

REPORT DCR Project

Search Engine using positional index and stop words

Introduction

This system is designed to efficiently index and search through various document types, including PDF, TXT, HTM and HTML files. The key features include using stop words, positional indexing, and calculating term frequency-inverse document frequency (TF-IDF) scores.

Stop Words

Stop words are often excluded from indexing and search operations to improve performance and relevance. In my system, I use a predefined set of English stop words to filter out terms that do not contribute significantly to the meaning of the documents.

List of Common Stop Words:

```
# List of common English stop words
STOP_WORDS = set([
    'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll", "you'd", 'your',
    'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', "she's", 'her', 'hers', 'herself', 'it',
    "it's", 'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this',
    'that', "that'll", 'these', 'those', 'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had',
    'having', 'do', 'does', 'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while',
    'of', 'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', 'above',
    'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', 'then', 'once',
    'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', 'most', 'other', 'some',
    'such', 'no', 'nor', 'not', 'only', 'own', 'same', 'so', 'than', 'too', 'very', 's', 't', 'can', 'will', 'just',
    'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', 've', 'y', 'ain', 'aren', "aren't", 'couldn',
    "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn',
    "isn't", 'ma', 'mightn', "mightn't", 'mustn', "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't",
    'wasn', "wasn't", 'weren', "weren't", 'won', "won't", 'wouldn', "wouldn't"
])
```

By excluding these words, we focus on indexing and searching for more meaningful content within the documents.

Positional Indexing

Positional indexing is a powerful technique used in my system to keep track of the positions of words within documents. This allows for advanced search capabilities, including phrase searches and proximity queries.

How Positional Indexing Works:

1. Tokenisation: Text is tokenised into individual words, excluding stop words.
2. Index Creation: For each document, the positions of each token are recorded.
3. Posting Lists: These positions are stored in posting lists, which map each token to its occurrences in the documents.

For example :

Document 1: "The quick brown fox jumps over the lazy dog"

Token: "quick" -> Positions: [1]

Token: "fox" -> Positions: [3]

This indexing method ensures that we can quickly locate the exact occurrences of terms within the documents.

Term Frequency-Inverse Document Frequency (TF-IDF)

TF-IDF is a statistical measure used to evaluate the importance of a word in a document relative to a collection of documents (corpus). My system uses TF-IDF to rank documents based on their relevance to the search query.

TF-IDF Calculation:

1. **Term Frequency (TF):** Measures how frequently a term appears in a document.

$$TF(t, d) = \frac{\text{number of times } t \text{ appears in } d}{\text{total number of terms in } d}$$

$$IDF(t) = \log \frac{N}{1 + df}$$

$$TF - IDF(t, d) = TF(t, d) * IDF(t)$$

2. **Inverse Document Frequency (IDF):** Measures how important a term is in the entire corpus.

3. **TF-IDF Score:** Combines TF and IDF to give the final score.

Example: Term: "fox"

Document Frequency (DF): 1 (appears in 1 out of 10 documents)

IDF: $\log(10 / 1) = 1$

TF in Document 1: 1 / 9

TF-IDF: $(1/9) * 1 \approx 0.111$

Using TF-IDF, effectively ranks documents based on their relevance, improving the accuracy of search results.

Results (snapshots)

Searching the word “music”

1. As one word which is not a phrase just 1 word:

The output shows us Document IDs, the file addresses, total sequences and the relevance score 😊

2. now search for a phrase with random words related to social science, try to add some words related to it and see the output (second, third and fourth images)

Every project in the world has strengths and weaknesses as do my project.

Strengths:

Stop Words Removal:

- Implementing stop words removal helps in reducing the noise and improving the relevance of the search results.

Positional Indexing:

- Keeping track of the positions of words within documents allows for advanced search capabilities like phrase searches and proximity queries.

TF-IDF Calculation:

- Using TF-IDF for scoring search results is a robust method for evaluating the importance of a word in a document relative to the entire corpus.

Index Persistence:

- Saving the index to disk using `pickle` allows for efficient reloading and avoids the need to reprocess documents each time.

Weaknesses:

Error Handling:

- There's limited error handling in the code. For instance, handling cases where a file cannot be opened or read, or when the input directory is invalid.

User Interface:

- The search functionality is interactive through the console, which might not be user-friendly for non-technical users. Consider creating a simple web interface using Flask or Django.

Sincerely yours.