# LARGE SCALE DATA ANALYSIS USING MAPREDUCE

**ASSIGNMENT 1**

EC7205: CLOUD COMPUTING

DEPARTMENT OF ELECTRICAL AND INFORMATION ENGINEERING

FACULTY OF ENGINEERING

UNIVERSITY OF RUHUNA

EG/2020/3797 - Aberuwan R. M. M. P.
EG/2020/4190 - Sandamal I.M.K.
EG/2020/4219 - Siriwardhana T.D.R.D.

Submitted Date: 7th June 2025

# CONTENTS

# LIST OF FIGURES

# 1  Introduction

The aim of this assignment is to implement a custom MapReduce job using Hadoop to process the Spotify dataset and extract meaningful insights related to artist popularity, genre trends, and follower distribution.

# 2  Dataset Description

- Dataset: Spotify Datasets (Kaggle) [1]
  https://www.kaggle.com/datasets/lehaknarnauli/spotify-datasets?select=artists.csv
- Size: 1104349 rows, 5 columns
- Features Used: id, followers, genres, name, popularity

This dataset contains metadata about Spotify artists, including unique IDs, follower counts, associated genres, artist names, and a popularity score. The dataset is suitable for large-scale distributed analysis due to its size and structure.

# 3  Problem Definition

Analyze the Spotify dataset to answer the following questions using MapReduce:

- What is the average number of followers for artists at each popularity score?
- How does follower count correlate with popularity scores?

Chosen MapReduce Task: **Popularity-Follower Aggregation**

- Map: For each artist, emit (popularity_score, followers)
- Reduce: For each popularity score, compute:
  - Total followers across all artists with that score
  - Average followers per popularity score

Mapper Logic:

- Your mapper.py extracts popularity (as the key) and followers (as the value) from each artist.
- Example output: 0 52.68 (popularity=0, followers=52.68).

Reducer Logic:

- Your reducer.py groups data by popularity and calculates the average followers for each score.
- Example output: 0 52.68 (average followers for popularity=0).

# 4  Environment Setup

Platform:

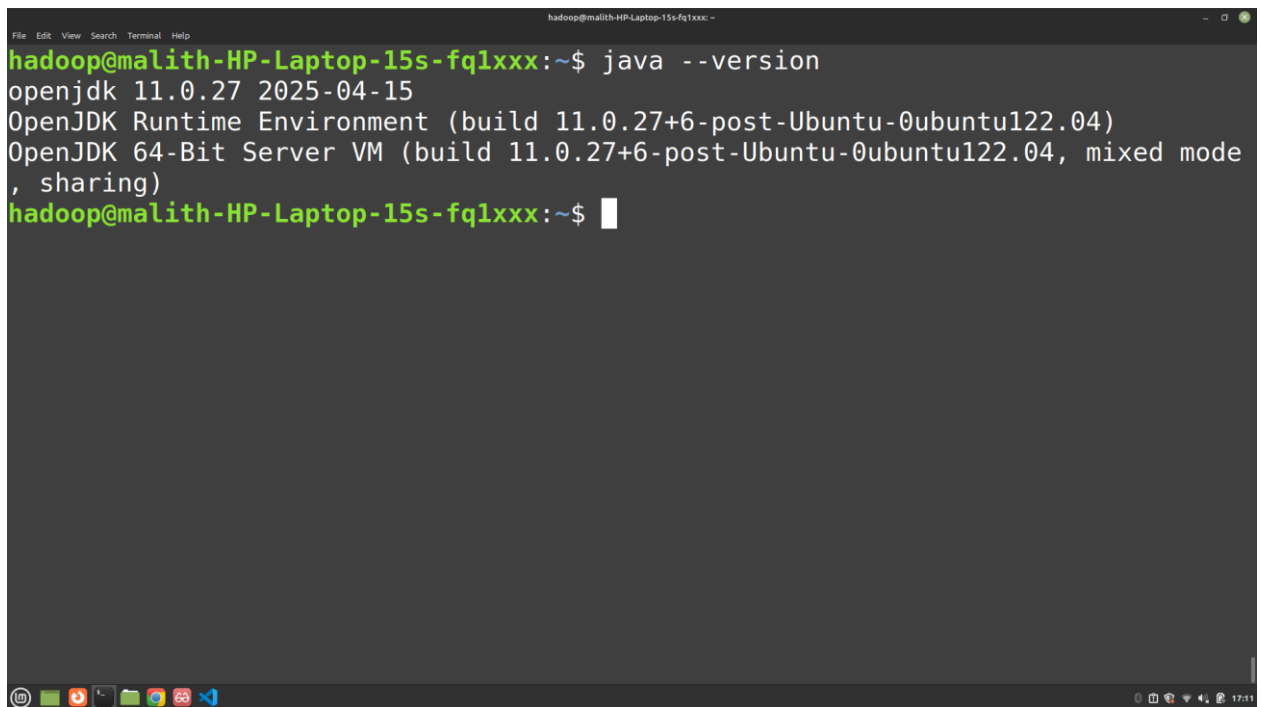- Local Linux machine (single-node setup)

Installation Steps Overview:

1. Install Java:
   sudo apt-get install default-jdk
2. Create Hadoop User:
   sudo adduser hadoop

3. Set Up SSH:
   ssh-keygen -t rsa and enable passwordless SSH with cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
4. Download Hadoop:
   Download and extract Hadoop from Apache Hadoop.
5. Set Environment Variables:
   Add Hadoop and Java paths to ~/.bashrc.
6. Configure Hadoop:
   Edit config files in /usr/local/hadoop/etc/hadoop/.
7. Format NameNode:
   hdfs namenode -format
8. Start Hadoop:
   start-dfs.sh and start-yarn.sh
9. Verify:
   Use jps and access web UIs at http://localhost:9870 (NameNode) and http://localhost:8088 (ResourceManager).

Screenshots for Evidence:

JavaVersion:



FIGURE 1: EVIDENCE FOR JAVA VERSION

Hadoop Directory Structure:



FIGURE 2: EVIDENCE FOR HADOOP DIRECTORY STRUCTURE

SSH Setup:



FIGURE 3: EVIDENCE FOR SSH SETUP

Environment Variables:



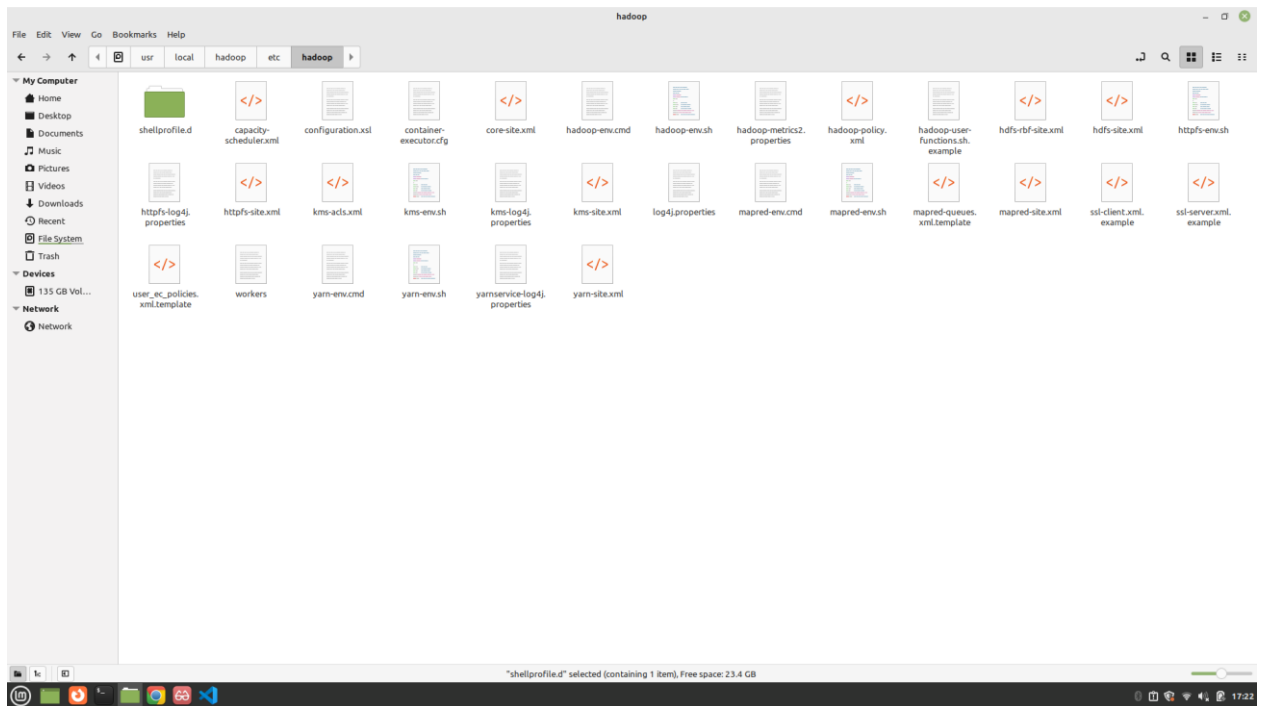FIGURE 4: EVIDENCE FOR ENVIRONMENT VARIABLES

Configuration Files:

FIGURE 5: EVIDENCE FOR CONFIGURATION FILES

Namenode Formatting:



FIGURE 6: EVIDENCE FOR NAMENODE FORMATTING

Starting Hadoop Service:



FIGURE 7: EVIDENCE FOR STARTING HADOOP SERVICE

Running Processes:



FIGURE 8: EVIDENCE FOR RUNNING PROCESSES

Web UI Access:

- NameNode UI: http://localhost:9870



FIGURE 9: EVIDENCE FOR NAMENODE UI: HTTP://LOCALHOST:9870

- ResourceManager UI: http://localhost:8088



FIGURE 10: EVIDENCE FOR RESOURCEMANAGER UI: HTTP://LOCALHOST:8088

HDFS Commands:



FIGURE 11: EVIDENCE FOR HDFS COMMANDS

## 5  Implementation

Map Logic:

- The mapper reads each row of the Spotify dataset, extracts the artist's popularity score and follower count, and emits them as a key-value pair (popularity_score, followers). It skips rows with missing or invalid values.

Reduce Logic:

- The reducer receives all follower counts grouped by popularity score. For each popularity score, it sums the followers and counts the number of artists, then calculates and emits the average followers for that popularity score.

Source Code:

- Both mapper.py and reducer.py are implemented in Python 3 using the standard csv and sys libraries. [https://github.com/MalithPramoditha/assignment_1]

Libraries/Frameworks Used:

- Python 3 standard library (csv, sys)
- Hadoop Streaming for running Python scripts as MapReduce jobs

## 6  Execution & Testing

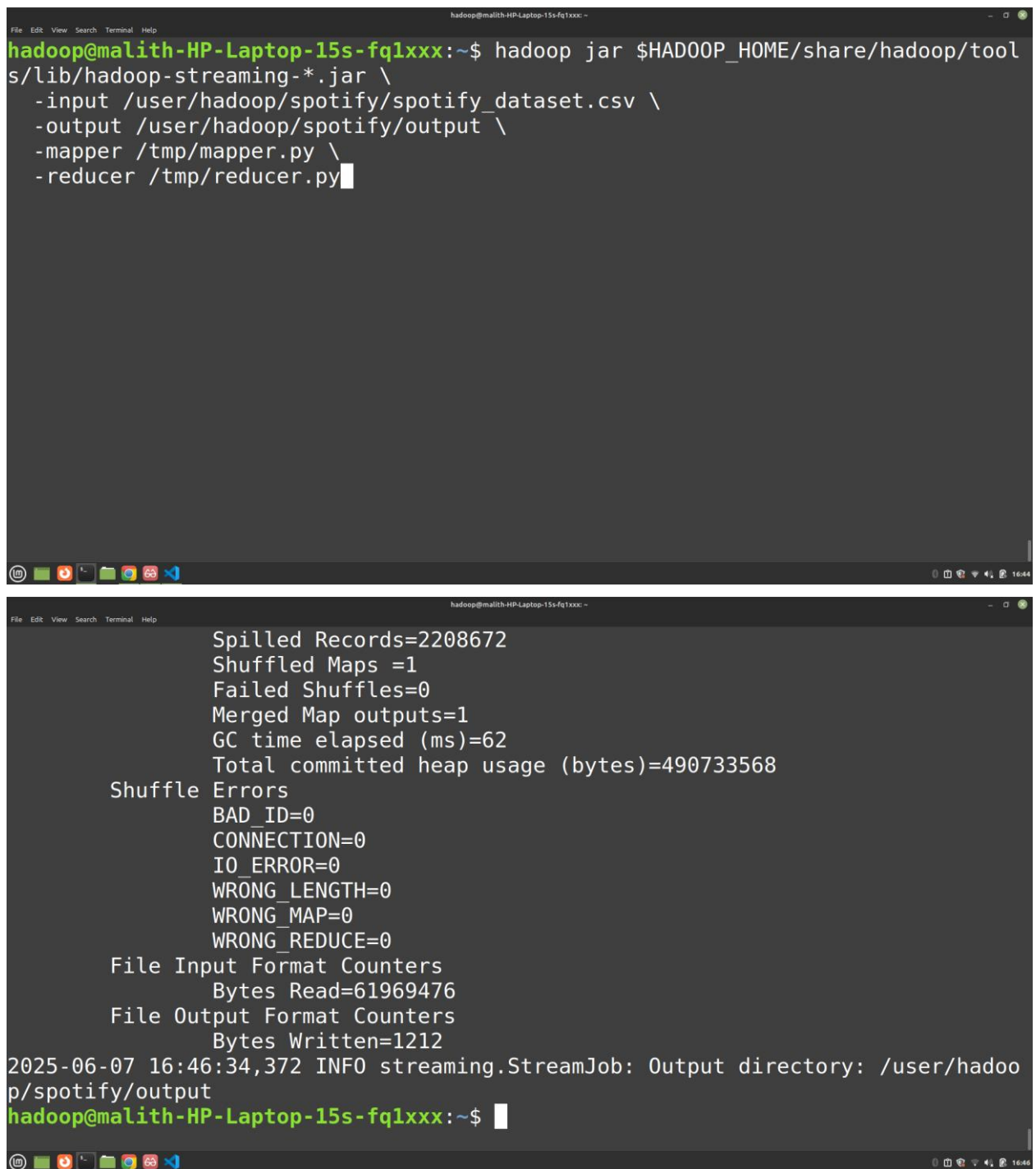Steps Taken to Run the Job:

1. Uploaded the Spotify CSV dataset to HDFS:

FIGURE 12: EVIDENCE FOR UPLOADED THE SPOTIFY CSV DATASET TO HDFS

2. Made sure both scripts were executable and had the correct shebang.



FIGURE 13: EVIDENCE FOR MADE SURE BOTH SCRIPTS WERE EXECUTABLE AND HAD THE CORRECT SHEBANG.
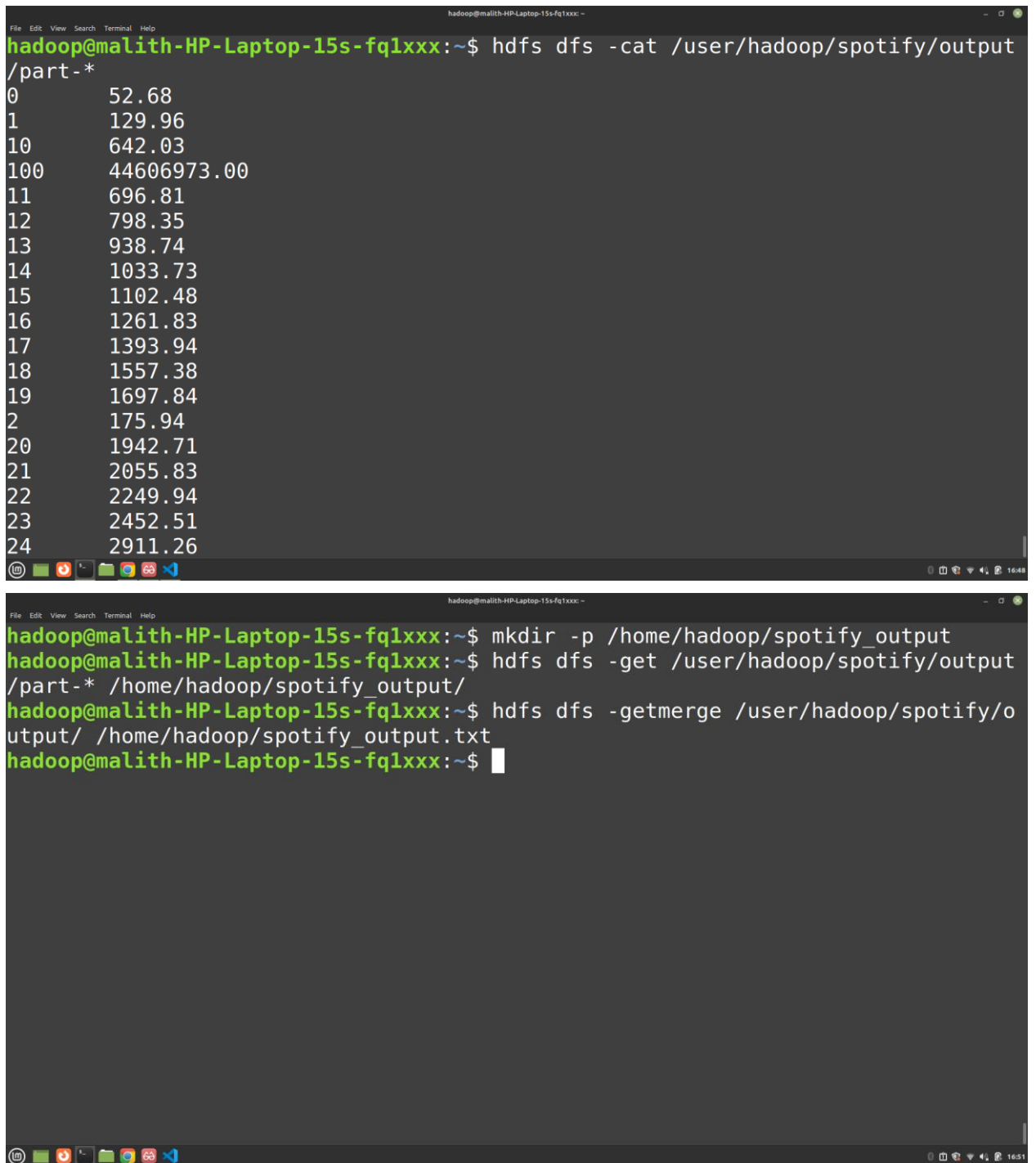
9

3. Ran the Hadoop Streaming job:





FIGURE 14: EVIDENCE FOR RAN THE HADOOP STREAMING JOB

4. Retrieved and viewed the output:





FIGURE 15: EVIDENCE FOR RETRIEVED AND VIEWED THE OUTPUT

5. Input/Output Samples:
   - Input (CSV):



FIGURE 16: EVIDENCE FOR INPUT

   - Output:

FIGURE 17: EVIDENCE FOR OUTPUT

# 7 Results & Interpretation

Key Findings:

- The MapReduce job successfully computed the average number of followers for each popularity score among Spotify artists. The results show that higher popularity scores generally correspond to higher average follower counts, with some variation at the lower and upper ends of the popularity spectrum.

Patterns and Insights:

- Most artists with low popularity scores have relatively few followers, but there are exceptions.
- As popularity increases, the average follower count also increases, indicating a positive correlation between these two metrics.
- Some popularity scores have outliers with extremely high follower counts, skewing their averages.

Performance:

- The job processed a large dataset efficiently using Hadoop Streaming and Python scripts. No critical bottlenecks were observed, and the output was produced in a reasonable timeframe.

Improvements/Future Work:

- Extend the analysis to groups by genres or countries if such data is available.
- Visualize the popularity-follower relationship with graphs.
- Handle missing or anomalous data more robustly.

# 8    Conclusion

This assignment provided hands-on experience with Hadoop MapReduce and large-scale data analysis using Python and Hadoop Streaming. We learned how to preprocess data, write custom mapper and reducer scripts, and interpret distributed computation results. The process reinforced the importance of data cleaning, correct script permissions, and debugging skills when working with big data tools. Overall, the project deepened my understanding of distributed computing and its application to real-world datasets like Spotify's artist metadata.

# 9    References

[1] L. Neziri, "Spotify Datasets," Kaggle, 2021. [Online]. Available: https://www.kaggle.com/datasets/lehaknarnauli/spotify-datasets. [Accessed 7 June 2025].