

BigData Assignment 2

Anastasia Pichugina DS-02

Methodology.

The workflow is exactly the same as described in the assignment description:

- 1) .parquet file is processed to separate .txt files,
- 2) all data is moved to HDFS,
- 3) MapReduce pipeline computes required statistics, and saves in HDFS,
- 4) app.py connects to the cassandra-server, creates the keyspace and tables, and inserts data from the reducer's output,
- 5) query.py reads the input query, takes data from the Cassandra, calculates the BM25 score and prints the resulting queries.

Some implementation details:

- docker-compose up starts the application. The service cluster-master has the entrypoint app.sh that runs all scripts in the correct order: start-services.sh, prepare_data.sh, index.sh, and search.sh.
- start-services.sh starts all services required for Hadoop components running other scripts: HDFS, Yarn, web UI for MapReduce, and prepares Spark.
- prepare_data.sh: loads a.parquet to the HDFS, prepares data using util.py files:
 - prepare_data1.py: goes through the main files and creates new separate .txt files. All files are loaded to the HDFS /data folder as <doc_id>_<doc_title>.txt.
 - creates the RDD object from those files and creates {doc_id}\t{title}\t{content} which will be saved to the /index/data/ in HDFS (part-00000 and _SUCCESS) using saveAsTextFile().

Note! index.sh doesn't accept any arguments. To change amount of documents collected change the n parameter in the prepare_data1.py.

- MapReduce pipeline and data loading to the Cassandra is conducted in the index.sh script. The script runs mapper1.py, reducer1.py, and app.py:
 - mapper1.py: takes the input from the "/index/data" in HDFS, for each documents it processes terms and prints to the stdout info about terms and documents. The prints have 'tags': 'TF' and 'DOC'. Each row starts with one of them so that reducer could distinguish document and terms data.
 - reducer1.py: processes mapper1.py output. Based on tags it collects all required statistics: term frequencies - how many times document x meets the term y, document frequencies - how many documents contain word x, document data (doc_id - doc_length), and global statistics so that query.py doesn't compute them every time: document count and average document length. The reducer's output is collected in the /tmp/index_output folder (part-00000 and _SUCCESS files).

- app.py connects to the Cassandra-server, creates the keyspace “search_engine” and tables, takes the /tmp/index_output/part-00000 content and inserts it into tables.
- Then search.sh is running. This script must be run with an argument - input query. SparkSession is created and connected to the cassandra server. The input query is tokenized using regex. Cassandra tables are loaded and converted to the RDD objects:
 - terms_rdd: rdd with tuples of “term” and “document frequency”,
 - term_frequencies_rdd: a list of tuples, where each tuple is ((doc_id, term), count). Since term frequencies are defined as “amount of term appearance id doc”, each tuple (doc_id, term) is unique,
 - documents_rdd: rdd with tuples of “doc_id” and “doc_length”
 - stats isn’t converted to RDD but all needed values are extracted.
- IDF map is created for each term in the query. BM25 scores are computed and top 10 documents are printed.

Cassandra schema:

Keyspace: “”

Tables:

terms: term text Primary Key, df int. Since the BM25 formula requires document frequency for each term, this table is created.

term_frequencies: term text, doc_id int, tf int, Primary Key is (term, doc_id). Here the primary key is purple, because Term Frequency is defined by term and doc_id

documents: doc_id int Primary Key, length int. For each document we need to know its length, thus this table keeps these values.

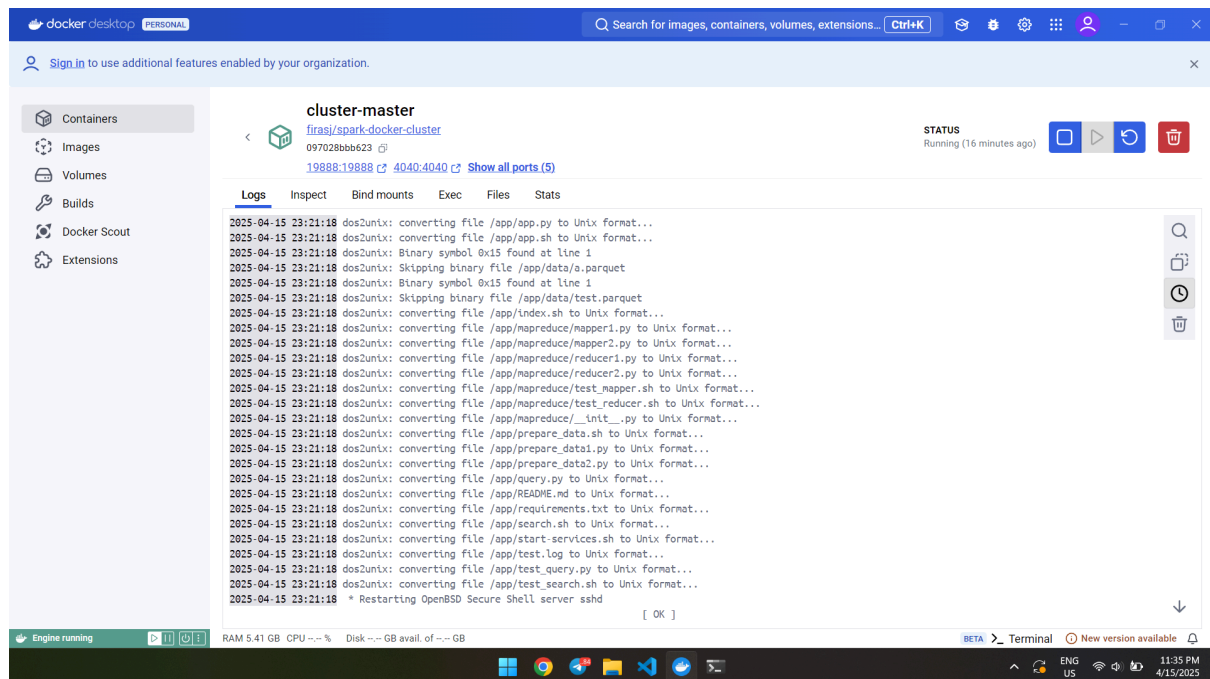
global_stats: key text Primary Key, value int. Number of documents and average document length are constant, so we don’t need to compute them every time processing input query, so this table keeps 2 of these values.

Demonstration.

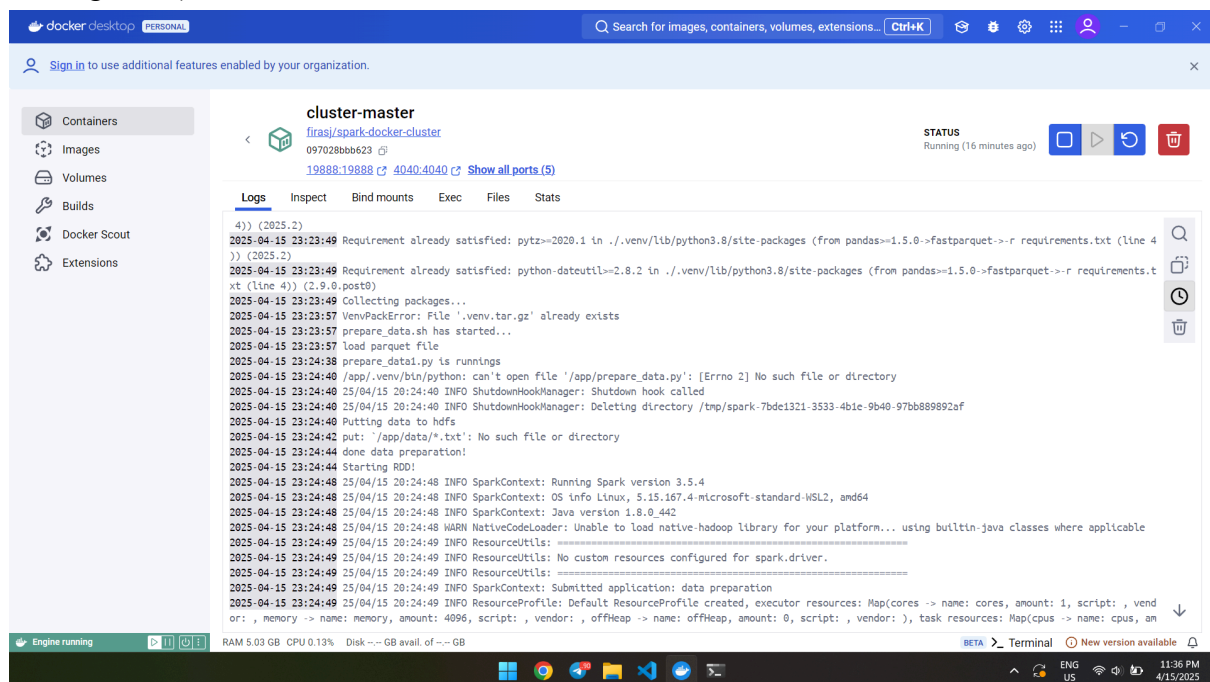
The whole project can be run via docker-compose:

docker-compose up –build

This command will run app.sh that starts all other scripts. While loading you can see logs about converting files:



This is done to avoid errors in Windows-written files which have ‘\r\n’ at the ends of lines. app.sh sets the virtual environment and installs all packages from requirements.txt, runs prepare_data.sh (here you can see logs like “load parquet file”, “prepare_data1.py is running”, etc.).



As a result hdfs will collect all needed data ready for MapReduce pipeline. Next, other scripts will be run. As a result, found documents will be saved to the output.txt.


My results:

I can see all files in HDFS, results of Mapreduce, jobs, Cassandra tables. However, my output.txt is empty, and I assume join on RDDs in the query.py works somehow wrong.

Here I tested the system on test.parquet and global stats are small. But the last output shows the actual stats for a.parquet.

```
root@cluster-master: /app
key | value
-----+-----
docs_total | 5
avg_length | 6.2
(2 rows)
cqlsh:search_engine> select * from stats;
key | value
-----+-----
docs_total | 5
avg_length | 6.2
(2 rows)
cqlsh:search_engine>
PS C:\Users\79876> docker exec -it cassandra-server cqlsh
Connected to Test Cluster at 127.0.0.1:9042
[cqlsh 6.2.0 | Cassandra 5.0.4 | CQL spec 3.4.7 | Native protocol v5]
Use HELP for help.
cqlsh> use search_engine;
cqlsh:search_engine> select * from stats;
key | value
-----+-----
docs_total | 978
avg_length | 573.98975
(2 rows)
cqlsh:search_engine>
PS C:\Users\79876> docker exec -it cassandra-server cqlsh
Connected to Test Cluster at 127.0.0.1:9042
[cqlsh 6.2.0 | Cassandra 5.0.4 | CQL spec 3.4.7 | Native protocol v5]
Use HELP for help.
cqlsh> use serach_engine;
```

Here you can see stages of the application that are finished:



Cluster

About

Nodes

Node Labels

Applications

NEW

NEW SAVING

SUBMITTED

ACCEPTED

RUNNING

FINISHED

FAILED

KILLED

Scheduler

Tools

Cluster Metrics

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Used Resources
2	0	0	2	0	<memory:0 B, vCores:0>

Cluster Nodes Metrics

Active Nodes	Decommissioning Nodes	Decommissioned Nodes
1	0	0

Scheduler Metrics

Scheduler Type	Scheduling Resource Type	Minimum Allocation
Capacity Scheduler	[memory-mb (unit=Mi), vcores]	<memory:1024, vCores:1>

Show 20 entries

ID	User	Name	Application Type	Application Tags	Queue	Application Priority	StartTime	LaunchTime	FinishTime
application_1744748497377_0002	root	TFIDF search	SPARK		default	0	Tue Apr 15 23:25:38 +0300 2025	Tue Apr 15 23:25:39 +0300 2025	Tue Apr 15 23:26:04 +0300 2025
application_1744748497377_0001	root	streamjob4268397438201556323.jar	MAPREDUCE		default	0	Tue Apr 15 23:24:54 +0300 2025	Tue Apr 15 23:24:55 +0300 2025	Tue Apr 15 23:25:04 +0300 2025

Showing 1 to 2 of 2 entries

Here is the screen of the search.sh logs. As I said, output.txt is empty

```
root@cluster-master: /app  Windows PowerShell  Windows PowerShell  + - x
25/04/15 20:47:09 INFO TaskSetManager: Finished task 9.0 in stage 3.0 (TID 37) in 70 ms on cluster-slave-1 (executor 1) (9/14)
25/04/15 20:47:09 INFO TaskSetManager: Starting task 11.0 in stage 3.0 (TID 39) (cluster-slave-1, executor 2, partition 11, PROCESS_LOCAL, 8828 bytes)
25/04/15 20:47:09 INFO TaskSetManager: Finished task 8.0 in stage 3.0 (TID 36) in 73 ms on cluster-slave-1 (executor 2) (10/14)
25/04/15 20:47:09 INFO TaskSetManager: Starting task 12.0 in stage 3.0 (TID 40) (cluster-slave-1, executor 1, partition 12, PROCESS_LOCAL, 8828 bytes)
25/04/15 20:47:09 INFO TaskSetManager: Finished task 10.0 in stage 3.0 (TID 38) in 61 ms on cluster-slave-1 (executor 1) (11/14)
25/04/15 20:47:09 INFO TaskSetManager: Starting task 13.0 in stage 3.0 (TID 41) (cluster-slave-1, executor 2, partition 13, PROCESS_LOCAL, 8828 bytes)
25/04/15 20:47:09 INFO TaskSetManager: Finished task 11.0 in stage 3.0 (TID 39) in 69 ms on cluster-slave-1 (executor 2) (12/14)
25/04/15 20:47:09 INFO TaskSetManager: Finished task 12.0 in stage 3.0 (TID 40) in 58 ms on cluster-slave-1 (executor 1) (13/14)
25/04/15 20:47:09 INFO TaskSetManager: Finished task 13.0 in stage 3.0 (TID 41) in 59 ms on cluster-slave-1 (executor 2) (14/14)
25/04/15 20:47:09 INFO YarnScheduler: Removed TaskSet 3.0, whose tasks have all completed, from pool
25/04/15 20:47:09 INFO DAGScheduler: ResultStage 3 (takeOrdered at /app/query.py:76) finished in 0.882 s
25/04/15 20:47:09 INFO DAGScheduler: Job 2 is finished. Cancelling potential speculative or zombie tasks for this job
25/04/15 20:47:09 INFO YarnScheduler: Killing all running tasks in stage 3: Stage finished
25/04/15 20:47:09 INFO DAGScheduler: Job 2 finished: takeOrdered at /app/query.py:76, took 1.499595 s
25/04/15 20:47:09 INFO SparkContext: Starting job: collectAsMap at /app/query.py:78
25/04/15 20:47:09 INFO DAGScheduler: Job 3 finished: collectAsMap at /app/query.py:78, took 0.000256 s
25/04/15 20:47:09 INFO SparkContext: SparkContext is stopping with exitCode 0.
25/04/15 20:47:09 INFO SparkUI: Stopped Spark web UI at http://cluster-master:4040
25/04/15 20:47:09 INFO YarnClientSchedulerBackend: Interrupting monitor thread
25/04/15 20:47:09 INFO YarnClientSchedulerBackend: Shutting down all executors
25/04/15 20:47:09 INFO YarnSchedulerBackend$YarnDriverEndpoint: Asking each executor to shut down
25/04/15 20:47:09 INFO YarnClientSchedulerBackend: YARN client scheduler backend Stopped
25/04/15 20:47:09 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
25/04/15 20:47:09 INFO MemoryStore: MemoryStore cleared
25/04/15 20:47:09 INFO BlockManager: BlockManager stopped
25/04/15 20:47:09 INFO BlockManagerMaster: BlockManagerMaster stopped
25/04/15 20:47:09 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
25/04/15 20:47:09 INFO SparkContext: Successfully stopped SparkContext
25/04/15 20:47:10 INFO ShutdownHookManager: Shutdown hook called
25/04/15 20:47:10 INFO ShutdownHookManager: Deleting directory /tmp/spark-30d7800d-19f3-4ba4-9c52-c5b4caf04c9f
25/04/15 20:47:10 INFO ShutdownHookManager: Deleting directory /tmp/spark-b03ad22f-23f4-4221-ba20-9a70ab61a6b5
25/04/15 20:47:10 INFO ShutdownHookManager: Deleting directory /tmp/spark-b03ad22f-23f4-4221-ba20-9a70ab61a6b5/pyspark-5d15c560-715e-460c-867b-c697edc108f5
25/04/15 20:47:10 INFO CassandraConnector: Disconnected from Cassandra cluster.
25/04/15 20:47:10 INFO SerialShutdownHooks: Successfully executed shutdown hook: Clearing session cache for C* connector
root@cluster-master:/app# ls
README.md  app.sh  index.sh  ouput.txt  prepare_data.py  query.py  search.sh  start-services.sh
app.py     main.py  prepare_data.sh  requirements.txt  spark-warehouse  test.log
root@cluster-master:/app# cat output.txt
cat: output.txt: No such file or directory
root@cluster-master:/app# cat ouput.txt
root@cluster-master:/app#
```

I also added “tail -f /dev/null” to the end of the app.sh file. So you can run
docker exec -it cluster-master bash
./search.sh “your query”