Project 3: Demo

Christopher Di Giacomo

40133600

10 November 2023

COMP479 2232 D

Dr. Sabine Bergler

# Table of Contents

# Environment/Setup

Directories:

- **/reports:** This directory contains the report and demo for this project
- **/src:** This directory contains the *subProject1.py* and *subProject2.py* files
- **/indexes:** This directory contains the indexes in JSON format. It is (re)generated at runtime for either *.py* files.
- **/Corpus:** This directory contains all of the reuters collection *.sgm* files

Notes:

- The current working directory must be the parent directory of the submission folder. Although the /reports subdirectory will be present in the project folder structure, within the text editor of choice, this will not interfere with any computation.
- If using an extension based IDE like VS Code, you must run the pipeline.py file as a Python file, and not via any code runner extensions.

# Subproject1: First-Run Walkthrough

Running *subProject1.py* will trigger the *main()* function call and result in the following sequence of operations until terminal output.

## Main

The main function below executes, retrieving the collection, building both indexes, and calling *getStats()* to display the performance data.

```python
def main():
    """
    Main function to execute the token retrieval, index building, and statistics generation process.
    It retrieves tokens, builds indexes using both naive and SPIMI approaches, and then prints statistics.
    """

    # Measure time taken to create token streams
    start = time.time()
    tokenStream = getTokens()
    end = time.time()
    tokensTime = round(end - start, 3)

    # Build SPIMI and Naive indexes
    spimiIndex, spimiTime = spimi(tokenStream, 10000)
    naiveIndex, naiveTime = naive(tokenStream, 10000)

    # Generate and print statistics
    getStats(naiveIndex, naiveTime, spimiIndex, spimiTime, tokensTime)
```

## getTokens()

The first function called in the main is the *getTokens()* function, displayed below. As in the previous report, it iterates through all .sgm files of the collection and calls the tokenizer() function on each document. The resulting tokenized collection is returned, in variable tokensList.

```python
def getTokens():
    """
    Retrieves and tokenizes documents from a specified corpus. Each .sgm file in the
    corpus is parsed, and documents within are tokenized. Tokens are associated with their
    document ID.

    Returns:
        list: A list where each element represents the tokens of a document, indexed by document ID.
    """

    tokensList = []
    # Inform about the start of token retrieval process
    print(f'Retrieving tokens ...')
    # Loop through each file in the corpus
    for file_number in range(22):
        with open(f'./Corpus/reut2-{str(file_number).zfill(3)}.sgm', 'r', encoding='windows-1252') as f:
            # Parse the file content using BeautifulSoup
            soup = BeautifulSoup(f, 'html.parser')
            documents = soup.find_all('text')

            for document in documents:
                doc_type = document.get('type')
                text_parts = []

                # Different handling for 'BRIEF' type documents
                if doc_type == 'BRIEF':
                    title = document.title.get_text() if document.title else ""
                    text_parts.append(title)
                else:
                    title = document.title.get_text() if document.title else ""
                    body = document.body.get_text() if document.body else ""
                    text_parts.append(title)
                    text_parts.append(body)

                # Combine title and body for tokenization
                full_text = ' '.join(text_parts)
                tokens = tokenizer(full_text)
                tokensList.append(tokens)

    # Indicate completion of token retrieval
    print(f'Tokens retrieved!')
    return tokensList
```

## tokenizer()

The tokenizer() function is displayed below. As in project 2, it returns the tokenized inputted string with respect to a regex pattern, using NLTK's regexp_tokenize().

```python
4   def tokenizer(text):
5       """
6       This function takes in a text and returns a list of tokens. It uses a regular expression to
7       handle acronyms, abbreviations, and general word patterns including words with hyphens or apostrophes.
8
9       Args:
10          text (str): The text to be tokenized.
11
12      Returns:
13          list: A list of tokens extracted from the input text.
14      """
15
16      # Pattern captures abbreviations/acronyms (e.g., U.S.A.) and general word forms.
17      pattern = r'\b(?:[A-Z]{1,2}\.)*[A-Z]{1,2}\.?\b|\b\w+(?:[-\']\w+)*\b'
18      return nltk.regexp_tokenize(text, pattern)
```

As mentioned, once the getTokens() function has tokenized all .sgm files using tokenizer(), it returns a list of tokenized documents which are in the form of lists of tokens. This is returned to the main method, which has calculated the elapsed time for this phase of subproject 1. Next, it begins by constructing the spimiIndex, whose implementation can be found below.

## spimi()

The function builds the inverted index in a single pass using the tokensList, returned by getTokens(). It also has an optional parameter for subproject 1, which halts the construction

```python
69   def spimi(tokensList, testCount=None):
70       """
71       Builds an inverted index using the Single-Pass In-Memory Indexing (SPIMI) algorithm.
72       It processes tokenized documents, generating a dictionary where each term is associated
73       with a list of tuples containing document IDs and term frequency.
74
75       Args:
76           tokensList (list): A list of lists, each containing tokens from a document.
77           testCount (int, optional): A threshold for the number of term-document pairs to process.
78
79       Returns:
80           tuple: A tuple containing the inverted index and the time taken to build it.
81       """
82       index = {}
83       docID = 1
84       pairCount = 0
85       # Print collection size and start building index
86       print(f'Collection size: {len(tokensList)}')
87       print(f'Building SPIMI index ...')
88       startTime = time.time()
89       for tokens in tokensList:
90           for token in tokens:
91               pairCount += 1
92               # Update or add token in the index
93               if token in index:
94                   postingsList = index[token]
95                   if postingsList[-1][0] == docID:
96                       postingsList[-1] = (docID, postingsList[-1][1] + 1)
97                   else:
98                       postingsList.append((docID, 1))
99               else:
100                  index[token] = [(docID, 1)]
101
102              if testCount and pairCount >= testCount:
103                  break
104          docID += 1
105          if testCount and pairCount >= testCount:
106              break
107      # Sort index by keys (terms) before saving
108      index = dict(sorted(index.items(), key=lambda item: item[0]))
109      endTime = time.time()
110
111      print(f'SPIMI index built!')
112      save2json(index, 'spimi.json')
113      return index, endTime - startTime
```

after the first *testCount* terms have been processed. In main() we set testCount to 10 000. The function builds the index in one pass over the data within tokensList. It adds terms to the index, along with a postings list of the docIDs it occurs in. The postings are in the form of docID-term frequency tuples, and everytime a term reoccurs in the same document, the term frequency for that document is incremented. Throughout the construction process the function is being timed, and when it returns the final index is also returns the elapsed time.

## save2json()

Just before returning, the function calls save2json() which is seen below. It simply outputs the dictionary structure to a JSON file object within the /indexes directory.

```python
def save2json(data, filename):
    """
    Saves a given data object to a JSON file.

    Args:
        data: The data to be saved.
        filename (str): The name of the file to save the data in.
    """

    # Create directory if it doesn't exist
    if not os.path.exists('indexes'):
        os.makedirs('indexes')

    # Save data to a JSON file
    with open(f'./indexes/{filename}', 'w') as f:
        json.dump(data, f, indent=4)
```

## naïve()

Back in main(), the spimi index and its construction time is retrieved is retrieved, and naïve() is now called, which can be seen below. It has the same parameters as spimi(), which work the same way. However, this function stores term-docID pairs in a list F, which it then sorts and

```python
132    def naive(tokensList, testCount = None):
133        """
134        Builds an inverted index using a naive approach. It processes term-document pairs,
135        sorts them, and then groups them to create postings lists.
136
137        Args:
138            tokensList (list): A list of tokenized documents.
139            testCount (int, optional): A threshold for the number of term-document pairs to process.
140
141        Returns:
142            tuple: A tuple containing the inverted index and the time taken to build it.
143        """
144
145        F = []  # List to store term-docID pairs
146        docID = 1
147        print(f'Building Naive index ...')
148
149        startTime = time.time()
150        # Create term-docID pairs
151        for tokens in tokensList:
152            for token in tokens:
153                F.append((token, docID))
154                if testCount and len(F) == testCount:
155                    break
156            docID += 1
157            if testCount and len(F) == testCount:
158                break
159
160        # Sort and remove exact duplicates
161        F.sort()
162        F = list(dict.fromkeys(F))
163
164        # Create postings lists
165        index = {}
166        for term, docID in F:
167            if term in index:
168                if index[term][-1] != docID:
169                    index[term].append(docID)
170            else:
171                index[term] = [docID]
172
173        endTime = time.time()
174        print(f'Naive index built!')
175        save2json(index, 'naive.json')
176        return index, endTime - startTime
```

cleans for duplicates, all before actually constructing the index. When, building the index, it works exactly as does spimi(), although it does not maintain term frequency. It also writes the index to a JSON file, and returns both the index and the elapsed time.

Back in main(), both indexes are now built, and both of their elapsed times have been retrieved.

## getStats()
The next function call is to getStats() which can be seen below. This function simply takes all of the previously gathered data as input parameters, and displays it in a proper manner via

```
178  def getStats(naiveIndex, naiveTime, spimiIndex, spimiTime, tokensTime):
179      """
180      Prints statistics comparing the Naive and SPIMI indexing methods. It displays the time taken
181      to create token streams, build indexes, and the sizes of the created indexes.
182
183      Args:
184          naiveIndex (dict): The inverted index created using the naive approach.
185          naiveTime (float): The time taken to build the naive index.
186          spimiIndex (dict): The inverted index created using the SPIMI approach.
187          spimiTime (float): The time taken to build the SPIMI index.
188          tokensTime (float): The time taken to create token streams.
189      """
190      print('\n========= STATISTICS =========\n')
191      # Token streams creation time
192      print(f'Token streams creation time (sec): {tokensTime}')
193      # SPIMI Index Construction Statistics
194      print('\n--------------------\n| SPIMI Construction |\n--------------------')
195      print(f'\nTime taken for SPIMI index to process 10 000 terms (ms): {round(spimiTime * 1000, 3)}')
196      print(f'Size of SPIMI index: {len(spimiIndex)} terms')
197      # Naive Index Construction Statistics
198      print('\n--------------------\n| Naive Construction |\n--------------------')
199      print(f'\nTime taken for Naive index to process 10 000 terms (ms): {round(naiveTime * 1000, 3)}')
200      print(f'Size of Naive index: {len(naiveIndex)} terms')
201      # Differences between SPIMI and Naive Methods
202      print('\n---------------\n| Differences |\n---------------')
203      time_diff = naiveTime - spimiTime
204      time_diff_percent = (time_diff / naiveTime) * 100
205      print(f'\nTime difference (ms): {round(time_diff * 1000, 3)}')
206      print(f'Time difference (%): {round(time_diff_percent, 3)}%')
```

terminal output. It displays the elapsed time for both index construction techniques, as well as the time difference in ms and %, as follows.

```
Retrieving tokens ...
Tokens retrieved!
Collection size: 21578
Building SPIMI index ...
SPIMI index built!
Building Naive index ...
Naive index built!

========= STATISTICS =========

Token streams creation time (sec): 10.812

--------------------
| SPIMI Construction |
--------------------

Time taken for SPIMI index to process 10 000 terms (ms): 4.77
Size of SPIMI index: 2769 terms

--------------------
| Naive Construction |
--------------------

Time taken for Naive index to process 10 000 terms (ms): 6.836
Size of Naive index: 2769 terms

---------------
| Differences |
---------------

Time difference (ms): 2.066
Time difference (%): 30.224%
```

# Subproject 2: First-Run Walkthrough

Running *subProject2.py* will trigger the *main()* function call and result in the following sequence of operations until terminal output.

## Main

The main function below executes, retrieving the collection, building both indexes, and calling the queryManager() to handle continuous prompts for queries. The function begins by calling the same three initial functions as in the previous subproject. It calls getTokens(), as well as spimi() and naïve(), however it does not set thresholds for the indexes, and neither does it bother retrieving their construction times. Since these have all been seen in subproject1, we will skip to the last function call to queryManager()

```
397    def main():
398        """
399        The main function to initiate the search engine application. It retrieves token streams, builds indexes,
400        and launches the query manager to handle user queries.
401
402        This function is the starting point of the application.
403        """
404        tokenStreams = getTokens()
405        spimiIndex, _ = spimi(tokenStreams)
406        naiveIndex, _ = naive(tokenStreams)
407        queryManager(spimiIndex, naiveIndex, tokenStreams)
```

## queryManager()

One of the main functions which handles all functionality to this subproject is the queryManager(), which is seen bellow. This function is responsible for continuously prompting the user for input parameters to the auxiliary functions necessary to the overall functioning of subproject2. It continuously prompts the user for a query, until 'q' is entered, which terminates the engine. For a single term query, the function would skip the prompt for operations, but in the case of a multi-term query the first function call would be to getOperation().

```
371    def queryManager(spimiIndex, naiveIndex, collection):
372        """
373        Manages the querying process. It prompts the user to enter queries and handles the retrieval
374        and display of results using the chosen index and ranking method.
375
376        Args:
377            spimiIndex (dict): The SPIMI index.
378            naiveIndex (dict): The Naive index.
379            collection (list): The document collection, used for BM25 ranking.
380
381        This function continues to prompt for queries until the user chooses to quit.
382        """
383        print('\n====== WELCOME TO THE REUTERS21578 SEARCH ENGINE ======\n')
384        print('\n** Enter \'q\' to quit **\n')
385        print('** Queries do not require boolean operators **\n')
386        while True:
387            query = input('Enter your query: ')
388            if query == 'q':
389                break
390            numKeywords = len(tokenizer(query))
391            operation = getOperation() if numKeywords > 1 else None
392            spimi, naive = getIndexes(spimiIndex, naiveIndex)
393            queryTermRanking, bm25Ranking = getRanking(operation, spimi)
394            queryTest(query=query, spimiIndex=spimi, naiveIndex=naive, collection=collection, operation=operation, queryTermRanking=queryTermRanking, bm25Ranking=bm25Ranking)
395        print('\n====== SEE YOU LATER :P ======\n')
```

## getOperation()

This function, as seen below, is simply responsible for continuously prompting the user for a Boolean operation (AND, OR) until a valid input is received.

```
294   def getOperation():
295       """
296       Prompts the user to choose a boolean operation ('AND' or 'OR') for the query processing.
297
298       Returns:
299           str: The chosen boolean operation ('AND' or 'OR').
300       """
301       valid = False
302       while not valid:
303           operation = input('Enter your operation (\'AND\', \'OR\') : ')
304           # Check if the entered operation is valid
305           if operation in ['AND', 'OR']:
306               valid = True
307           else:
308               print(f'Sorry, {operation} is not a valid operation.')
309       return operation
```

Once back in queryManager() with an operation retrieved, the next function call is to getIndexes() which can be seen below. This function works the same way in principle, prompting the user until a valid input is entered. The user could select to query either index, or both. The function will return the required, pre-built indexes depending on the selection.

```
311   def getIndexes(spimiIndex, naiveIndex):
312       """
313       Prompts the user to choose which index(es) to query – SPIMI, Naive, or both.
314
315       Args:
316           spimiIndex (dict): The SPIMI index.
317           naiveIndex (dict): The Naive index.
318
319       Returns:
320           tuple: A tuple of the selected indexes. None is used for unselected indexes.
321       """
322       valid = False
323       while not valid:
324           answer = input('Which index would you like to query? \'n\' = naive, \'s\' = spimi, \'b\' = both : ')
325           # Validate the user input and return the appropriate indexes
326           if answer in ['n', 's', 'b']:
327               valid = True
328           else:
329               print(f'Sorry, {answer} is not a valid answer.')
330       # Return the appropriate combination of indexes based on user choice
331       if answer == 's':
332           return spimiIndex, None
333       elif answer == 'n':
334           return None, naiveIndex
335       elif answer == 'b':
336           return spimiIndex, naiveIndex
```

Once back in queryManager(), with the selected indexes, the next function call is to getRanking(), which principally works the same as the previous two functions. It prompts the user until it receives a valid ranking selection. The user can select query term ranking if they previously selected their operation as *OR*. They can select bm25 ranking if at least one of their selected indexes is the spimi index, or they can input *None* for no ranking. The function returns Boolean flags for each ranking technique.

```
338    def getRanking(operation, spimiIndex):
339        """
340        Prompts the user to choose a ranking method for the query results.
341
342        Args:
343            operation (str): The chosen boolean operation ('AND' or 'OR').
344            spimiIndex (dict): The SPIMI index.
345
346        Returns:
347            tuple: A tuple (queryTermRank, bm25) indicating whether query term ranking or BM25 should be used.
348        """
349        queryTermRank = False
350        bm25 = False
351        valid = False
352        while not valid:
353            answer = input('Which ranking would you like? \'q\' = query term ranking, \'b\' = bm25, \'n\' = none : ')
354            # Validate the user input and set the ranking method accordingly
355            if answer in ['q', 'b', 'n']:
356                if answer == 'q' and operation != 'OR':
357                    print('Sorry, query term ranking is only for multi-term OR queries.')
358                elif answer == 'b' and not spimiIndex:
359                    print('Sorry, bm25 ranking is only available for the spimi index.')
360                else:
361                    valid = True
362                    if answer == 'q':
363                        queryTermRank = True
364                    elif answer == 'b':
365                        bm25 = True
366            else:
367                print(f'Sorry, {answer} is not a valid answer.')
368
369        return queryTermRank, bm25
```

## queryTest()

Once back in queryManager(), now with the ranking flags set, queryManager() makes its last function call of the current iteration to queryTest(). The queryTest() function as shown below, is the function responsible for processing the queries depending on the inputted runtime parameters, which is why it is the longest. However it works rather simply. It begins by calling the tokenizer() on the query to retrieve a list of the tokenized query terms. If the list is empty then that means the inputted query was invalid. If it has a length of 1 then that means it is a single term query, and if more than 1 it is a multi-term query. The function begins by verifying if the naïve index has been initialized, if so, it steps into the if block and verifies the nature of the query. If it's a single term query, it queries the naïve index for the query term and outputs the

```
8    def queryTest(query, spimiIndex=None, naiveIndex=None, collection=None, operation='OR', queryTermRanking=False, bm25Ranking=False):
9        """
10       Processes a search query against provided SPIMI and/or Naive indexes. It supports single or multi-term
11       queries with AND/OR operations and optional ranking (Query Term Ranking or BM25).
12
13       Args:
14           query (str): The search query.
15           spimiIndex (dict): The SPIMI index (optional).
16           naiveIndex (dict): The Naive index (optional).
17           collection (list): The document collection (used for BM25 ranking).
18           operation (str): The boolean operation to apply ('AND', 'OR').
19           queryTermRanking (bool): Flag to use Query Term Ranking.
20           bm25Ranking (bool): Flag to use BM25 ranking.
21
22       The function tokenizes the query, checks the appropriate index for each term, and then applies the
23       specified operation (AND/OR) to compute the final result set. It also handles ranking if specified.
24       """
25       # Tokenize the input query
26       queryTerms = tokenizer(query)
27
28       # Check if the query is empty after tokenization
29       if len(queryTerms) == 0:
30           print(f'Sorry, \'{query}\' is not a valid query')
31           return
32
33       # Process with Naive Index if available
34       if naiveIndex:
35           print("\n====== NAIVE INDEX ======\n")
36           # Process single-term queries
37           if len(queryTerms) == 1:
38               result = naiveIndex.get(queryTerms[0])  # Get postings for the single query term.
39           else:
40               # Process multi-term queries
41               postingsLists = []
42               result = None
43               for term in queryTerms:
44                   postingsList = naiveIndex.get(term)  # Get postings for each query term.
45                   postingsLists.append(postingsList)
46               # Apply boolean operations
47               if operation == 'AND':
48                   result = conjunction(postingsLists)
49               elif operation == 'OR':
50                   result = disjunction(postingsLists, queryTermRanking)
51
52           # Print the results from Naive Index
53           if result:
54               print(f"QUERY: \'{query}\'\nOPERATION: {operation}\nRANKING: {'Query Term Ranking' if queryTermRanking else 'None'} \nPOSTINGS LIST SIZE: {len(result)}\nPOSTINGS LIST: {str(result)}\n")
55           else:
56               print(f"QUERY: \'{query}\' - Your query returned no results with operation: {operation}.\n")
```

result. In the case of a multi-term query, the function iterates through the query terms, and retrieves their postings lists.

## conjunction()

It then calls either conjunction() or disjunction() depending on the operation chosen at runtime. If AND, then conjunction is called.

```python
189    def conjunction(postingsLists):
190        """
191        Performs an AND operation on a list of postings lists, returning the intersection of these lists.
192
193        Args:
194            postingsLists (list of lists): A list where each element is a postings list (list of docIDs).
195
196        Returns:
197            list: The intersection of the postings lists, representing documents that contain all terms.
198        """
199        # Return None if any of the postings lists is None
200        if None in postingsLists:
201            return None
202
203        # Sort postings lists by length for efficient intersection
204        sortedPostingsLists = sorted(postingsLists, key=len)
205        # Start with the shortest postings list
206        finalPostingsList = sortedPostingsLists[0]
207
208        # Intersect with each subsequent postings list
209        for postingsList in sortedPostingsLists[1:]:
210            finalPostingsList = intersect(finalPostingsList, postingsList)
211            # If intersection is empty, no further processing is needed
212            if not finalPostingsList:
213                break
214        return finalPostingsList
```

## Intersect()

The conjunction function calls the intersect() function on each of the postings lists that have been passed as a parameter beginning with the two smallest. It does so two at a time, and combines the results in order. The intersection function can be seen below, which functions as a regular intersection, combining postings lists by common docIDs while avoiding duplicates and incrementing with respect to the list positioned at the smallest docID, until the end of one of the lists has been reached.

```python
248    def intersect(pList1, pList2):
249        """
250        Computes the intersection of two postings lists.
251
252        Args:
253            pList1 (list): The first postings list.
254            pList2 (list): The second postings list.
255
256        Returns:
257            list: The intersection of pList1 and pList2.
258        """
259        result = []
260        i = j = 0
261        # Iterate through both lists to find common elements
262        while i < len(pList1) and j < len(pList2):
263            if pList1[i] == pList2[j]:
264                result.append(pList1[i])
265                i += 1
266                j += 1
267            elif pList1[i] < pList2[j]:
268                i += 1
269            else:
270                j += 1
271        return result
```

Back in the queryTest() function, in the case of a conjunction, the resulting list of docIDs is returned and displayed.

## disjunction()

However, if the selected operation is *OR* then the disjunction() function is called which can be seen below. It simply unions the postings lists in order and without duplicates, unless they are all *None*.

```python
216    def disjunction(postingsLists, queryTermRanking=False):
217        """
218        Performs an OR operation on a list of postings lists, returning the union of these lists. Can also
219        apply query term ranking based on term frequency across the postings lists.
220
221        Args:
222            postingsLists (list of lists): A list where each element is a postings list (list of docIDs).
223            queryTermRanking (bool): Flag to apply query term ranking.
224
225        Returns:
226            list: The union of the postings lists, with optional query term ranking applied.
227        """
228        # Filter out None postings lists (terms not found)
229        validPostingsLists = [plist for plist in postingsLists if plist is not None]
230
231        # If all postings lists are None, return None
232        if not validPostingsLists:
233            return None
234
235        # Flatten the list of lists into a single list using itertools.chain
236        unionPostings = list(chain(*validPostingsLists))
237
238        # If query term ranking is not applied, remove duplicates and sort
239        if not queryTermRanking:
240            finalPostings = sorted(set(unionPostings))
241        else:
242            # Sort by document ID and apply query term ranking
243            finalPostings = sorted(unionPostings)
244            finalPostings = queryTermRank(finalPostings)
245
246        return finalPostings
```

## queryTermRank()

In the case of queryTermRanking flag is set to true, then this is handled in the disjunction() function by calling the queryTermRank() function on the resulting disjunct set. The queryTermRank() can be seen below, where it maintains a local dictionary of docIDs, incrementing each of their counts everytime they are encountered in the postingsList. It then sorts the dictionary with respect to value, and returns the dictionary as a list of ordered docID-count tuples.

```python
273    def queryTermRank(postingsList):
274        """
275        Applies query term ranking to a postings list. Ranks documents based on the number of query terms they contain.
276
277        Args:
278            postingsList (list): A postings list (list of docIDs).
279
280        Returns:
281            list: A list of tuples (docID, score) where score is the count of query terms in the document.
282        """
283        # Count the frequency of each document ID in the postings list
284        docFreq = {}
285        for docID in postingsList:
286            docFreq[docID] = docFreq.get(docID, 0) + 1
287
288        # Sort documents by frequency (score) in descending order
289        rankedDocs = sorted(docFreq.items(), key=lambda item: item[1], reverse=True)
290        return [(docID, score) for docID, score in rankedDocs]
```

Back in queryTest(), now with the resulting disjunct set of docIDs in the case of an *OR* operation, the results are displayed to terminal. The function then breaks from the *if naïve* if block, and checks if the spimi index parameter was initialized. If yes, then it enters the if block and goes through the same steps as above for the naïve if block. **The only difference**, is two main points.

## convertPostingsLists()

Firstly, since the spimi index's postings lists are lists of tuples and not regular docIDs, it calls a function to normalize the postings lists to a simple list of docIDs. The function can be seen below. It does so before calling either conjunction() or disjunction().

```python
175    def convertPostingsLists(postingsLists):
176        """
177        Convert a list of postings lists with (docID, tf) tuples to a list of lists of docIDs.
178        If a postings list is None, it remains None.
179
180        Args:
181            postingsLists (list): A list of postings lists, where each postings list contains tuples of (docID, tf).
182
183        Returns:
184            list: A list of lists, where each inner list contains docIDs or is None.
185        """
186        # Convert each postings list to a list of docIDs, or keep as None if postings list is None
187        return [[docID for docID, tf in postingsList] if postingsList is not None else None for postingsList in postingsLists]
```

## bm25()

Secondly, the only other difference is that it checks if the bm25 flag has been set. If yes, it makes a function call to the bm25() function below. It maintains the scores of documents in a local dictionary rankedResults. It does so by retrieving and iterating over the docIDs in each of the queryTerms' postings lists. When the docID matches one of the docIDs in the result parameter, then it computes the score, and stores/updates it in the rankedResults dictionary. Once all queryTerms have been iterated over the dictionary is ordered by score and returned as a sorted list of docID-score tuples.

```python
126    def bm25(queryTerms, result, index, collection, k1, b):
127        """
128        Applies the BM25 ranking formula to the result set of a query. It calculates scores for documents
129        based on term frequency, document frequency, and document length.
130
131        Args:
132            queryTerms (list): List of tokenized query terms.
133            result (list): List of document IDs obtained from the initial search.
134            index (dict): The search index.
135            collection (list): The document collection, to compute document lengths.
136            k1 (float): BM25 parameter k1.
137            b (float): BM25 parameter b.
138
139        Returns:
140            list: A sorted list of tuples (docID, score), ranked according to BM25.
141        """
142        N = len(collection)  # Total number of documents
143        L_total = sum(len(doc) for doc in collection)  # Total length of all documents
144        L_avg = L_total / N  # Average document length
145
146        rankedResults = {}
147
148        if result is None:
149            return None
150
151        # Calculate BM25 score for each term in each document
152        for term in queryTerms:
153            postingsList = index.get(term, [])
154            df = len(postingsList)  # Document frequency
155            idf = math.log(N / df) if df != 0 else 0  # Inverse document frequency
156
157            for docID, tf in postingsList:
158                if docID and docID in result:
159                    L_d = len(collection[int(docID)])  # Document length for the current docID
160                    # BM25 formula
161                    score = idf * ((tf * (k1 + 1)) / (k1 * ((1 - b) + b * (L_d / L_avg)) + tf))
162                    # Accumulate score for each document
163                    if docID not in rankedResults:
164                        rankedResults[docID] = 0.0
165                    rankedResults[docID] += score
166
167        # Check if all scores are zero
168        if all(score == 0.0 for score in rankedResults.values()):
169            return None
170
171        # Sort the results by score in descending order
172        sortedRankedList = sorted(rankedResults.items(), key=lambda item: item[1], reverse=True)
173        return sortedRankedList
```

Once this is complete, the queryManager() then prompts the user for another query, restarting the whole process.

*** For output logs, please refer to Appendix A and Appendix B of 4013360_report.pdf