

Assignment 2

Arina Goncharova B22-DS-02

April 2025

1 Methodology

1.1 Data preparation

In this stage the data is formatted to include document id, title and content as tab separated values. After that it is stored in hdf5.

1.2 Indexing

I decided to use only 1 MapReduce pipeline because it was enough for me.

- mapper1.py

Processes the files from the folder `/index/data` and produces the lines in the format: word \t document id. .

- reducer1.py

Takes as input the lines produced by mapper, reduces them to find values needed for BM25 calculation and stores them in Cassandra tables.

1.2.1 Cassandra schema

I decided to use improved version of tf-idf formula - BM25.

$$BM25(q, d) = \sum_{t \in q} \log \frac{N}{df(t)} \frac{k_1 \cdot tf(t, d)}{k_1[(1 - b) + b \cdot \frac{dl(d)}{dl_{avg}}] + tf(t, d)}$$

To calculate it I need:

- table `term_doc_freq` - $tf(t, d)$

This table is required to get the number of times the given term occurs in the given document.

```
CREATE TABLE IF NOT EXISTS term_doc_freq (
    term TEXT,
    document_id TEXT,
    term_in_doc_frequency INT,
    PRIMARY KEY (term, document_id)
);
```

Since Cassandra is not efficient in aggregations, I cannot use term_doc_freq table for the $df(t)$, $dl(d)$, dl_{avg} and N calculation. Therefore, I decided to create separate table for this values.

- table term_freq - $df(t)$ This table stores the number of documents for each term where it occurred.

```
CREATE TABLE IF NOT EXISTS term_freq (
    term TEXT PRIMARY KEY,
    term_freq INT
);
```

- table doc_length - $dl(d)$

This table stores the length of each document.

```
CREATE TABLE IF NOT EXISTS doc_length (
    document_id TEXT PRIMARY KEY,
    document_length INT
);
```

- table general_stats - dl_{avg} and N

This table stores the average document length and the total number of documents in the corpus. I decided to keep the average length of the document directly, since the total length of the document is not needed in this task. If it is needed, we can simply multiply the average length by the total number of documents

```
CREATE TABLE IF NOT EXISTS general_stats (
    statistics_name TEXT PRIMARY KEY,
    statistics_value DOUBLE
);
```

1.3 Ranking

Spark RDD allows us to parallelize document processing and scoring, which makes query searching more efficient when we have many large documents.

1.3.1 Score calculation

To calculate score of the document for the given query I am using BM25 score. I added $\epsilon = 1e-6$ to the denominators to avoid division by 0.

```
# epsilon is added to avoid division by 0
def compute_BM25(total_num_of_documents, dl_avg, q_words,
term_freq_dict, term_doc_freq_dict, doc_length_dict, document_id,
k=1, b=0.75, epsilon=1e-6):
    result = 0
```

```

if document_id in doc_length_dict:
    dl = doc_length_dict[document_id]
else:
    dl = 0

for word in q_words:
    if word in term_freq_dict:
        df_t = len(term_doc_freq_dict[word])
        if document_id in term_doc_freq_dict[word]:
            tf_t_d = term_freq_dict[word][document_id]
        else:
            tf_t_d = 0
    else:
        df_t = 0
        tf_t_d = 0
    result += math.log(total_num_of_documents/(df_t+epsilon)) *
    (k + 1) * tf_t_d / (k*(1-b + b*dl/dl_avg) + tf_t_d + epsilon)
return result

```

2 Demonstration

2.1 Run the repository

2.1.1 Use ready search engine with predefined query

1. Clone the repository

```
git clone https://github.com/arinagoncharova2005/big-
data-assignment2-2025.git
```

2. Go to main directory

```
cd big-data-assignment2-2025
```

3. Start the app

```
docker-compose up
```

2.1.2 Search documents with custom query

1. Open master node terminal ans start the script

```
docker exec -it cluster-master
bash search.sh "<your query>"
```

2.1.3 Index document

1. Open master node terminal ans start the script

```
docker exec -it cluster-master  
bash index.sh "<path to file>"(optional)
```

2.2 Results

2.2.1 Uploading data to hdfs





2.2.2 Mapper

To test the mapper I passed a small document in the format produced by data preparation step. File with this example can be found in the repository in folder mapreduce with name `sample_mapper_input.txt`.



2.2.3 Reducer

I passed a small document to the reducer to check that it works. This file with example can be found in the repository in folder mapreduce with name sample_mapper_output.txt.

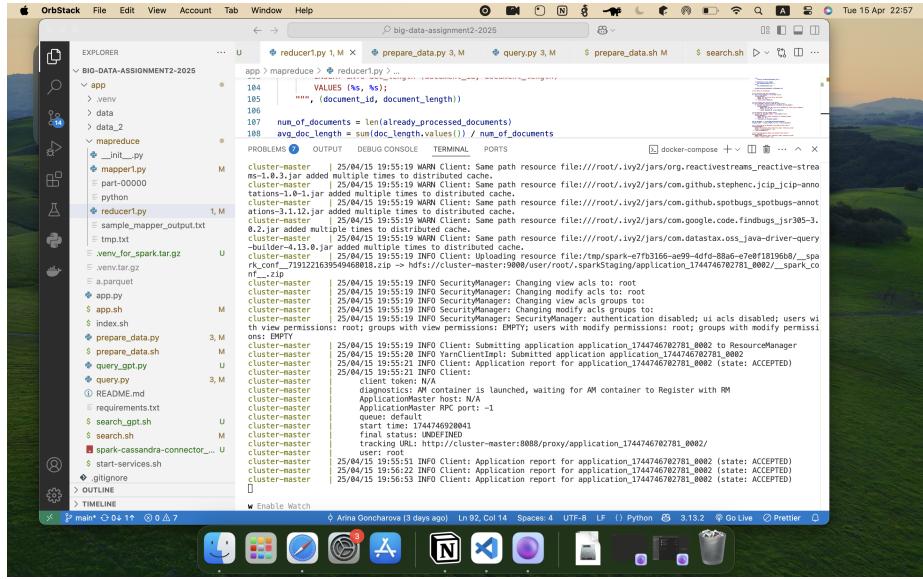


2.2.4 Indexing

```
cd big-data-assignment2-2025
cd app
bash prepare_data.sh
bash index.sh
```

2.2.5 Search engine execution

I have a problem with execution of indexing on slave node. I suppose that problem occurs because of my arm processor.



But on master it runs (replaced –master yarn with –master local).

2.2.6 Reflection

During completion of this assignment I understood the purpose of using MapReduce framework and Spark RDD.

2.3 Difficulties and solutions

- Use OrbStack instead of Docker Desktop.

This program offers optimized Docker usage for macOS.

- Passing python environment to slave node.

Added arguments to the script to pass the archive with the environment

- The problem with execution on slave nodes

I spend several evenings trying to debug my code and understand the problem. But only at the last evening I found that on linux the program executes but on my laptop it is just logging ACCEPTED.

Screenshot of a code editor showing a project structure and a terminal log.

Project Structure:

```

BIG-DATA-ASSIGNMENT2-2025
├── app
│   ├── mapper.py
│   ├── reducer.py
│   ├── query.py
│   ├── requirements.txt
│   ├── search_gpt.sh
│   ├── search.sh
│   └── spark-cassandra-connector...
├── data
└── mapreduce
    ├── __init__.py
    ├── mapper.py
    ├── part-00000
    ├── python
    ├── sample_mapper_input.txt
    ├── sample_mapper_output.txt
    ├── tmp.txt
    └── venv.tar.gz

```

Code Editor (query.py):

```

1  #!/usr/bin/env python3
2  from pyspark.sql import SparkSession
3  from cassandra.cluster import Cluster
4  import sys
5  import re
6  import math
7
8  # get words from the text
9  def tokenize_text(text):
10     return re.findall(r'\w+', text.lower())
11
12 print(sys.argv[1])
13
14 # epsilon is added to avoid division by 0
15 def compute_BWQ(total_num_of_documents, dl_avg, q_words, term_freq_dict, term_doc_freq_dict, doc_length_dict, document_id, k=1,
16                 result=0):
17     if document_id in doc_length_dict:
18         dl = doc_length_dict[document_id]
19     else:
20         dl = 0
21
22     for word in q_words:
23         if word in term_freq_dict:
24             df_t = len(term_doc_freq_dict[word])
25             if document_id in term_doc_freq_dict[word]:
26                 result += (df_t / dl) * math.log((total_num_of_documents / df_t) + epsilon)
27
28     return result
29
30
31 if __name__ == '__main__':
32     document_id = int(sys.argv[1])
33     q_words = tokenize_text(sys.argv[2])
34     total_num_of_documents = int(sys.argv[3])
35     dl_avg = float(sys.argv[4])
36     q_words = set(q_words)
37     term_freq_dict = {}
38     term_doc_freq_dict = {}
39     doc_length_dict = {}
40
41     for word in q_words:
42         if word in term_freq_dict:
43             term_freq_dict[word] += 1
44         else:
45             term_freq_dict[word] = 1
46
47         if word in term_doc_freq_dict:
48             term_doc_freq_dict[word].add(document_id)
49         else:
50             term_doc_freq_dict[word] = {document_id}
51
52         if document_id in doc_length_dict:
53             doc_length_dict[document_id] += 1
54         else:
55             doc_length_dict[document_id] = 1
56
57     result = compute_BWQ(total_num_of_documents, dl_avg, q_words, term_freq_dict, term_doc_freq_dict, doc_length_dict, document_id)
58
59     print(result)

```

Terminal Log (docker-compose logs):

```

cluster-master | 25/04/15 20:27:15 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)
cluster-master | 25/04/15 20:27:15 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)
cluster-master | 25/04/15 20:28:46 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)
cluster-master | 25/04/15 20:29:16 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)
cluster-master | 25/04/15 20:30:17 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)
cluster-master | 25/04/15 20:31:17 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)
cluster-master | 25/04/15 20:31:48 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)
cluster-master | 25/04/15 20:32:08 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)
cluster-master | 25/04/15 20:32:48 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)
cluster-master | 25/04/15 20:33:19 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)
cluster-master | 25/04/15 20:34:19 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)
cluster-master | 25/04/15 20:34:49 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)
cluster-master | 25/04/15 20:35:19 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)
cluster-master | 25/04/15 20:36:28 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)
cluster-master | 25/04/15 20:38:51 INFO Client: Application report for application_1744746702781_0002 (state: ACCEPTED)

```