

# Final Project

## Group 02

### Students Name:

1. Abdullah Alhoraibi

2. Mohamad Meaari

3. Ran Arino

In [1]: 1 `import os`

In [2]: 1 *#This function displays the menu as follows*  
2 *# 1. Indexing*  
3 *# 3. Exit*  
4  
5 `def printMenu():`  
6  `print('Menu:')`  
7  `print('Please enter 1 for indexing and 3 to exit')`  
8  `result_menu = int(input('1. Indexzing\n3. Exit\n'))`  
9  `return result_menu`  
10 `printMenu()`

Menu:  
Please enter 1 for indexing and 3 to exit  
1. Indexzing  
3. Exit  
1

Out[2]: 1

```

In [3]: 1 # This function takes a text file as input and replaces all punctuations int
        2 # Input: text
        3 # Output: text with no punctuations
        4
        5 def punctuationsRemoval(text: str) -> str:
        6     """
        7     (str) -> str
        8     Return a text after replacing any punctuations to blank, " ".
        9
       10     >>>punctuationsRemoval('abcdefg##@@higklmn,*()%opqr""stu')
       11     'abcdefg  higklmn  opqr  stu'
       12     """
       13     PUNCTUATIONS = ["!", "(", ")", "-", " ", "\n", "[", "]", "{", "}", ";", "
       14
       15     clear_text = ""
       16     for word in text:
       17         clear_text += word if word not in PUNCTUATIONS else " "
       18         # if a variable "word" is not punctuation, then append "word" it
       19
       20     return clear_text
       21
       22 punctuationsRemoval('abcdefg##@@higklmn,*()%opqr""stu')

```

Out[3]: 'abcdefg higklmn opqr stu'

```

In [4]: 1 # This function takes a text as input and removes all stopwords.
        2 # Input: text
        3 # Output: text with nostop words
        4 def stopWordRemoval(text: str) -> list:
        5     """
        6     (str) -> list
        7     Return word list based on the input data(txt format), which are excluded
        8     Before running this function, should be removed punctuations from the te
        9
       10     >>>stopWordRemoval("The monkeys jump on the bed.")
       11     ['monkeys', 'jump', 'bed.']
       12     """
       13     word_list = list(filter(lambda word: word != "", text.split(" ")))
       14     # excluding "" from the text after splitting by the blnk
       15
       16     with open('Stop_Words.txt', 'r') as f: # Load stopwords file
       17         Stop_Words = f.read()
       18
       19     remove_list = [i.strip("'").strip('"') for i in Stop_Words.split(", ")]
       20     clean_list = list(filter(lambda word: word.lower() not in remove_list, w
       21
       22     return clean_list
       23
       24 stopWordRemoval("The monkeys jump on the bed.")

```

Out[4]: ['monkeys', 'jump', 'bed.']

The following functions takes a clean text as argument, and appends TermDocFreqFile with the list of terms, the document in which they appear and their frequencies. termDocFreqFile format: 3 columns, values separated with space Term doc# freq Ontario 1 5 years 1 1 sugar 1 2

Input: cleanText, and document id output: termDocFreqFile format: 3 columns, values separated with space

```
In [5]: 1 def appendTermDocFreq(docid: int, cleanText: str, termDocFreqFile):
2         """
3         Appends TermDocFreqFile with the term(lowercase), the document number, a
4         The format is like below.
5
6         ontario 1 2\n
7         government 1 3\n
8         """
9         term_freq = {} # format -> {term: frequency}
10        for word in cleanText:
11            word = word.lower() # every word changes lower case
12            if word not in term_freq:
13                term_freq[word] = 1
14            else:
15                term_freq[word] += 1
16
17        append_text = ''
18        for k, v in term_freq.items():
19            append_text += '{} {} {}\n'.format(k, docid, v)
20
21        termDocFreqFile.write(append_text)
```

The following function reads termDocFreqFile line by line and append the global index that go from terms as keys to list of documents that contain them with their frequencies as val. Appending the index works as follows: - for line in termDocFreqFile, - if the term does not exist in index, add the term as key, the value will be a dictionary containing docid:freq as key:val in index - if the term already exists in index, append the val (which is a dictionary) with docid:freq input: termDocFreqFile output: fill in the index structure global variable defined in top of the module.

Note:

In order to read or write the latin alphabets (like é), we set the following statement in open() function: **encoding='utf-8'**.

```
In [6]: 1 def genIndex(termDocFreqFile):
2         index_file = {}
3         # Format -> {term_01: {doc#: freq, doc#: freq,...}, term_02:{doc#: freq,
4
5         termDocFreqFile = open("TermDocFreq.txt", 'r', encoding='utf-8')
6         for line in termDocFreqFile: # read text document line by line
7             read = line[:-1].split(" ") # apply split method after removing the
8             if read[0] not in index_file: # read[0], read[1], read[2] = term, do
9                 index_file[read[0]] = {read[1]: read[2]}
10            else:
11                index_file[read[0]][read[1]] = read[2]
12
13        return index_file
```

The following function reads all the text files in the folder 'dataset', appends them to a list, and returns the list. Input: None Output: a list of texts

```
In [7]: 1 def readFolderContent():
2         files = []
3         file_list = os.listdir('dataset')
4         for filename in sorted(file_list):
5             with open('dataset' + '/' + filename, 'r', encoding='utf-8') as infi
6                 files.append(infile.read())
7         return files
```

The following function creates necessary files needed in this project. For more information about this function review the flowchart given in the instructions.

```

In [9]: 1 def indexing():
2         termDocFreqFile = open("TermDocFreq.txt", 'w', encoding='utf-8')
3
4         # readFolderContent is called to create a list of files.
5         files = readFolderContent()
6         id=1
7         for file in files:
8             puncRemoved = punctuationsRemoval(file) # remove all punctuations
9             stopWordsRemoved = stopWordRemoval(puncRemoved) # remove all stop wo
10            appendTermDocFreq(id, stopWordsRemoved, termDocFreqFile) # Call app
11            id += 1
12
13            global global_index_file
14            global_index_file = genIndex(termDocFreqFile) # Call genIndex function t
15            termDocFreqFile.close()
16
17 def main():
18     option=printMenu()
19     if option == 1:
20         indexing()
21
22
23 if __name__ == "__main__":
24     main()

```

Menu:

Please enter 1 for indexing and 3 to exit

1. Indexzing

3. Exit

1

```

In [10]: 1 global_index_file

```

```

Out[10]: {'ontario': {'1': '2',
    '11': '1',
    '12': '1',
    '14': '1',
    '21': '1',
    '22': '1',
    '25': '1',
    '30': '2',
    '32': '1',
    '35': '2',
    '42': '4',
    '56': '2',
    '65': '6',
    '94': '8',
    '121': '1',
    '130': '4',
    '137': '2',
    '146': '1',
    '168': '2',
    '170': '1'

```

