17420, 17953, 17976, 18005, 18091, 18154, 18250, 18299, 18404, 19993, 20038, 20053, 20071, 20087, 20325, 20447, 20496, 20500, 20631, 20656, 20739, 20764, 20772, 20780, 2
0795, 20868, 20893, 20907, 20925, 21138, 21187, 21202, 21277, 21303, 21339, 21422, 21477, 21508, 21512, 21521, 21542, 21556, 21573]
- Process finished in 0.002406 seconds -
Please enter a query('E' for exit):
```

Figure 2 (Phrase query search for "James Baker", sorted by ID's.)

```
-----
Please enter a query('E' for exit):"fourth quarter performance"
"fourth quarter performance" is found at:
[18, 1019, 2130, 4165]
- Process finished in 0.003970 seconds -
Please enter a query('E' for exit):cocoa export shipment tonne
cocoa export shipment tonne is found at:
[(14372, 0.12379653032954416), (1, 0.1133386991237559), (5258, 0.09867527070539589)]
- Process finished in 0.002124 seconds -
Please enter a query('E' for exit):E
Thank you for using my search engine!
maral@192 Assignment2_tf_idf_cosine_search_engine %
```

Figure 3 (Free text query search for cocoa export shipment tonne, sorted by cosine similarity. First element of the tuples represents the document ID's, second element represent the cosine similarity.)