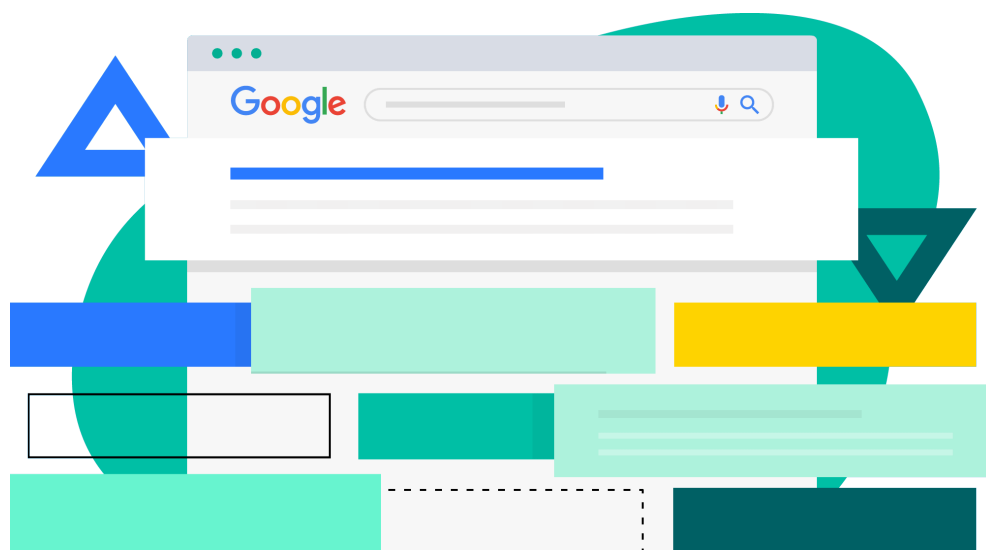


BM25 Ranking



سید محمد حسین هاشمی ۴۰۲۲۳۶۳۱۴۳

فروردین ۱۴۰۳

فهرست مطالب

۲	۱	کتابخانه‌های مورد نیاز
۲	۲	کلاس Inverted Index
۲	۳	متد <code>-init-</code>
۳	۴	متد <code>build_index</code>
۳	۵	متد <code>tokenize</code>
۴	۶	متد <code>bm25_score</code>
۴	۷	متد <code>search</code>
۵	۸	استفاده

۱ کتابخانه‌های مورد نیاز

```
BM25_ranking.py

# Import necessary modules
import re
import sys
import math
from collections import defaultdict
```

کتابخانه‌های مورد نیاز که در پروژه لود شده‌اند.

۲ کلاس Inverted Index

```
BM25_ranking.py

# Define a class to represent an inverted index
class InvertedIndex:
```

تمام کدهای مربوطه در این کلاس نوشته می‌شود

۳ متد __init__

```
BM25_ranking.py

def __init__(self, documents):
    # Initialize the inverted index with the given documents
    self.documents = documents
    # Create an empty index to store the token-to-document mapping
    self.index = defaultdict(list)
    # Create an empty dictionary to store the term frequencies
    self.term_freq = defaultdict(lambda: defaultdict(int))
    # Create an empty dictionary to store the document lengths
    self.doc_len = {}
    # Initialize the average document length to 0
    self.avg_dl = 0
    # Build the inverted index
    self.build_index()
```

بدر اینجا داکيومنت‌ها به عنوان ورودی دریافت شده و پس از تعریف متغیرهای مورد نیاز با فراخوانی `build_index`، Inverted Index ایجاد می‌شود.

۴ متد build_index

```
BM25_ranking.py

def build_index(self):
    # Iterate over the documents and tokenize them
    for doc_id, document in enumerate(self.documents):
        # Tokenize the document
        tokens = self.tokenize(document)
        # Store the length of the document
        self.doc_len[doc_id] = len(tokens)
        # Add to the total document length
        self.avg_d1 += len(tokens)
        # Iterate over the tokens and update the index and term frequencies
        for token in tokens:
            # Add the document ID to the token's posting list
            self.index[token].append(doc_id)
            # Increment the term frequency for the token in the document
            self.term_freq[token][doc_id] += 1
        # Calculate the average document length
    self.avg_d1 /= len(self.documents)
```

در این تابع Inverted Index ساخته می‌شود و علاوه بر آن طول هر داکيومنت به همراه میانگین طول داکيومنت‌ها محاسبه می‌شود.

۵ متد tokenize

```
BM25_ranking

def tokenize(self, text):
    # Tokenize the text by splitting on whitespace and converting to lowercase
    return re.split('[^a-zA-z]', text.lower())
```

در این متد از کلمات از متن استخراج می‌شوند.

۶ متد bm25__score

```
BM25_ranking.py

def bm25_score(self, query, doc_id):
    # Calculate the BM25 score for the query and document
    score = 0
    for token in self.tokenize(query):
        # Check if the token is in the index
        if token not in self.index:
            continue
        # Calculate the document frequency (df) of the token
        df = len(self.index[token])
        # Calculate the inverse document frequency (idf) of the token
        idf = math.log((len(self.documents) - df + 0.5) / (df + 0.5))
        # Calculate the term frequency (tf) of the token in the document
        tf = self.term_freq[token][doc_id]
        # Calculate the document length (dl) of the document
        dl = self.doc_len[doc_id]
        # Calculate the BM25 score for the token
        k1 = 1.2
        b = 0.75
        numerator = tf * (k1 + 1)
        denominator = tf + k1 * (1 - b + b * dl / self.avg_dl)
        score += idf * (numerator / denominator)
    return score
```

در اینجا داکيومنت‌ها بر اساس فرمول BR25 امتیازدهی می‌شوند که فرمول آن به‌شرح زیر است:

$$score(d) = \sum_{t \in q} \left(\frac{f_{t,d} \cdot (k_1 + 1)}{k_1 \cdot \left((1 - b) + b \cdot \frac{|d|}{avgdl} \right) + f_{t,d}} \right) \cdot \log \left(\frac{N - n_t + 0.5}{n_t + 0.5} \right) \cdot \frac{(k_2 + 1) \cdot f_{t,q}}{k_2 + f_{t,q}}$$

۷ متد search

```
BM25_ranking.py

def search(self, query):
    # Calculate the BM25 score for the query and each document
    scores = {}
    for doc_id in range(len(self.documents)):
        scores[doc_id] = self.bm25_score(query, doc_id)
    # Return the sorted scores in descending order
    return sorted(scores.items(), key=lambda x: x[1], reverse=True)
```

در این متد کلمات از رشته ورودی استخراج شده و پس از آن با استفاده از متد bm25__score دسته‌بندی می‌شوند و براساس امتیاز مرتب شده و برگشت داده می‌شود.

```

BM25_ranking.py

# Usage
if __name__ == '__main__':
    # Check if the script is being run from the command line
    if len(sys.argv) != 2:
        print("Usage: python BM25_ranking.py [file_name]")
        sys.exit(1)

    # Get the filename from the command-line argument
    filename = sys.argv[1]

    # Read the documents from the file
    documents = []

    with open(filename, encoding="utf8") as file:
        for line in file:
            documents.append(line.strip())

    # Create an instance of the InvertedIndex class
    ii = InvertedIndex(documents)

    # Prompt the user to enter a query
    query = input('Search: ')

    # Search for the query and get the top 10 results
    results = ii.search(query)
    print("Results for query '%s':" % query)
    for doc_id, score in results[0: 11]:
        print("Doc %d: score: %.5f" % (doc_id, score))

```

برای استفاده از کلاس گفته شده این کد نوشته شده که بعد از فراخوانی کلاس ایجاد شده و پس از خواندن داکيومنت‌ها و ایجاد Inverted Index در آن با استفاده از رنکینگ BM25 جستجو انجام می‌شود و ۱۰ نتیجه اول به همراه امتیاز خروجی داده می‌شود. برای مثال برای ورودی مانند تصویر زیر:

```

first document
second document
third document|

```

خروجی مانند تصویر زیر تولید می‌شود.

```
Search: third
Results for query 'third':
Doc 2: score: 0.51083
Doc 0: score: 0.00000
Doc 1: score: 0.00000
```