# CSCI 5408

# DATA MANAGEMENT AND WAREHOUSING

# LAB ASSIGNMENT - 5

**Banner ID:** B00952865

**GitLab Assignment Link:**
https://git.cs.dal.ca/apurohit/CSCI5408_F23_B00952865_AdityaMaheshbhai_Puro hit/-/tree/main/Lab5

# Table of Contents

# Spark Set-up

**Step-1:** Goto Cloud Dataproc service in Google cloud platform using the search option and click create cluster.
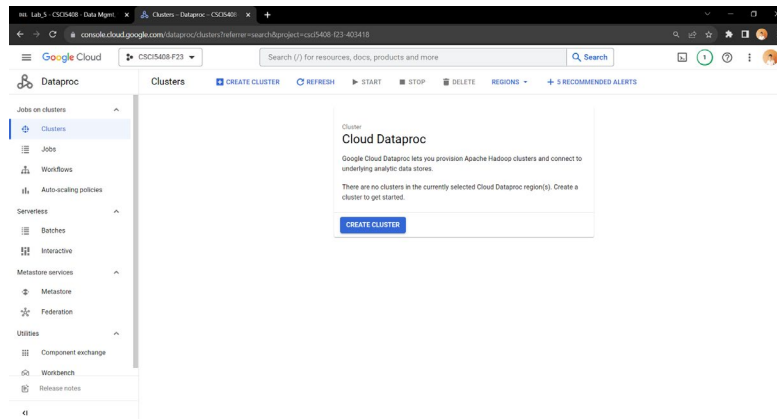


*Figure 1: GCP Dataproc homescreen [1].*

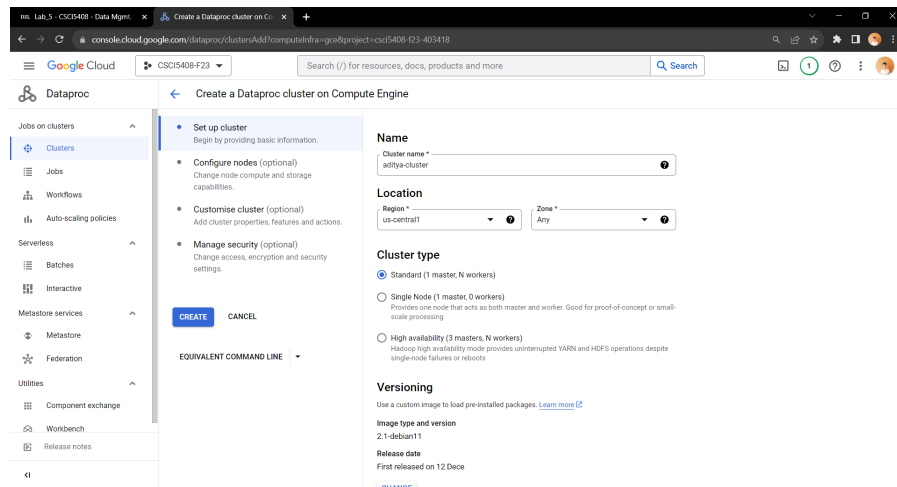**Step-2:** Give any convenient name to the cluster.



*Figure 2: Creating cluster [1].*

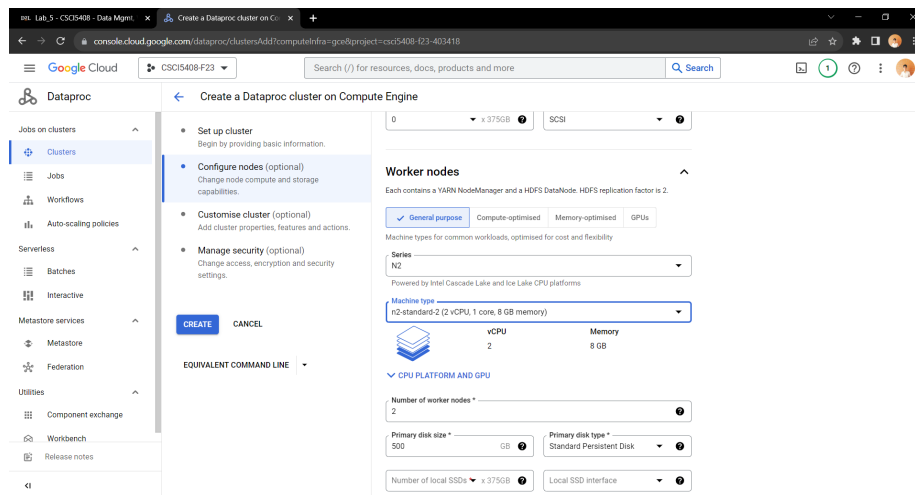**Step-3:** Reduced the vCPUs of worked nodes to 2, due to usage limits of my account.



*Figure 3: Worker nodes vCPU reduced [1].*

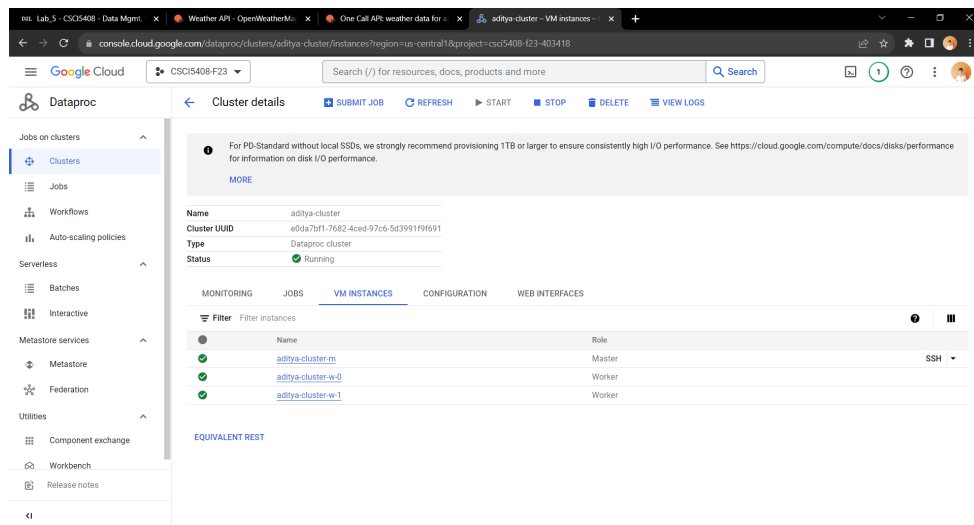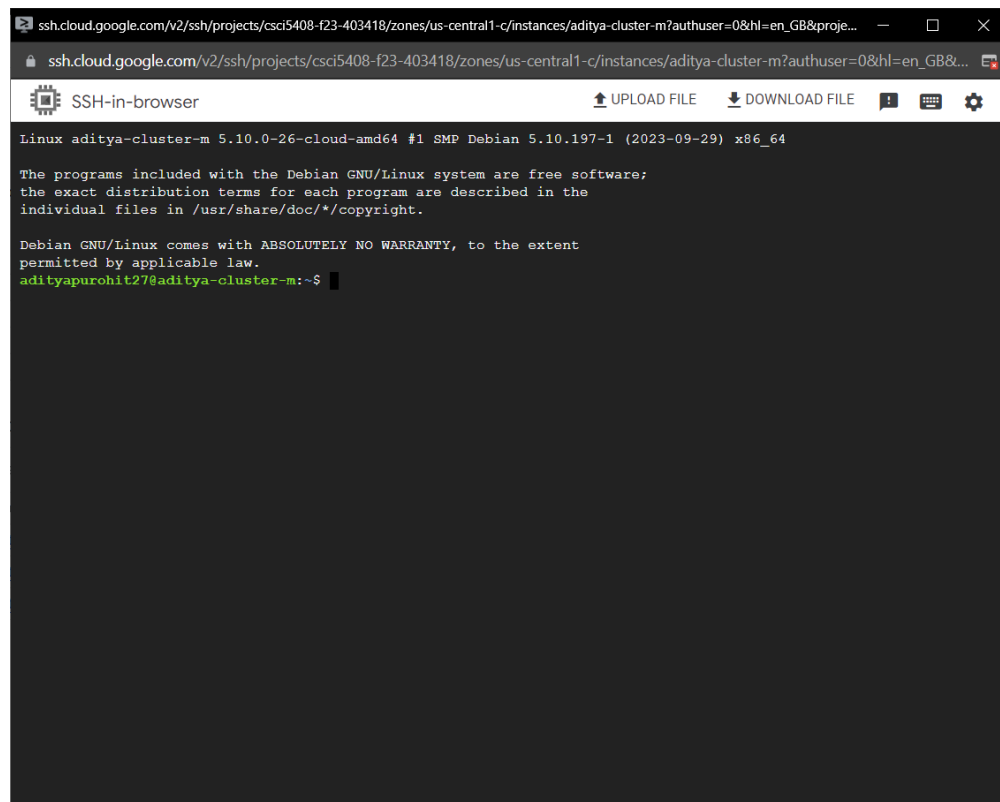**Step-4:** After a while, you should see all the nodes of cluster running.



*Figure 4: Cluster running [1].*

**Step-5:** Connect to master node using ssh.



*Figure 5: Connected to master node using SSH [1].*

# Java Program

Firstly, I have added 2 spark dependencies (core and sql) so that I can use them in my java program.
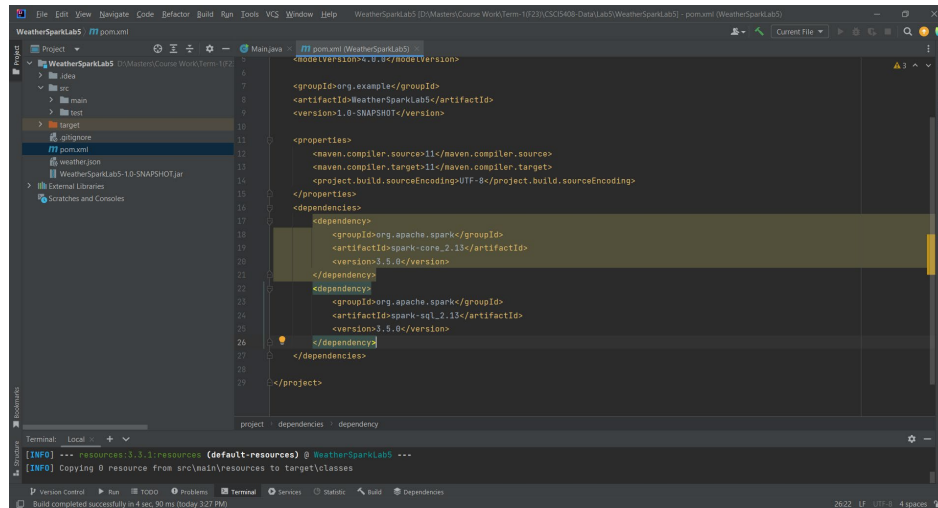


*Figure 6: Adding spark dependencies in pom.xml [2] [3] [4] [5].*

After that, I wrote the below program to read the weather.json file, filter the data where feels_like during day is less than 15. Exclusion of current, minutely, and hourly fields was not needed as they were not present in the original weather.json file.
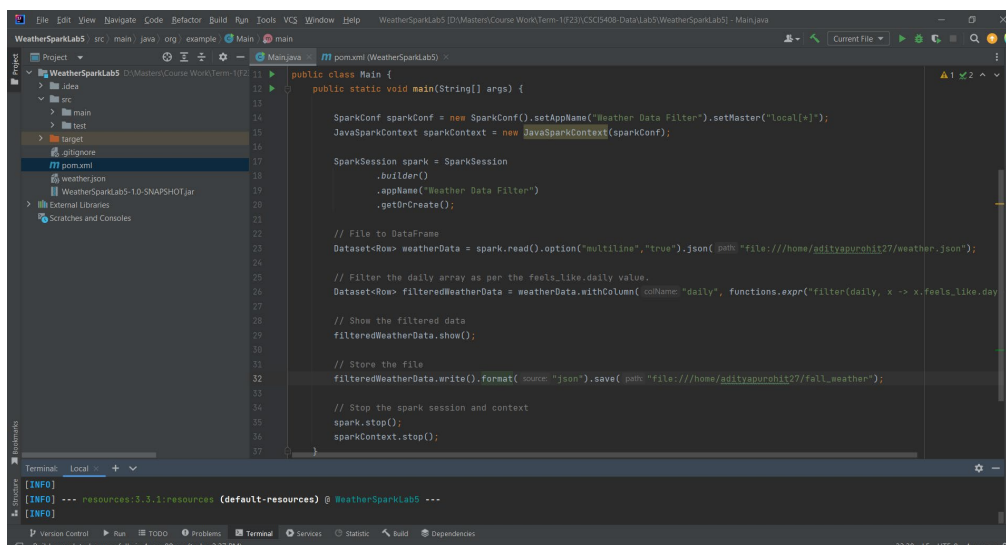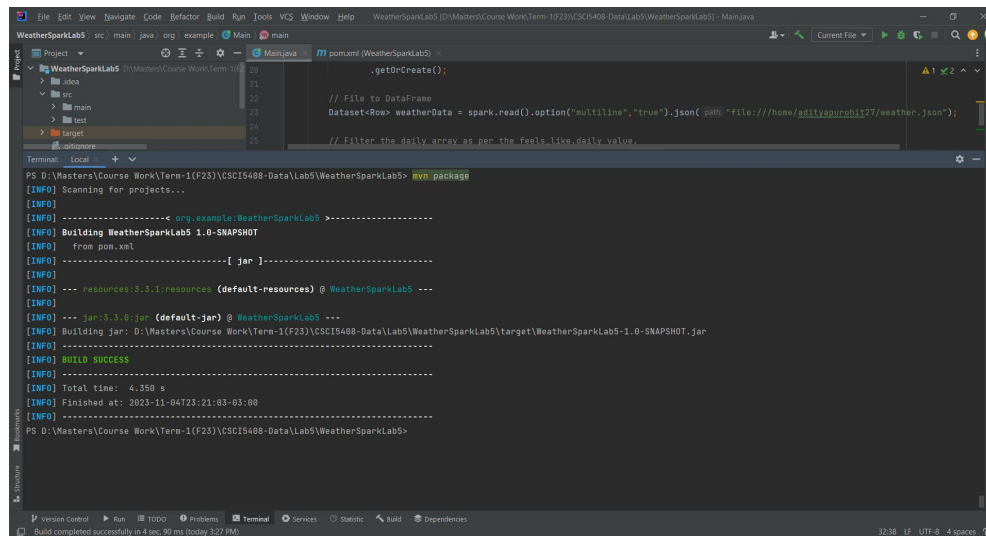


*Figure 7: Overview of the java program that uses spark dependencies [2] [3] [4] [5] [6].*

In the program I have firstly set the JavaSparkContext and started the SparkSession so that spark can do the further processing. Then, I loaded the weather.json from the master node into a spark

DataSet. After that, I used functions.expr of spark sql, to filter the "daily" array inside the weather.json. x inside the expr represents each json object inside the "daily" array and the x.feels_like.day is use to access the day temperature and check it with the 15 with less than operator to get only the needed entries from the "daily" array. This filter the daily array. Lastly, the filtered dataset it saved in a json under the fall_weather directory and then the spark session and context are closed [7] [8].

## Execution of Java Program

**Step-1:** Run mvn package command to create a .jar file of java program.



Figure 8: Build the jar file [2] [3] [4] [5].

**Step-2:** Upload the jar file and weather.json file (from the teams channel) into the master node, using the ssh window.
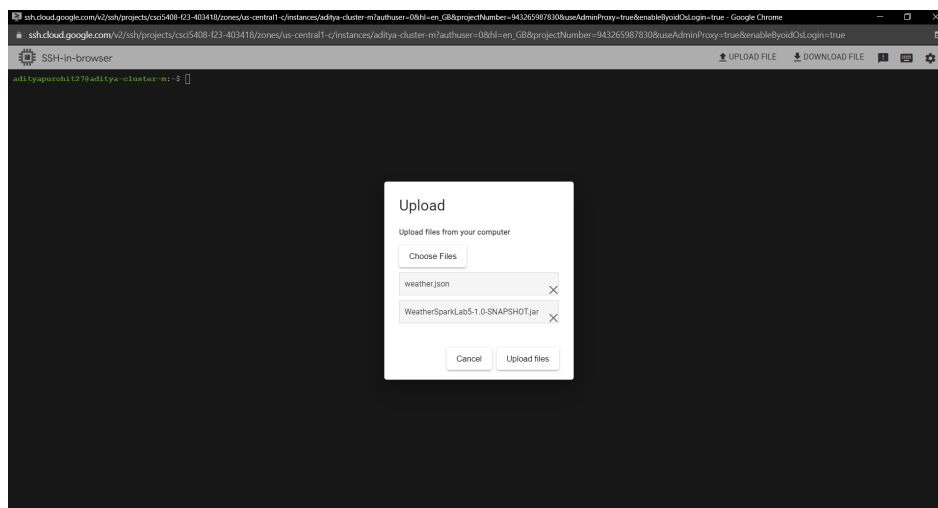


Figure 9: Upload the JAR and JSON [1].

**Step-3:** Run the jar file using the spark command: "spark-submit --class org.example.Main WeatherSparkLab5-1.0-SNAPSHOT.jar". Here the –class tag specifies the starting point of your java program.
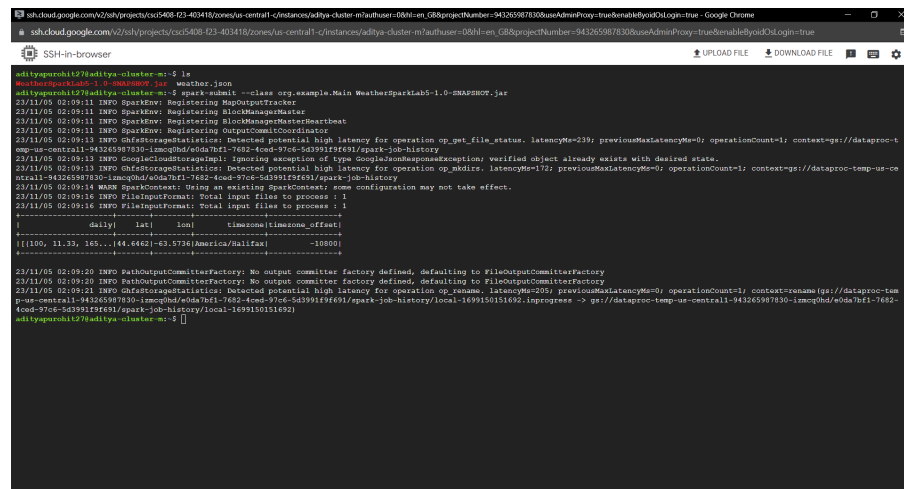


*Figure 10: Run the JAR file [1].*

**Step-4:** Check the output json file under the fall_weather folder and download it and give it a suitable name using linux "mv" command.
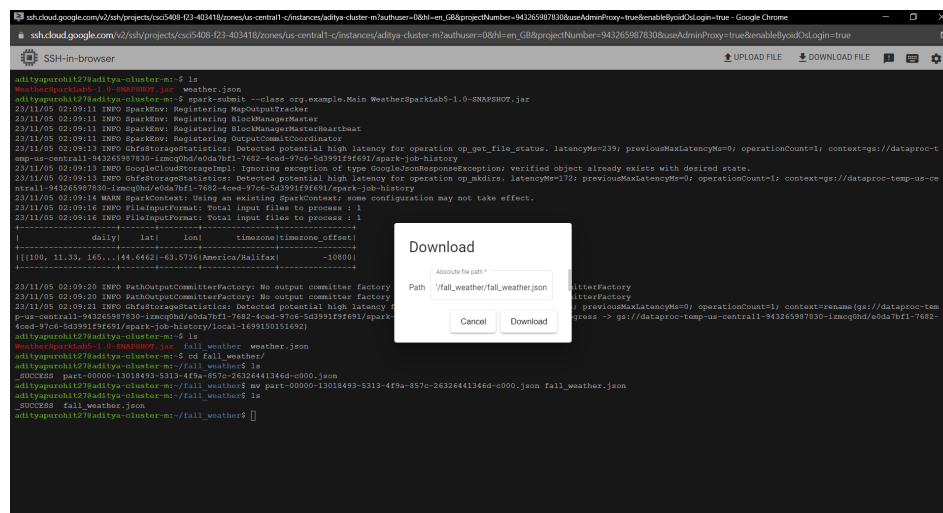


*Figure 11: Rename the output JSON and download it [1].*

*Figure 12: Successful download of the fall_weather.json file [1].*

# References:

[1]  "Dataproc," *Google Cloud*. [Online]. Available: https://cloud.google.com/dataproc?hl=en. [Accessed: 05-Nov-2023].

[2]  "IntelliJ IDEA – the leading Java and Kotlin IDE," *JetBrains*. [Online]. Available: https://www.jetbrains.com/idea/. [Accessed: 05-Nov-2023].

[3]  B. Porter, J. van Zyl, and O. Lamy, "Welcome to Apache maven," *Apache.org*. [Online]. Available: https://maven.apache.org/. [Accessed: 05-Nov-2023].

[4]  "spark-core," *Mvnrepository.com*. [Online]. Available: https://mvnrepository.com/artifact/org.apache.spark/spark-core. [Accessed: 05-Nov-2023].

[5]  "spark-sql," *Mvnrepository.com*. [Online]. Available: https://mvnrepository.com/artifact/org.apache.spark/spark-sql. [Accessed: 05-Nov-2023].

[6]  "Java | Oracle," *Java.com*. [Online]. Available: https://www.java.com/en/. [Accessed: 05-Nov-2023].

[7]  "Getting Started," *Apache.org*. [Online]. Available: https://spark.apache.org/docs/latest/sql-getting-started.html. [Accessed: 05-Nov-2023].

[8]  "Functions (spark 3.0.2 JavaDoc)," Apache.org, 16-Feb-2021. [Online]. Available: https://spark.apache.org/docs/3.0.2/api/java/org/apache/spark/sql/functions.html. [Accessed: 05-Nov-2023].