

Inverted Index and B-tree-based Information Retrieval System

Andrea Frigatti

June 19, 2024

1 Introduction

This report details the implementation of an information retrieval system that utilizes an inverted index and B-tree structures to efficiently handle various types of search queries. The implementation is spread across three Python scripts, each contributing to different stages of the process: constructing the inverted index, building and managing B-trees, and executing search queries.

2 Implementation Details

2.1 Inverted Index Construction

The first script, `inverted_index.py`, is responsible for reading a dataset of movie plots, cleaning the text, and constructing an inverted index. Key steps in this script include:

- Reading and processing the data.
- Removing punctuation, normalizing text, deleting a list of stopwords.
- Constructing the inverted index and saving it along with the document data.

```
1 def load_data(file_path , stopwords):
2     df = pd.read_csv(file_path)
3     plots = df["Plot"]
4
5     for id, plot in plots.items():
6         plot = plot.translate(str.maketrans('', '', string.punctuation))
7         plot = re.sub(r'\W+', ' ', plot)
8
9         words = set(plot.lower().split())
10        for word in words:
11            if word not in stopwords:
12                posting_lists.setdefault(word, []).append(id)
13
14        document_data[id] = df.loc[id].to_dict()
15
16    with open(save_file_path , "w" , encoding='utf-8') as file:
17        for word in posting_lists:
18            file.write(f"{word}: {posting_lists[word]}\n")
19
20    save_document_data_pickle(document_data , document_pickle_file_path)
```

Listing 1: Inverted Index Construction

Advantages:

- **Efficient Text Processing:** The script effectively removes punctuation and normalizes the text, ensuring uniform and reliable indexing.
- **Use of Stopwords:** Excluding common stopwords keeps the inverted index concise and focused on meaningful terms.
- **Data Persistence:** Saving the inverted index and document data to disk allows for quick reloads and use in subsequent processes.

2.2 B-trees Construction and Management

The second script, `btree.py`, defines a B-tree structure and builds B-trees from the inverted index file. Key steps include:

- Defining the `Node` and `BTree` classes.
- Building B-trees from the inverted index file (Search Tree and Reversed Search Tree).
- Saving the B-trees for future use in `retrieving.py`.

```
1 class Node:
2     def __init__(self, is_leaf=False):
3         self.is_leaf = is_leaf
4         self.keys = []
5         self.values = []
6         self.children = []
7
8 class BTree:
9     def __init__(self, degree):
10         self.root = Node(is_leaf=True)
11         self.degree = degree
12
13     def build_from_file(self, file_path, inverted=False):
14         with open(file_path, 'r', encoding='utf-8') as file:
15             for line in file:
16                 key, values_str = line.strip().split(': ')
17                 values = list(map(int, values_str.strip('[]').split(', ')))
18                 if inverted:
19                     key = key[::-1]
20                 self.insert(key, values)
21
22     def insert(self, key, values):
23         if len(self.root.keys) == (2 * self.degree) - 1:
24             new_root = Node()
25             new_root.children.append(self.root)
26             self._split_child(new_root, 0)
27             self.root = new_root
28
29         self._insert_non_full(self.root, key, values)
30
31     def _insert_non_full(self, node, key, values):
32         i = len(node.keys) - 1
33         if node.is_leaf:
34             node.keys.append(None)
35             node.values.append(None)
36             while i >= 0 and key < node.keys[i]:
37                 node.keys[i + 1] = node.keys[i]
38                 node.values[i + 1] = node.values[i]
39                 i -= 1
```

```

40         node.keys[i + 1] = key
41         node.values[i + 1] = values
42     else:
43         while i >= 0 and key < node.keys[i]:
44             i -= 1
45         i += 1
46         if len(node.children[i].keys) == (2 * self.degree) - 1:
47             self._split_child(node, i)
48             if key > node.keys[i]:
49                 i += 1
50             self._insert_non_full(node.children[i], key, values)
51
52     def _split_child(self, parent, i):
53         degree = self.degree
54         child = parent.children[i]
55         new_child = Node(is_leaf=child.is_leaf)
56
57         parent.keys.insert(i, child.keys[degree - 1])
58         parent.values.insert(i, child.values[degree - 1])
59         parent.children.insert(i + 1, new_child)
60
61         new_child.keys = child.keys[degree:]
62         new_child.values = child.values[degree:]
63         child.keys = child.keys[:degree - 1]
64         child.values = child.values[:degree - 1]
65
66         if not child.is_leaf:
67             new_child.children = child.children[degree:]
68             child.children = child.children[:degree]

```

Listing 2: B-tree Construction

Advantages:

- **Efficient Data Storage:** B-trees provide efficient storage and retrieval, supporting quick searches.
- **Scalability:** B-trees handle large datasets effectively by maintaining balanced structures.
- **Support for Inverted Searches:** Constructing both normal and inverted B-trees enables efficient handling of several types of queries.

2.3 Query Processing and Data Retrieval

The final script, `retrieving.py`, loads the B-trees and processes various types of queries to retrieve relevant documents. Key functionalities include:

- Loading B-trees and document data.
- Handling single word, conjunctive, disjunctive, prefix, suffix, and wildcard queries.
- Combining results from multiple query types.

```

1 def search(self, key):
2     return self._search(self.root, key)
3
4 def _search(self, node, key):
5     i = 0
6     while i < len(node.keys) and key > node.keys[i]:
7         i += 1
8     if i < len(node.keys) and key == node.keys[i]:
9         return node.values[i]

```

```

10     elif node.is_leaf:
11         return None
12     else:
13         return self._search(node.children[i], key)

```

Listing 3: Single Word Query

2.3.1 Single Word Query

A single word query searches for documents containing the specified term.

```

1 def single_word_query(btree, word):
2     return btree.search(word)

```

Listing 4: Single Word Query

2.3.2 Conjunctive Query

A conjunctive query returns documents containing all the specified terms.

```

1 def conjunctive_query(btree, words):
2     results = [set(btree.search(word)) for word in words]
3     return list(set.intersection(*results))
4
5 def intersect(lists):
6     if not lists:
7         return []
8
9     sorted_lists = sorted(lists, key=len)
10    return list(reduce(set.intersection, map(set, sorted_lists)))

```

Listing 5: Conjunctive Query

2.3.3 Disjunctive Query

A disjunctive query returns documents containing any of the specified terms.

```

1 def disjunctive_query(btree, words):
2     results = [set(btree.search(word)) for word in words]
3     return list(set.union(*results))
4
5 def union(lists):
6     return list(set(chain.from_iterable(lists)))

```

Listing 6: Disjunctive Query

2.3.4 Prefix Query

A prefix query searches for documents containing words that start with the specified prefix.

```

1 def prefix_query(btree, prefix):
2     return btree.prefix_search(prefix)
3
4 def search_prefix(self, prefix):
5     results = []
6     self._search_prefix(self.root, prefix, results)
7     return results
8 def _search_prefix(self, node, prefix, results):
9     i = 0
10    while i < len(node.keys):
11        if node.keys[i].startswith(prefix):
12            results.append((node.keys[i], node.values[i]))

```

```

13         if node.keys[i] > prefix and not node.is_leaf:
14             self._search_prefix(node.children[i], prefix, results)
15         i += 1
16         if not node.is_leaf:
17             self._search_prefix(node.children[i], prefix, results)

```

Listing 7: Prefix Query

2.3.5 Suffix Query

A suffix query searches for documents containing words that end with the specified suffix.

```

1 def suffix_query(inverted_btree, suffix):
2     return inverted_btree.prefix_search(suffix[::-1])
3
4 def search_suffix(self, suffix):
5     reversed_suffix = suffix[::-1]
6     results = []
7     self._search_prefix(self.root, reversed_suffix, results)
8     return [(key[::-1], values) for key, values in results]

```

Listing 8: Suffix Query

2.3.6 Wildcard Query

A wildcard query searches for documents containing words that match the specified pattern.

```

1 def wildcard_query(btree, inv_btree, pattern):
2     if '*' in pattern:
3         prefix, suffix = pattern.split('*')
4         prefix_results = prefix_query(btree, prefix)
5         suffix_results = suffix_query(inv_btree, suffix[::-1])
6         return list(set(prefix_results).intersection(suffix_results))

```

Listing 9: Wildcard Query

Advantages:

- **Flexible Query Handling:** The system supports a wide range of query types, making it versatile for different search needs.
- **Efficient Retrieval:** Using B-trees for search operations ensures that queries are handled quickly, even with large datasets.

3 Conclusion

The described implementation leverages the strengths of inverted indexes and B-trees to create a robust and efficient information retrieval system. The use of text normalization, stopwords removal, and persistent storage ensures that the system is both efficient and scalable. The support for various query types makes this system highly versatile and powerful for different search needs.