# IR ASSIGNMENT 1

PATEL DHRUV (MT23056)

## Q1) Preproccesing:

```
remove HTML Tags

1   from bs4 import BeautifulSoup
2   import os
3
4   def remove_html_tags(text):
5       soup = BeautifulSoup(text, "html.parser")
6       return soup.get_text()
7
8   def process_files(input_directory):
9       file_list = os.listdir(input_directory)
10
11      for file_name in file_list:
12          if file_name.endswith(".txt"):
13              file_path = os.path.join(input_directory, file_name)
14
15              with open(file_path, 'r', encoding='utf-8') as file:
16                  content = file.read()
17                  cleaned_content = remove_html_tags(content)
18
19              with open(file_path, 'w', encoding='utf-8') as file:
20                  file.write(cleaned_content)
21
22              print(f"Processed: {file_name}")
23
24  input_directory = "/content/drive/MyDrive/text_files"
25
26  process_files(input_directory)
```

This Python script utilizes the BeautifulSoup library to remove HTML tags from the content of text files within a specified directory. The function `remove_html_tags` takes a text input, parses it using BeautifulSoup, and extracts the plain text content by discarding any HTML tags. The `process_files` function iterates through the files in the input directory, specifically targeting those with a ".txt" extension. For each text file, it reads the content, removes HTML tags using the aforementioned function, and then overwrites the original file with the cleaned content. The script provides a simple yet effective way to preprocess text files, eliminating HTML markup and producing cleaner versions of the documents, which can be beneficial for various text analysis tasks. Additionally, it prints a message for each processed file to indicate the progress.

```
1   from nltk.tokenize import word_tokenize
2   import nltk
3   nltk.download('punkt')
4
5   print("Before Lowercase:\n")
6   Print('/content/drive/MyDrive/text_files/')
7
8
9   for file in dir_list:
10      input_file_path = '/content/drive/MyDrive/text_files/' + file
11
12      with open(input_file_path, 'r') as input_file:
13          lines = input_file.readlines()
14          tokenized_lines = [word_tokenize(line) for line in lines]
15
16      with open(input_file_path, 'w') as output_file:
17          for tokenized_line in tokenized_lines:
18              output_file.write(' '.join(tokenized_line) + '\n')
19
20  print("After Lowercase:\n")
21  Print('/content/drive/MyDrive/text_files/')
22
```

This Python script uses the Natural Language Toolkit (nltk) to tokenize words in text files. First, it imports the necessary modules, including `word_tokenize` from nltk, and downloads the required punkt tokenizer. Then, the script attempts to print "Before Lowercase" and the content of the directory '/content/drive/MyDrive/text_files/' using a non-existent function `Print`, which might be a typo, as the correct function name is `print`.

Next, it iterates through the list of files in the directory, and for each file, it reads the content, tokenizes the words in each line using nltk's `word_tokenize` function, and writes the tokenized lines back to the same file. Essentially, this script tokenizes the words in each line of the text files and updates the files with the tokenized content. However, there's an issue with the function name 'Print' instead of 'print', and the script prints the same directory path both before and after tokenization without displaying the actual content. Additionally, the script may need some modifications to correctly execute.

```
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
import nltk
nltk.download('stopwords')


print("Before Lowercase:\n")
Print('/content/drive/MyDrive/text_files/')

def remove_stopwords(data):
    stop_words = set(stopwords.words("english"))
    tokens = word_tokenize(data)
    without_stopwords = [w for w in tokens if w.lower() not in stop_words]
    return without_stopwords

for file in dir_list:
    input_file_path = '/content/drive/MyDrive/text_files/' + file

    with open(input_file_path, 'r') as input_file:
        lines = input_file.readlines()
        modified_lines = [" ".join(remove_stopwords(line)) for line in lines]

    with open(input_file_path, 'w') as output_file:
        for modified_line in modified_lines:
            output_file.write(modified_line + '\n')


print("After Lowercase:\n")
Print('/content/drive/MyDrive/text_files/')
```

This Python script utilizes the Natural Language Toolkit (nltk) to remove stopwords from text files stored in the directory '/content/drive/MyDrive/text_files/'. It first imports the necessary modules, including the stopwords list and the word_tokenize function. The script then attempts to print "Before Lowercase" and the content of the specified directory using a non-existent function `Print` (likely a typo, should be `print`). The core functionality involves iterating through each file in the directory, reading its content, tokenizing the words in each line using nltk's `word_tokenize` function, and then removing stopwords using a custom function `remove_stopwords`. The modified lines without stopwords are then written back to the same file. This process effectively cleans the text data by eliminating common English stopwords. However, the script may need adjustments, such as fixing the 'Print' typo and ensuring the correct execution of the intended functionality.

```
4) punctuations

1    import string
2
3    print("Before Lowercase:\n")
4    Print('/content/drive/MyDrive/text_files/')
5
6
7    for file in dir_list:
8        input_file_path = '/content/drive/MyDrive/text_files/' + file
9
10       with open(input_file_path, 'r') as input_file:
11           data = input_file.read()
12
13       translator = str.maketrans('', '', string.punctuation)
14       data_without_punctuation = data.translate(translator)
15
16       with open(input_file_path, 'w') as output_file:
17           output_file.write(data_without_punctuation)
18
19   print("After Lowercase:\n")
20   Print('/content/drive/MyDrive/text_files/')
21
```

This Python script aims to preprocess text files in the directory '/content/drive/MyDrive/text_files/' by removing punctuation from the content. It first attempts to print "Before Lowercase" and the content of the specified directory using a non-existent function `Print` (likely a typo, should be `print`). The script then iterates through each file in the directory, reads its content, and utilizes the `str.maketrans` and `translate` methods to remove punctuation from the text. The modified content without punctuation is then written back to the same file. The script appears to focus on cleaning the text data by eliminating punctuation marks, which can be beneficial for various natural language processing tasks. However, it may need corrections, such as fixing the 'Print' typo and ensuring the correct execution of the intended functionality.

```
5) blank space tokens

[ ]   1   print("Before Lowercase:\n")
      2   Print('/content/drive/MyDrive/text_files/')
      3
      4
      5   for file in dir_list:
      6       input_file_path = '/content/drive/MyDrive/text_files/' + file
      7
      8       with open(input_file_path, 'r') as input_file:
      9           lines = input_file.readlines()
     10           modified_lines = [" ".join(line.split()) for line in lines]
     11
     12       with open(input_file_path, 'w') as output_file:
     13           for modified_line in modified_lines:
     14               output_file.write(modified_line + '\n')
     15
     16   print("After Lowercase:\n")
     17   Print('/content/drive/MyDrive/text_files/')
     18
```

This Python script aims to preprocess text files in the directory '/content/drive/MyDrive/text_files/' by removing extra whitespaces and converting multiple consecutive whitespaces into a single space. It begins by attempting to print "Before Lowercase" and the content of the specified directory using a non-existent function `Print` (likely a typo, should be `print`). The script then iterates through each file in the directory, reads its content, and applies a modification to each line by using a list comprehension that splits each line into words and joins them back together with a single space in between, effectively removing extra whitespaces. The modified lines are then written back to the same file. This process helps standardize the formatting of the text data by eliminating unnecessary whitespaces. However, the script may require corrections, such as fixing the 'Print' typo and ensuring the correct execution of the intended functionality.

## Q2) unigram indexing:

```
1    import os
2    import pickle
3
4    data_folder = "/content/drive/MyDrive/IR/text_files"
5    files_list = os.listdir(data_folder)
6    data = {}
7
8    for current_file in files_list:
9        with open(os.path.join(data_folder, current_file), "r") as file_content:
10           data[current_file] = file_content.read()
11
12
13   inverted_index_dict = {}
14
15   for current_file, file_text in data.items():
16       words = file_text.split()
17
18       for word in words:
19           if word not in inverted_index_dict:
20               inverted_index_dict[word] = set()
21
22           inverted_index_dict[word].add(current_file)
23
24   inverted_index_dict = {word: list(files) for word, files in inverted_index_dict.items()}
25
```

```
inverted_index_dict = {word: list(files) for word, files in inverted_index_dict.items()}


with open("/content/drive/MyDrive/Q2_pickle.pkl", "wb") as file_obj:
    pickle.dump(inverted_index_dict, file_obj)


with open("/content/drive/MyDrive/Q2_pickle.pkl", "rb") as file_obj:
    loaded_unigram_inverted_index = pickle.load(file_obj)

print(inverted_index_dict)
```
events': ['file57.txt'], 'busking': ['file475.txt'], 'laying': ['file163.txt', 'file475.txt'],

This Python code establishes an inverted index for a collection of text files stored in a specified folder. It first reads the content of each file into a dictionary where filenames serve as keys and content as values. Subsequently, it constructs an inverted index, associating each word in the files with a set of filenames where it appears. The resulting index is then converted from sets to lists for serialization. The code utilizes the pickle module to save the inverted index as a binary file and later loads it back. Finally, the original inverted index is printed, though there's a variable naming inconsistency as the loaded index is stored in loaded_unigram_inverted_index but not used in the print statement.

**Load Pickle file and save pickle file**

```
[4]  1  def load_pickle(filename):
     2      with open(filename, 'rb') as f:
     3          return pickle.load(f)
     4
     5  def save_pickle(obj, filename):
     6      with open(filename, 'wb') as f:
     7          pickle.dump(obj, f)
```

Loaded pickle file and open it.

```python
def docs(query, operations, unigram_index, all_files):
    lem = WordNetLemmatizer()
    query = query.lower()
    query = re.sub('[^A-Z a-z ]+', ' ', query)
    query = query.split()
    tokens = [j for j in query]
    stop_words = set(stopwords.words("english"))
    tokens = [lem.lemmatize(token) for token in tokens if token not in stop_words and token not in string.punctuation and token.strip()]

    print(tokens)

    combined_elements = []
    for token, operation in zip(tokens, operations):
        combined_elements.extend([token, operation])
    combined_elements.extend(tokens[len(operations):] + operations[len(tokens):])
    query_ans = ' '.join(combined_elements)

    query_tokens = tokens

    res_set = [list(unigram_index[token]) for token in query_tokens if token in unigram_index]

    result = []
    compare = 0
    i = 0

    for operations in operations:
        if res_set and i < len(res_set):
            # print(res_set)
            if operations == "AND":
                if i == 0:
```

```python
def file_names(num_files):
    return [f"file{i + 1}.txt" for i in range(num_files)]

def process_query(query, operations, inverted_index, files):
    result, number_of_comparisons, query_ans = docs(query, operations, inverted_index, files)
    query_statement = " ".join(operations)
    sorted_result = sorted(result)
    return query_statement, len(result), sorted_result, query_ans

if __name__ == "__main__":
    n = int(input("Enter the number of queries: ").strip())
    files = file_names(999)

    for i in range(n):
        query = input(f"Enter query {i + 1}: ").strip()
        operations = input("Enter operations separated by commas: ").strip().split(",")

        query_statement, num_retrieved, retrieved_names, query_ans = process_query(query, operations, loaded_unigram_inverted_index, file

        # print(f"\nQuery {i + 1}: {query_statement}")
        print(f"\nQuery {i + 1}: {query_ans}")

        print("Number of documents retrieved for query", i + 1, ":", num_retrieved)
        print("Names of the documents retrieved for query", i + 1, ":", retrieved_names)
        print()
```

This Python code defines functions for basic information retrieval operations on a set of text documents. The functions include `union`, `intersect`, `complement`, `andd`, and `orr` for set operations. The `docs` function processes a user query, tokenizes and lemmatizes it, and then applies the specified operations on the inverted index of the documents. The program takes user input for the number of queries, and for each query, it prompts the user to enter a search query and corresponding operations (e.g., AND, OR) separated by commas. The `process_query` function utilizes the `docs` function to execute the query, and the results, including the query statement, the number of retrieved documents, and their names, are then printed for each query. The code demonstrates a basic information retrieval system with set-based query operations on a collection of documents.

# Q3) Positional Indexing

```
[ ]   1   def positional(directory_path, num_docs):
      2       dict = {}
      3
      4       for doc_id in range(1, num_docs + 1):
      5           path = os.path.join(directory_path, f"file{doc_id}.txt")
      6
      7           if os.path.exists(path):
      8               with open(path, 'r') as file:
      9                   terms = file.read().split()
      10
      11                  for position, term in enumerate(terms, start=1):
      12                      if term not in dict:
      13                          dict[term] = {'docs': {doc_id: [position]}}
      14                      else:
      15                          dict[term]['docs'].setdefault(doc_id, []).append(position)
      16          else:
      17              print(f"File file{doc_id}.txt does not exist.")
      18
      19      return dict
      20
      21  # Path to the directory containing the text files
      22  directory_path = "/content/drive/MyDrive/IR/text_files/"
      23  # Number of documents
      24  num_docs = 999
      25
      26  # Build the positional index
      27  dict = positional(directory_path, num_docs)
      28
      29  # Print the positional index
      30  print(dict)
      31
```

This Python code defines a function positional that builds a positional index for a set of text documents. The function takes two parameters: directory_path, which is the path to the directory containing the text files, and num_docs, which is the total number of documents to be processed.

The code iterates through each document (indexed from 1 to num_docs) in the specified directory. For each document, it reads the content, tokenizes it into terms, and then constructs a positional index. The positional index is stored in a dictionary (dict), where each term is a key. For each term, the corresponding value is another dictionary with the key 'docs' representing the document IDs and their associated positions within the document.If a term is encountered for the first time, a new entry is added to the main dictionary with the term as the key. If the term already exists, the document ID and position information are appended to the existing entry.

```python
25   def retrived_docs(term, dict):
26       result_docs = set()
27
28       for term in term:
29           if term in dict:
30               docs_positions = dict[term]['docs']
31               result_docs.update(docs_positions.keys())
32
33       sorted_docs = sorted(result_docs)
34       return sorted_docs
35
36   # Process input queries
37   n_queries = int(input("Enter the number of queries: "))  # Read the number of queries
38   queries = [input("Enter query: ").strip() for _ in range(n_queries)]  # Read the queries
39
40   # Execute queries and get results
41   results = []
42
43   for i, query in enumerate(queries, start=1):
44       after_query = preprocess_query(query)
45       term = after_query.split()
46
47       if(len(term)>5):
48           print("Query contains more than 5 words")
49           continue
50
51
52       # Check if the first term is in the positional index
53       if term and term[0] not in dict:
54           print(f"Word '{term[0]}' not in dictionary for query {i}!")
55           continue
```

```python
62               break
63
64           new_index = {}
65           for doc_id in l_index:
66               if doc_id in dict[term]['docs']:
67                   found_positions = [pos for pos in l_index[doc_id] if pos + 1 in dict[term]['docs'][doc_id]]
68                   if found_positions:
69                       new_index[doc_id] = found_positions
70
71           l_index = new_index
72
73       documents_found = list(l_index.keys())
74       results.append(documents_found)
75
76   # Output results
77   for i, result in enumerate(results, start=1):
78       print(f"Number of docs retrieve for query {i} using positional index: {len(result)}")
79       if result:
80           print(f"Names of docs retrieve for query {i} using positional index: {', '.join(map(str, result))}")
81       else:
82           print("No documents found.")
83
```

```
Enter the number of queries: 2
Enter query: load is loaded
Enter query: great value
Number of docs retrieve for query 1 using positional index: 0
No documents found.
Number of docs retrieve for query 2 using positional index: 10
Names of docs retrieve for query 2 using positional index: 65, 103, 330, 466, 597, 748, 767, 789, 899, 993
```

This Python code defines a query processing system using a positional index built on a set of text documents. The program first includes a function preprocess_query to clean and preprocess user

input queries. This involves converting the query to lowercase, removing non-alphabetic characters, tokenizing, eliminating stop words and punctuation, and lemmatizing the tokens.

Another function, retrived_docs, is designed to retrieve documents based on the processed query terms from the positional index. It iterates through each term in the query, checking if it exists in the index, and then narrowing down the search by considering the positional information of each term in relation to the previous term. The retrieved documents are stored and returned.

The main part of the program takes user input for the number of queries, processes each query using the preprocess_query function, and then utilizes the retrived_docs function to retrieve relevant documents based on the positional index. The results are printed, including the number of retrieved documents and their names.