

IR Assignment 1

Aditya Dixit(MT21004)

Shivang Kanwar(MT21084)

Libraries Used: nltk, pandas, numpy, os, pickle, re

Question 1

(a) Carry out the suitable preprocessing steps on the given dataset.

We have done our assignment on google colab. So after downloading the dataset we have extracted it from zip file and uploaded it on google drive. After that, we have mounted our drive on colab.

Steps in Preprocessing are:-

- Converting to lowercase
- Removing underscore
- Converting number to text
- Removing punctuation
- Removing stopword
- Removing blank space
- Removing single character
- Removing non-ASCII words
- Doing Lemmatization

```
1 ##### PART A #####
2 #convert text to lowercase
3 def convert_lower_case(text):
4     return np.char.lower(text)
5
6 #remove underscore
7 def rem_underscore(text):
8     t=str(text)
9     return t.replace("_","")
10
11 #Remove stopwords from text
12 def remove_stop_words(text):
13     stop_words = stopwords.words('english')
14     token = word_tokenize(str(text))
15     text2 = ""
16     for t in token:
17         if t not in stop_words:
18             text2 = text2+" "+t
19     return np.char.strip(text2)
20
21 #Remove punctuation marks from text
22 def remove_punctuation(text):
23     symbols = "/:;!\"#+-.&()*^_`<=>?@[\\
24     for i in range(len(symbols)):
25         text = np.char.replace(text, symbols[i], ' ')
26         text = np.char.replace(text, " ", " ")
27     text = np.char.replace(text, ',', ' ')
28     text = np.char.replace(text, '"', "")
29     return text
30
```

```

31 #Remove blank space
32 def remove_blank_space(text):
33     t=str(text)
34     return " ".join(t.split())
35
36 #Remove single characters
37 def remove_single_char(text):
38     words = word_tokenize(str(text))
39     text2 = ""
40     for w in words:
41         if len(w) > 1:
42             text2 = text2 + " " + w
43     return np.char.strip(text2)
44
45 #convert numbers to text
46 def convert_number_to_text(text):
47     text = np.char.replace(text, "0", " zero ") #replace 0 by zero in text.
48     text = np.char.replace(text, "1", " one ") #replace 1 by one in text.
49     text = np.char.replace(text, "2", " two ") #replace 2 by two in textd.
50     text = np.char.replace(text, "3", " three ") #replace 3 by three in text.
51     text = np.char.replace(text, "4", " four ") #replace 4 by four in text.
52     text = np.char.replace(text, "5", " five ") #replace 5 by five in text.
53     text = np.char.replace(text, "6", " six ") #replace 6 by six in text.
54     text = np.char.replace(text, "7", " seven ") #replace 7 by seven in text.
55     text = np.char.replace(text, "8", " eight ") #replace 8 by eight in text.
56     text = np.char.replace(text, "9", " nine ") #replace 9 by nine in text.
57     return text
58

```

```

59 #performing lemmatization on data
60 def lem(text):
61     lemmatizer = WordNetLemmatizer()
62     tokens = word_tokenize(str(text))
63     text2 = ""
64     for w in tokens:
65         text2 = text2 + " " + lemmatizer.lemmatize(w)
66     return np.char.strip(text2)
67
68 #removing non-ascii characters
69 def remove_non_ascii(text):
70     t=str(text)
71     return t.encode("ascii","ignore").decode("utf-8","ignore")
72
73 #perform preprocessing
74 def preprocess_data(text):
75     text = convert_lower_case(text)
76     text = rem_underscore(text)
77     text = convert_number_to_text(text)
78     text = remove_punctuation(text)
79     text = remove_stop_words(text)
80     text = remove_blank_space(text)
81     text = remove_single_char(text)
82     text = remove_non_ascii(text)
83     text = lem(text)
84     return text

```

(b) Implement the unigram inverted index data structure.

Steps in creating inverted index are:-

- Taking the data from the file
- Cleaning data with preprocessing
- Tokenization of the data
- Creating index with the dictionary. Words as key and document number as values
- Storing the index in a pickle file

For eg:- Data:- "lion is happy"

Preprocessed:- "lion happy"

Tokens:-["lion","happy"]

Index:- {"lion":[2,3,45], "happy":[3,5,78]}

```
##### PART B #####
def create_inverted_index(paths):    #function create inverted index.
    store_filename = []              #list for storing all file names
    docID = 0                        #counter for documents
    tii2=dict()                      # dictionary for storing index with word as key and position as values

    #loop for calculating index for each file
    for path in paths:
        file = open(path, 'r', encoding= 'ISO-8859-1')
        text = " ".join(file.read().split())    # for removing extra spaces
        file.close()
        print(docID)
        preprocessed_data = preprocess_data(text)
        tokens = word_tokenize(str(preprocessed_data))    #tokenize the preprocessed text and convert into tokens

    #loop for storing data in dictionary
    for t in tokens:
        if t not in tii2.keys():
            tii2[t] = list()
            tii2[t].append(docID)
        else:
            if docID not in tii2[t]:
                tii2[t].append(docID)

        filename = os.path.basename(path)    #extract last filename from the path
        store_filename.append([filename])    #append into list
        docID += 1

    store_filename = pd.DataFrame(store_filename)    #convert filenames list of list into dataframe

    # save the inverted index and filenames in pickle file so that can we build further
    with open('fn.pickle', 'wb') as handle:
        pickle.dump(tii2, handle, protocol=pickle.HIGHEST_PROTOCOL)

    store_filename.to_pickle("sf")
```

(c) Provide support for the following queries-

(i) x OR y

(ii) x AND y

(iii) x AND NOT y

(iv) x OR NOT y

We have implemented 5 functions AND, OR, NOT, AND NOT, OR NOT to perform operations on queries.

NOT, OR OPERATION

```
##### PART C #####
def not_operation(var):
    #function to perform not operation
    v = set(range(len(paths)))
    return list(v.difference(var))

def or_operation(var1,var2):
    #function to perform or operation
    #store length of var1 and var2 into n and m.
    n = len(var1)
    m = len(var2)
    Result = []
    comparisons = 0
    i = 0
    j = 0

    while i<n and j<m:
        #perform union in two set X and Y and count number of comparision.
        if var1[i]==var2[j]:
            Result.append(var1[i])
            i=i+1
            j=j+1
            comparisons=comparisons+1
        elif var1[i]<var2[j]:
            Result.append(var1[i])
            i=i+1
            comparisons=comparisons+1
        else:
            Result.append(var2[j])
            j=j+1
            comparisons=comparisons+1

    while i<n:
        #when only var1 set left then append into result without any comparision.
        Result.append(var1[i])
        i=i+1

    while j<m:
        #when only var2 set left then append into result without any comparision.
        Result.append(var2[j])
        j=j+1

    return Result, comparisons
    #return union of var1 and var2 and number of comparision.
```

AND, AND NOT, OR NOT OPERATION

```
def and_operation(var1,var2):          #function to perform and operation
    n = len(var1)                     #store length of var1 and var2 into n and m.
    m = len(var2)
    Result = []
    comparisons =0
    i = 0
    j = 0

    while i<n and j<m:                #perform union in two set var1 and var2 and count number of comparison.
        if var1[i]==var2[j]:
            Result.append(var1[i])
            i=i+1
            j=j+1
            comparisons=comparisons+1
        elif var1[i]<var2[j]:
            i=i+1
            comparisons=comparisons+1
        else:
            j=j+1
            comparisons=comparisons+1
    return Result, comparisons        #return intersection of var1 and var2 and number of comparison.

def or_not_operation(var1, var2):      #function to perform or Not operation
    var2 = not_operation(var2)        #find not of var2.
    Result, comparisons = or_operation(var1, var2) #after finding not of var2 then perform OR of var1 and var2 and find no of comparison.
    return Result, comparisons        #return the result after var1 or Not var2 and number of comparison.

def and_not_operation(var1, var2):    #function to perform and not operation
    var2 = not_operation(var2)        #find not of var2.
    Result, comparisons = and_operation(var1, var2) #after finding not of Y then perform and of var1 and var2 and find number of comparison.
    return Result, comparisons        #return the result after var1 or Not var2 and number of comparison.
```

Next we have implemented a function called `cal_result()` which drives the main query processing. It takes document list for each word and perform operation on them. It uses the `and`, `or`, `not`, `and not`, `or not` function by calling them and then storing the result.

```
1 def cal_result(Query_docIds):
2     op_pos = 0                      #iterator for positions
3     Total_comp = 0                  #to store total comaprison
4     X = Query_docIds[0]
5     for operand in Query_docIds[1:]:
6         Y = operand
7
8         operator = operations[op_pos]
9         op_pos = op_pos+1
10
11         if(operator == 'OR'):        #if operator is OR then call Function_or and find resultant list and number of comparison.
12             X, No_of_compare = or_operation(X, Y)
13             Total_comp = Total_comp + No_of_compare
14
15         elif(operator == 'AND'):      #if operator is AND then call Function_or and find resultant list and number of comparison.
16             X, No_of_compare = and_operation(X, Y)
17             Total_comp = Total_comp + No_of_compare
18
19         elif(operator == 'OR NOT'):   #if operator is OR NOT then call Function_or and find resultant list and number of comparison.
20             X, No_of_compare = or_not_operation(X, Y)
21             Total_comp = Total_comp + No_of_compare
22
23         elif(operator == 'AND NOT'):  #if operator is AND NOT then call Function_or and find resultant list and number of comparison.
24             X, No_of_compare = and_not_operation(X, Y)
25             Total_comp = Total_comp + No_of_compare
26
27         else:
28             print("You Entered Wrong Operations.....!!!")
29
```

Main function is the starting function of the program. This function is executed first. Firstly it stores all paths in a path list, which is then passed to the `create_inverted_index()` function to create index. After this input to the program from user is taken. Query and operations to be performed are taken and after doing operations result is shown to the user.

CREATING INDEX

```
1 if __name__ == "__main__":
2
3     files_path = "/content/drive/MyDrive/Humor,Hist,Media,Food"
4
5     # find all files those are stored into files_path and store into path.
6     paths = []
7     for (dirpath, dirnames, filenames) in os.walk(str(files_path)):
8         for i in filenames:
9             paths.append(str(dirpath)+str("/") + i)
10
11     create_inverted_index(paths)    #for creating index
12     store_filename = []            #for storing filenames
13
14     # read the index and filename
15     with open('fn.pickle', 'rb') as handle:
16         inverted_index = pickle.load(handle)
17     store_filename=pd.read_pickle("sf")
18
```

TAKING INPUT AND PERFORMING OPERATIONS AND PRINTING RESULT

```
18
19 # Our input query
20 N = int(input("Enter Number of Queries you want to test:-"))
21 while N>0:
22     Query = input("Enter the query:-")
23     operations = input("Enter operations with comma between them:-")
24     operations=operations.upper()
25     operations=operations.split(",")
26     operations = [i.strip() for i in operations]
27
28 # Query Pre-processing
29 pre_processed_query = preprocess_data(Query)
30 pre_processed_query = pre_processed_query.flatten()
31 pre_processed_query = np.char.split(pre_processed_query[0])
32 pre_processed_query = pre_processed_query.tolist()
33
34 #add word documents id to query document ids in sorted order
35 Query_docIds = []
36
37 for r in pre_processed_query:
38     if r not in inverted_index.keys():
39         print(r+" is not in any document. Please enter valid query")
40         break
41     word_docIds = inverted_index[r]
42     Query_docIds.append(sorted(word_docIds))
43
44 #check if processed query documents are 1 less than number of operations
45 if(len(Query_docIds)-1 == len(operations)):
46     cal_result(Query_docIds)
47 else:
48     print("Wrong Input")
49     N-=1
```

OUTPUT:-

Enter Number of Queries you want to test:-2
Enter the query:-lion stood thoughtfully for a moment
Enter operations with comma between them:-or,or,or

Number of documents matched: 211

No. of comparisons required: 336

Documents retrieved are:-

['conan.txt', 'coyote.txt', 'incarnel.hum', 'cartoon_.txt', 'ivan.hum', 'lbinter.hum', 'murphys.txt', 'li

Enter the query:-telephone,paved, roads

Enter operations with comma between them:-or not,and not

Number of documents matched: 996

No. of comparisons required: 2244

Documents retrieved are:-

['bbq.txt', 'gohome.hum', 'mothers.txt', 'harmful.hum', 'hate.hum', 'herb!.hum', 'mowers.txt', 'hell.jok'

Question 2

(a) Carry out the following preprocessing steps on the given dataset

Before preprocessing, we manually extracted the folder and then retrieved all files and stored their paths in the list.

```
✓ [9] files_path = "/content/drive/MyDrive/Humor,Hist,Media,Food"
0s

#Store path of all files in paths
paths = []
for (dirpath, dirnames, filenames) in os.walk(str(files_path)):
    for i in filenames:
        paths.append(str(dirpath)+str("/") + i)

print(len(paths))

1133
```

- (i). In preprocessing, first converted the text of each file into lower case.
- (ii). Then converted this text into tokens using word_tokenize from nltk library.
- (iii). Removed all stopwords from these tokens using stopwords of nltk library.
- (iv). Removed punctuations from above filtered tokens using re.sub() which replaced all the char except a-z A-Z 0-9 _
- (v). Finally removed all the blank spaces from these tokens.

```
##### part a #####

#Convert the text to lower case
def convert_lower_case(text):
    return np.char.lower(text)

#Perform word tokenization
def word_tokenization(text):
    tokens = word_tokenize(str(text))
    return tokens

#Remove stopwords from tokens
def remove_stop_words(tokens):
    stop_words = stopwords.words('english')
    filtered_tokens = []
    for word in tokens:
        if word not in stop_words:
            filtered_tokens.append(word)
    return filtered_tokens

#Remove punctuation marks from tokens
def remove_punctuation(words):
    filtered_words = []
    for word in words:
        filtered_word = re.sub(r'^\w\s$', ' ', word) #\w for string containing char a-z,A-Z,0-9,_, \s for whitespaces, ^for except these
        if filtered_word != ' ':
            filtered_words.append(filtered_word)
    return filtered_words

#Remove blank space tokens
def remove_blank_space_tokens(words):
    words = ' '.join(words).split()
    return words

#Perform preprocessing
def preprocess_data(text):
    text = convert_lower_case(text)
    tokens = word_tokenization(text)
    filtered_tokens = remove_stop_words(tokens)
    filtered_tokens = remove_punctuation(filtered_tokens)
    filtered_tokens = remove_blank_space_tokens(filtered_tokens)
    return filtered_tokens
```

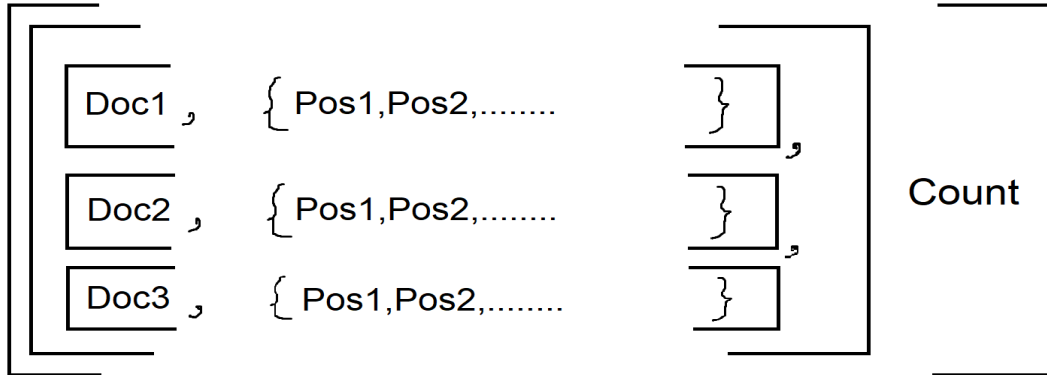
(b) Implement the positional index data structure

Created the positional index using the dictionary that has a list as a value for each token.

The list itself has two elements:

- (i) the list of lists having doc_no and tokens positions set in that doc
- (ii) the total occurrence of the token in all the docs

Also maintained dictionary doc_map to store doc_no, filename as key-value pair.



So iterated for each file, read the text, and preprocessed it. Then for all the tokens, updated the positional_index

```
positional_index = {}
doc_map = {}

for path in paths:
    file = open(path, 'r', encoding='ISO-8859-1')
    text = file.read().strip()
    file.close()
    tokens = preprocess_data(text)

    for position, term in enumerate(tokens):
        if term in positional_index:
            positional_index[term][1] += 1
            posting = positional_index[term][0]
            pos = find_pos_of_list_having_doc(posting, doc_no)
            if pos != "":
                positional_index[term][0][pos][1].add(position)
        else:
            positional_index[term][0].append([])
            positional_index[term][0][-1].append(doc_no)
            positional_index[term][0][-1].append(set())
            positional_index[term][0][-1][1].add(position)

    filename = os.path.basename(path)
    doc_map[doc_no] = filename
    doc_no += 1

df = pd.DataFrame(positional_index)
display(df)
```

#posting list of token

#if doc already added
#add position to set
#if doc added first time
#create new list
#append doc_no to the newly appended list(-1 is position of new list)
#append set for positions to the newly appended list
#add position to this newly appended set

#new list as a value for token
#list at 1st pos for storing posting lists for each doc
#occurrence count of token at 2nd position
#create new list
#append doc_no to the newly appended list(-1 is position of new list)
#append set for positions to the newly appended list
#add position to this newly appended set

Positional index after processing all docs and corresponding tokens.

	newsgroups	talk	bizarre	rigler	dao	nrc	ca	michael	subject	t	b	boxed	edition	message	id	1992dec11	033233	26164	sol	uvic	reply	t
0																						
1	145	412	65	7	6	7	1176	405	567	8326	1222	3	75	500	147	1	1	1	16	5	172	20
2 rows x 71867 columns																						
<div><div></div></div>																						

(c) Provide support for the searching of phrase queries

For processing phrase query search, first, we preprocessed the query with the same steps as we did for file processing. Now if the filtered query tokens have only one word then print the docs in its posting list and the count of these docs. For phrase query for more than one word we first find the list of (document, position) pair for first word using `find_doc_position_pair_list()` helper method. Now iterate for all remaining words in the phrase query and find if the next word is also in same document but has position one more than that of the previous word. If so then append the doc_no in the result doc list and finally return it as a set.

```
def positional_func(first_word_doc_position_set, query_tokens):
    matched_docs = []
    for a in first_word_doc_position_set:
        doc = a[0]
        pos = a[1]
        token_count = 0

        for token in query_tokens:
            pos = pos+1
            token_posting = positional_index[token][0]
            token_docs = [a[0] for a in token_posting]
            if doc in token_docs:
                doc_positions = find_word_positions_in_doc(token_posting, doc)
                if pos in doc_positions:
                    token_count += 1
            else:
                token_count += 1
                break
        if token_count == len(query_tokens):
            matched_docs.append(a[0])

    return set(matched_docs)
```

```

##### part c #####

def check_tokens_in_index(tokens):
    for token in tokens:
        if token not in positional_index:
            return False
    return True

def find_word_positions_in_doc(posting_list, doc):
    for a in posting_list:
        if a[0] == doc:
            return a[1]          #return set of word positions in doc
    return {}

def find_doc_position_pair_list(word):
    doc_position_pair_list = []
    word_postings = positional_index[word][0]
    for a in word_postings:
        for position in a[1]:      #for each position of doc a[0]
            doc_position_pair_list.append((a[0], position))
    return doc_position_pair_list

```

Outputs:



starter_method()

Enter phrase query:



starter_method()

Enter phrase query: University of Maryland Baltimore County
 The number of documents retrieved: 2
 The list of document names retrieved: ['top10st1.txt', 'top10st2.txt']



starter_method()

Enter phrase query: welcome
 The number of documents retrieved: 94
 The list of document names retrieved: ['st_silic.txt', 'top10st2.txt', 'top10st1.txt', 'top10.txt',