

IRTM Homework2 Report

B10705054 劉知微

1 Execution Environment

The code is executed in an Anaconda Jupyter Notebook environment.

2 Programming Language

The programming language used is Python 3.

3 Execution Method

To run the code in a Jupyter Notebook, no additional environment setup is required. However, if you are using VS Code or another IDE, you may need to install the necessary packages, such as `numpy` and `nltk`. The following modules should be included in the code:

```
1 from nltk.stem import PorterStemmer
2 import os
3 import numpy as np
4 from collections import defaultdict
5 import math
```

Please ensure that the `data` folder and the `stopwords.txt` file (which is included in the folder) are prepared in the same directory as the code. The folder and stopwords file paths should be defined as follows:

```
1 document_folder = "./data"
2 stopwords_file = './stopwords.txt'
```

Lastly, to execute the code in Jupyter Notebook, simply click the "Run" button (play icon) in the toolbar section. This will run the code in the active cell.

4 Workflow

4.1 Tokenization

To tokenize the text, I use the built-in function `split()` to separate the words based on spaces. The resulting words are saved as tokens. Then, the `lower()` function is applied to convert all letters in the tokens to lowercase. I utilize the `PorterStemmer()` from `nltk.stem` to apply the Porter Stemming algorithm. For each word in the tokens, if the word is alphanumeric and not in the stopwords set, I stem the word and save it into the list `filtered_tokens`.

```
1 # Function for tokenization, lowercasing, and stemming
2 def tokenization(text):
3     tokens = text.split()
4     tokens = [word.lower() for word in tokens]
5     porter_stemmer = PorterStemmer()
6     # Filter tokens: stem words, keep only alphanumeric tokens, and
7     # remove stopwords
8     filtered_tokens = [
9         porter_stemmer.stem(word)
10        for word in tokens
11        if word.isalnum() and word not in stop_words
12    ]
13    return filtered_tokens
```

I read the files from the folder `data`. Each file path is generated using `os.path.join()`, and the file is opened for reading. Every file goes through the previously mentioned `tokenization()` function, and the tokenized content is appended to the list `documents`.

```
1 # List to store tokenized documents
2 documents = []
3
4 # Read and process each file in the folder
5 for i in range(1, 1096):
6     file_path = os.path.join(document_folder, f"{i}.txt")
7     with open(file_path, 'r', encoding='utf-8') as file:
8         # Read the file content
9         content = file.read()
10
11        # Tokenize the content
12        processed_content = tokenization(content)
```

```
13
14     # Append the tokenized content to the documents list
15     documents.append(processed_content)
```

4.2 Constructing a Dictionary

I used the `defaultdict` class from the `collections` module to create a dictionary. First, I constructed a set of unique terms from each document, which is an essential step for counting document frequency (DF). This ensures that if a document contains two or more identical words, the DF count is not inflated. For each term, I incremented its corresponding value in `term_df`. The `defaultdict` is ideal for this process as it automatically initializes new terms with a default count of zero.

```
1 term_df = defaultdict(int)
2 for doc in documents:
3     unique_terms = set(doc) # Create a set of unique terms from each
        document
4     for term in unique_terms:
5         term_df[term] += 1 # Increment the DF count for each unique
            term
```

Next, I sorted the terms alphabetically and assigned an index to each term.

```
1 sorted_terms = sorted(term_df.items(), key=lambda x: x[0]) # Sort
        terms alphabetically
2 dictionary_entries = []
3 for idx, (term, df) in enumerate(sorted_terms, start=1): # Assign
        index starting from 1
4     dictionary_entries.append((idx, term, df)) # Create entries with
        index, term, and DF value
```

Lastly, I saved the terms to a file called `dictionary.txt`. The dictionary contains three columns: `t_index`, `term`, and the corresponding `df` value.

```
1 dictionary_file = os.path.join('./dictionary.txt') # Path to the
        output file
2 with open(dictionary_file, 'w', encoding='utf-8') as f:
3     # Write header
4     f.write(f"t_index\tterm\tidf\n")
5     # Write each dictionary entry
6     for entry in dictionary_entries:
```

```
7         f.write(f"{entry[0]}\t{entry[1]}\t{entry[2]}\n")
```

4.3 TF-IDF Unit Vector

I created a dictionary based on terms, with each entry containing the term index and its inverse document frequency (IDF) value. The `math.log` function is used to calculate the IDF, which is computed as $\log(n/df)$, where n is the total number of documents, and df is the document frequency of the term. Then, I stored both the IDF value and the index in the dictionary `idf_dict`.

```
1  # Calculate IDF value
2  N = 1095
3  idf_dict = {}
4  for idx, term, df in dictionary_entries:
5      idf_value = math.log10(N / df) if df > 0 else 0
6      idf_dict[term] = (idx, idf_value)
7      # Store the term index and its corresponding IDF value
```

Next, I calculated the term frequency (TF) by counting the occurrences of each term in a document using a `defaultdict` called `term_count`. The TF value is computed as the number of occurrences of the term divided by the total number of terms in the document. I then stored the product of TF and IDF in `tfidf_vector` for each document, indexed by the term index. Following that, I normalized the vectors to unit vectors. Finally, I sorted the dictionary by term index and appended the sorted TF-IDF vector for each document to the list `tfidf_results`.

```
1  # Create a list to hold TF-IDF results for each document
2  tfidf_results = []
3
4  for doc in documents:
5      # Count term frequencies
6      term_count = defaultdict(int)
7      for term in doc:
8          term_count[term] += 1
9      # Total number of terms in the document
10     total_terms = len(doc)
11     tfidf_vector = {}
12     for term, count in term_count.items():
13         tf = count / total_terms # Calculate TF
14         idx, idf_value = idf_dict.get(term, (0, 0))
```

```

15         tfidf_vector[idx] = tf * idf_value # Calculate TF-IDF
16
17     magnitude = np.linalg.norm(list(tfidf_vector.values())) #
        Compute the magnitude
18     if magnitude > 0:
19         tfidf_vector = {idx: value / magnitude for idx, value in
            tfidf_vector.items()} # Normalize values
20
21     # Sort tfidf_vector by idx (the key)
22     sorted_tfidf_vector = dict(sorted(tfidf_vector.items()))
23     tfidf_results.append(sorted_tfidf_vector)

```

Finally, I stored all the TF-IDF values in the output folder, with each document file containing two columns: `t_index` and `tf_idf` values.

```

1 output_folder = os.path.join( "./output/" )
2 os.makedirs(output_folder, exist_ok=True)
3 for i, tfidf in enumerate(tfidf_results):
4     filename = os.path.join(output_folder, f"{i + 1}.txt")
5     with open(filename, 'w', encoding='utf-8') as file:
6         file.write(f"t_index\ttf_idf\n") # Header
7         for idx, tfidf_value in tfidf.items():
8             file.write(f"{idx}\t{tfidf_value}\n") # Write each TF-IDF
            pair

```

4.4 Cosine Similarity

Cosine similarity is calculated by first creating a union of the terms that documents x and y contain. For each document, a vector is formed by checking if the document has a corresponding *tf-idf* value for each term. If a term is absent, the *tf-idf* value is set to zero. Once the vectors are created, the cosine similarity is computed using the formula:

$$\text{Cosine Similarity} = \vec{x} \cdot \vec{y}$$

where $\vec{x} \cdot \vec{y}$ represents the dot product of the two unit vectors.

```

1 def cosine(Docx, Docy):
2     tfidf_x = tfidf_results[Docx - 1]
3     tfidf_y = tfidf_results[Docy - 1]
4
5     # Get the union of keys from both TF-IDF dictionaries

```

```

6     keys = set(tfidf_x.keys()).union(set(tfidf_y.keys()))
7
8     # Create vectors for each document based on the keys
9     vector_x = np.array([tfidf_x.get(key, 0) for key in keys])
10    vector_y = np.array([tfidf_y.get(key, 0) for key in keys])
11
12    # Calculate cosine similarity
13    cosine_similarity = np.dot(vector_x, vector_y)
14    return cosine_similarity
15 similarity = cosine(1, 2) # Calculate similarity between Document 1
    and Document 2
16 print(f"Cosine similarity between Document 1 and Document 2:
    {similarity}")

```

Cosine similarity between Document 1 and Document 2: 0.1672930057087566