

Midterm Solutions Exercise 06

Nguyen Thi Phuong Thao - N21DCCN078

Nguyen Thi Minh Thu - N21DCCN082

Du Trong Nhan - N21DCCN060

Trinh Khanh Quan - N20DCCN057

May 17, 2025

Contents

1	Problem Statement	2
2	Implementation Details	2
2.1	Data Structures	2
2.2	Key Functions	3
2.2.1	CreateIndex(Dir, StopList)	3
2.2.2	Find(Word, Weight, N)	4
2.2.3	FindFromFile(WordFile, N)	5
2.3	Helper Functions	6
3	Installation Guide	6
3.1	Prerequisites	6
3.2	Installation Steps	6
4	Usage Instructions	6
4.1	Running the Program	6
4.2	Using the Program	7
4.3	Word File Format	7
5	Examples	7
5.1	Example 1: Creating an Index	7
5.2	Example 2: Single Word Search	8
5.3	Example 3: Word File Search	8
6	Notes and Limitations	9

1 Problem Statement

Implement a program that creates and uses inverted indexes with the following functionalities:

- (a) **CreateIndex(Dir, StopList)**: Takes a directory name and a file called StopList (in that directory) as input. It returns an inverted index as output. The DocTable includes all files in the directory Dir, except for the StopList file. The TermTable includes only all words occurring in the directory that start with the letter “c” (lower- or uppercase).
- (b) **Find(Word, Weight, N)**: Finds the top N documents in the index associated with the specified word, considering the given weight.
- (c) **FindFromFile(WordFile, N)**: Similar to the above, but instead of taking a single word as input, it takes a file called WordFile. This file has, on each line, a word (string) and a weight (float). It then attempts to find, using the inverted index, the top N matches for this query.

2 Implementation Details

The implementation consists of a Python program (`inverted_index.py`) that creates and searches an inverted index structure for text documents.

2.1 Data Structures

The program uses two main data structures:

1. **DocTable**: A dictionary mapping document IDs (integers) to document names (filenames)

```
1 DocTable = {  
2     0: "document1.txt",  
3     1: "document2.txt",  
4     # ...  
5 }
```

2. **TermTable**: An inverted index mapping terms to dictionaries that map document IDs to term frequencies

```
1 TermTable = {  
2     "computer": {0: 5, 2: 3}, # "computer" appears 5 times in doc  
3     "cloud": {1: 7, 4: 2},    # "cloud" appears 7 times in doc 1,  
4     # ...  
5 }
```

2.2 Key Functions

2.2.1 CreateIndex(Dir, StopList)

This function:

- Takes a directory path and a stop words file name as input
- Loads stop words from the specified file
- Processes each file in the directory (excluding the stop list file)
- Extracts words from each document
- Filters to include only words starting with 'c' (case insensitive)
- Removes stop words
- Counts word frequencies
- Builds the DocTable and TermTable
- Returns both tables

```
1 def CreateIndex(Dir: str, StopList: str) -> Tuple[Dict[int, str], Dict[
2     str, Dict[int, int]]:
3     # Initialize data structures
4     DocTable = {} # Maps document IDs to filenames
5     TermTable = defaultdict(lambda: defaultdict(int)) # Maps terms to
6     {doc_id: frequency}
7
8     # Load stop words
9     stop_words = set()
10    stop_list_path = os.path.join(Dir, StopList)
11    if os.path.exists(stop_list_path):
12        with open(stop_list_path, 'r', encoding='utf-8', errors='ignore')
13        as f:
14            stop_words = {word.strip().lower() for word in f}
15
16    # List of files to process (excluding the StopList file)
17    files_to_process = [f for f in os.listdir(Dir)
18                        if os.path.isfile(os.path.join(Dir, f)) and f
19                        != StopList]
20
21    # Process each file
22    for doc_id, filename in enumerate(files_to_process):
23        file_path = os.path.join(Dir, filename)
24
25        # Add file to DocTable
26        DocTable[doc_id] = filename
27
28        try:
29            # Process file contents
30            with open(file_path, 'r', encoding='utf-8', errors='ignore')
31            as f:
32                content = f.read().lower()
33
34            # Extract words (simple tokenization)
```

```

30         words = re.findall(r'\b\w+\b', content)
31
32         # Count words that start with 'c' and not in stop words
33         word_counts = Counter(word for word in words
34                                if word.startswith('c') and word not in
stop_words)
35
36         # Update TermTable with word frequencies
37         for word, count in word_counts.items():
38             TermTable[word][doc_id] = count
39     except Exception as e:
40         print(f"Error processing file {filename}: {e}")
41
42     # Convert nested defaultdicts to regular dicts for return
43     return DocTable, {term: dict(doc_freqs) for term, doc_freqs in
TermTable.items()}

```

2.2.2 Find(Word, Weight, N)

This function:

- Searches for a specific word in the TermTable
- Calculates document scores based on term frequency * weight
- Returns the top N documents sorted by score

```

1 def Find(Word: str, Weight: float, N: int, DocTable: Dict[int, str],
2          TermTable: Dict[str, Dict[int, int]]) -> List[Tuple[str, float
3          ]]:
4     word = Word.lower()
5     results = []
6
7     # Check if word exists in the index
8     if word not in TermTable:
9         print(f"Word '{word}' not found in the index.")
10        return []
11
12    # Get document IDs and their frequencies for the word
13    doc_freqs = TermTable[word]
14
15    # Calculate scores based on term frequency and term weight
16    for doc_id, freq in doc_freqs.items():
17        doc_name = DocTable[doc_id]
18        # Score is the product of term frequency and the term weight
19        score = freq * Weight
20        results.append((doc_name, score))
21
22    # Sort by score in descending order
23    results.sort(key=lambda x: x[1], reverse=True)
24
25    # Return top N results
26    return results[:N]

```

2.2.3 FindFromFile(WordFile, N)

This function:

- Reads words and their weights from a file
- Accumulates scores for documents across all query terms
- Returns the top N documents sorted by total score

```
1 def FindFromFile(WordFile: str, N: int, DocTable: Dict[int, str],
2                 TermTable: Dict[str, Dict[int, int]]) -> List[Tuple[str
3                 , float]]:
4     # Initialize a dictionary to accumulate scores for each document
5     document_scores = defaultdict(float)
6
7     try:
8         # Read words and weights from the file
9         with open(WordFile, 'r', encoding='utf-8', errors='ignore') as
10            f:
11                for line in f:
12                    line = line.strip()
13                    if not line:
14                        continue
15
16                # Parse word and weight from the line
17                parts = line.split()
18                if len(parts) >= 2:
19                    word = parts[0].lower()
20                    try:
21                        weight = float(parts[1])
22                    except ValueError:
23                        print(f"Invalid weight for word '{word}', using
24                        default weight 1.0")
25                        weight = 1.0
26
27                # Check if word exists in the index
28                if word not in TermTable:
29                    print(f"Word '{word}' not found in the index,
30                    skipping.")
31                    continue
32
33                # Get document IDs and frequencies for the word
34                doc_freqs = TermTable[word]
35
36                # Calculate and accumulate scores
37                for doc_id, freq in doc_freqs.items():
38                    doc_name = DocTable[doc_id]
39                    # Score is product of term frequency and query
40                    term weight
41
42                    score = freq * weight
43                    document_scores[doc_name] += score
44
45            else:
46                print(f"Ignoring invalid line: {line}")
47
48    # Convert scores to list of tuples
```

```

43     results = [(doc_name, score) for doc_name, score in
44         document_scores.items()]
45
46     # Sort by score in descending order
47     results.sort(key=lambda x: x[1], reverse=True)
48
49     # Return top N results
50     return results[:N]
51
52 except Exception as e:
53     print(f"Error processing word file '{WordFile}': {e}")
54     return []

```

2.3 Helper Functions

The program includes several helper functions:

- **display_index()**: Displays the contents of DocTable and TermTable
- **save_to_csv()**: Saves DocTable and TermTable to CSV files
- **main()**: The main program loop that handles user interaction

3 Installation Guide

3.1 Prerequisites

- Python 3.6 or higher

3.2 Installation Steps

1. Ensure Python is installed on your system:

```
1 python --version
```

2. Download the `inverted_index.py` file to your preferred location

4 Usage Instructions

4.1 Running the Program

1. Open a terminal or command prompt
2. Navigate to the directory containing the `inverted_index.py` file
3. Run the program:

```
1 python inverted_index.py
```

4.2 Using the Program

After running the program, follow these steps:

1. Specify Document Directory:

- When prompted, enter the full path to the directory containing your text documents

2. Specify StopList File:

- Enter the name of the file containing stop words (words to be excluded from indexing)

3. View the Created Index:

- The program will display the created DocTable and TermTable
- The results will also be saved to CSV files (`doc_table.csv` and `term_table.csv`)

4. Search Options:

- Option 1: Search by a single word
 - Enter the word, its weight, and the number of results to retrieve
- Option 2: Search using a word file
 - Enter the path to a file containing words and weights
 - Enter the number of results to retrieve
- Option 3: Exit the program

4.3 Word File Format

The word file should have the following format:

```
1 word1 weight1
2 word2 weight2
3 ...
```

Example:

```
1 computer 1.5
2 coding 2.0
3 cloud 0.8
```

5 Examples

5.1 Example 1: Creating an Index

Input:

- Directory: `/documents`
- StopList: `stopwords.txt`

Contents of Directory:

- `article1.txt`: "Computer science is a fascinating field. Computing power has increased."
- `article2.txt`: "Cloud computing is revolutionizing IT. Companies are moving to the cloud."
- `stopwords.txt`: "is a the to are"

Output:

```

1 Document Table:
2 -----
3 ID: 0 -> article1.txt
4 ID: 1 -> article2.txt
5
6 Term Table:
7 -----
8 'computer' appears in: [0] (article1.txt (freq: 1))
9 'computing' appears in: [0, 1] (article1.txt (freq: 1), article2.txt (
   freq: 1))
10 'cloud' appears in: [1] (article2.txt (freq: 2))
11 'companies' appears in: [1] (article2.txt (freq: 1))

```

5.2 Example 2: Single Word Search

Input:

- Word: "computing"
- Weight: 1.5
- N: 2

Output:

```

1 Top 2 documents for 'computing' (weight: 1.5):
2 -----
3 Document: article1.txt, Score: 1.5
4 Document: article2.txt, Score: 1.5

```

5.3 Example 3: Word File Search

Input:

- WordFile: `query.txt` containing:

```

1 computing 1.5
2 cloud 2.0

```

- N: 2

Output:

```

1 Top 2 documents for query in 'query.txt':
2 -----
3 Document: article2.txt, Score: 5.5
4 Document: article1.txt, Score: 1.5

```


6 Notes and Limitations

1. The current implementation only indexes words that start with the letter 'c' (case insensitive)
2. The program uses a simple tokenization approach (words separated by spaces, punctuation removed)
3. All words are converted to lowercase for indexing and searching