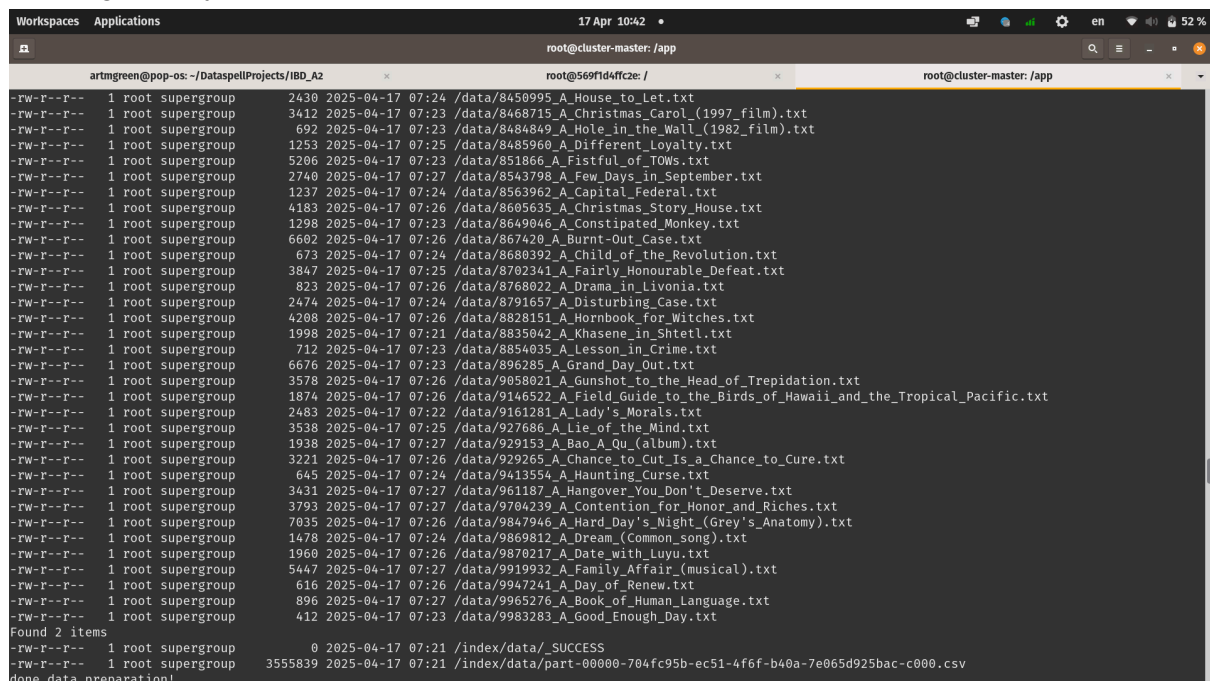# [S25] IBD A2 Report

Matevosian Artem, B22-DS-01

## Methodology: Data preparation

The prepare_data.py from the original assignment statement managed to satisfy my needs after uncommenting writing to csv operation, therefore, I directly used the script from the assignment - it reads parquet file, selects 1000 entries, creates docs + writes to csv in the format <doc_id>\t<title>\t<text>

- I spent an abnormal amount of time trying to fix Out-Of-Memory errors that appeared independently of spark vectorized reader configs, later these reduced in frequency when I started deleting all the files from the previous runs right before main script execution.

## Demonstration: Data preparation

Given that the script is directly taken from the assignment statement with minimal modifications for successful preparation on my low-resource machine, and the logs during its execution are extremely long, it does not appear feasible (or even needed) to print them out here. However, for the completeness purposes here is the screenshot of how the output end should generally look like whenever the preparation is successful:



NOTE: the screenshot is done while app.sh was running (so no commands were typed by hand here now, everything is in script)

# Methodology: Indexer

The result after running indexing: cassandra-server gets the tables as listed below
(P = partitioning key, C = clustering key)

| term_frequencies | | |
|---|---|---|
| term | TEXT | P |
| document_id | INT | C (ASC) |
| tf | INT | |

| document_frequencies | | |
|---|---|---|
| term | TEXT | P |
| df | INT | |

| doc_index | | |
|---|---|---|
| document_id | INT | P |
| title | TEXT | |
| length | INT | |

Why exactly these tables?
- term_frequencies = **tf(t, d)**
- document_frequencies = **df(t)**
- doc_index serves as mapping of document_id to title AND as **dl(d)**

For each table, there is a specific mapreduce job definition, ordered as below:
1. term_frequencies
2. document_frequencies
3. doc_index

The mapreduce jobs print out to /index/data in HDFS in the format directly transferable to cassandra, given the table definitions as above. There were multiple options to copy these outputs into Cassandra:
- (Selected by me) run consecutive INSERTS for every entry – able to run from wherever we have access to cassandra-server, no need to load the whole tsv locally when reading it lazily from HDFS
- run CQLSH COPY on TSVs (rejected since it supposedly requires to load outputs from HDFS into local filesystem)
- sstableloader after building same tables on the running machine (rejected for this very reason, cassandra-server should store the tables, I did not want to create them "locally" and then push somewhere else)

As a result, there are three tables defined as above.

NOTE: document_frequencies could be created much more quickly via
SELECT term, COUNT(document_id) FROM term_frequencies GROUP BY term;
However, according to Apache's docs, Cassandra's materialized views prohibit aggregate functions in their definitions. Also, printing the query result into a TSV for a subsequent CQLSH COPY may or may not be faster than running the second MapReduce pipeline, so I

stayed with the second pipeline. Finally, if SQL-based DBs had been available for the assignment (e.g. Postgres), it would have been possible to CREATE table AS SELECT statement as shown above.

# Demonstration: indexer's result

Right after the indexer's start the console should generally look like this:



For each of the MapReduce jobs the output in the end should look similar to this:

For the purposes of demonstration of indexer's results, this is the screenshot regarding Cassandra's tables (I executed bash in cassandra-server container to be able to CQLSH it):



# Methodology: Ranker

Previous approach: wrap everything into a BM25_Calculator class and broadcast it over the workers. Rejected since the class contains Cassandra session pointer, which cannot be packed ("pickled") to broadcast.

Current approach that one can see in query.py:
init_cassandra(): a function to give session credentials to whoever called it, e.g. workers to receive their own session (since it cannot be serialized and passed from main executor). Additionally it computes N and $dl_{avg}$ (see query.py)

compute_bm25(query_terms, doc_id):
1. receives its own Cassandra session
2. checks if the suggested document exists at all (and returns empty title + 0 score if it does not)
3. iterates over every unique term in query and computes "scary-looking formula" from the assignment statement (+0.5 in logarithm numerator and denominator to prevent numerical errors, given that logarithms do not act nicely with zeros) for every term, sums over the terms and returns the sum with doc info
4. Shutdowns session — I do not want to deal with errors that may or may not appear if it does not.

Top-level code:
1. parses query into terms
2. configures Spark

3. receives its own Cassandra session
4. finds all the documents that have non-zero chance to be relevant (tf(doc, term) ≠ 0)
5. computes BM25 in parallel (query is broadcasted, Cassandra session is created by workers on their own. Broadcasting N and dl_avg is a leftover for the testing purposes)

# Demonstration: Ranker

bash search.sh "clockwork orange milk"



bash search.sh "war peace"

NOTE: I cut off results with BM25=0, so less than 10 documents appearing is fine. Also, for these screenshots I typed the listed commands on my own into cluster-master's bash to be able to choose the most interesting ones and screen the outputs properly, but the same results should appear after running app.sh as well.