

## ReadMe and Report

To run this application you need to have these libraries, files and directories:

Task\_1:

Libraries: os, nltk, collections, json, numpy

Directory: scrap – corpus data

Task\_2:

Libraries: os, nltk, collections, itertools, numpy

Directory: scrap – corpus data

Task\_3:

Libraries: os

Directory: inverted\_index

Files: index\_1, index\_2, index\_3 inverted index files generated in task 1- a.

The code is in python 3.7.2. Make sure you have scrap named directory and provided corpus data in it.

To run each task run python Task\_X.py where X is the number of task. E.g. python Task\_1.py for task 1.

From design choice point of view I chose python inbuilt default dictionary as indexing data structure because of  $O(1)$  access. It has terms as keys and list as values and again within that list it has dictionary objects with document id: term position or term frequency based on the inverted index. For inverted index in task 1 it was

**$\{ "term" : [ \{ "document1": "term\_frequency" \}, "document2": "term\_frequency" ] ..... \}$**  and in

delta encoding it was

**$\{ "term" : [ \{ "document1": "first\_position, , d\_gap, , d\_gap" \}, "document2": "first\_position, d\_gap, d\_gap" ] ..... \}$**

Output:

**Task1** has output files for inverted index from 1-a named as **index\_1.txt**, **index\_2.txt**, **index\_3.txt** for unigram, bigram and trigram respectively in *inverted\_index* directory. It has [documentId, term] table in **document\_terms\_task1\_b.txt** and last document for task 1-d the delta encoded inverted index as **inverted\_index\_delta\_encoded.txt**

---

**Task2** has output files are in directory *Proximity\_query* and named following:  
Space , mission: k= 6 File name: **documents\_with\_space\_mission\_6.txt**  
Space , mission: k= 12 File name: **documents\_with\_space\_mission\_12.txt**  
earth , orbit: k= 5 File name: **documents\_with\_earth\_orbit\_5.txt**  
earth , orbit: k= 10 File name: **documents\_with\_earth\_orbit\_10.txt**

It is just a list of documents.

For design choice I used casefolding for each input.

---

**Task3** has output files in directory stats and are named following:  
Document frequency: **document\_frequency\_1.txt**, **document\_frequency\_2.txt**, **document\_frequency\_3.txt** for unigram, bigram and trigram respectively.  
Document frequency is sorted by terms and is in format:

**Term Dcoument1, Document2, Document3 : Frequency**

Term Frequency: **term\_frequency\_1.txt**, **term\_frequency\_2.txt**, **term\_frequency\_3.txt** for unigram, bigram and trigram respectively.

Term frequency is sorted by most to least and is in format :

**Term: Frequency**

Stop List: **stop\_list\_1.txt**, **stop\_list\_2.txt**, **stop\_list\_3.txt** for unigram, bigram and trigram respectively.

Stop list is in format:

**Term: Frequency**

Stop list explanation:

To generate a stop we can use various ways like tf-idf but after experimenting with few methods document frequency table gives good result. So, to choose a word as stop word in list we have to decide a cut off value in document frequency. After bit of experimentation I found that 800 is good value. It does not affect the context of the document and it mostly contain words which has only grammatical function or particular to the Wikipedia. E.g in trigram we have words like “from Wikipedia the”, “by using this” and “statement mobile

view”. We can see that this doesn’t have anything in common with the context of the document and its mostly for navigational purposes. Number of terms with increase in n decreases the list. E.g unigram has 160 stop words and trigram has 113 stop words.