# CSCI 5408

# DATA MANAGEMENT

# AND

# WAREHOUSING



# LAB ASSIGNMENT - 5

**Submitted By:**   Kenil Shaileshkumar Patel
                    (kenil.patel@dal.ca)
**Banner ID:**      B00954251
**Submitted On:**   November 04, 2023

**Gitlab Repository Link**

https://git.cs.dal.ca/kenil/csci5408_f23_b00954251_kenil_patel/-/tree/main/Lab5

**Table of Contents**

# Setting up Apache Spark on GCP instances

Created a Dataproc cluster on computer engine and also modified the nodes



**Figure 1:** Setting Up the Cluster.

The cluster for the weatherapi is created and it is successfully running.



**Figure 2:** weatherapi-cluster is running.

In the cluster, there are three VM instances in which one master node and other two worker node.



**Figure 3:** VM instances in the cluster.

## Spark Program (Python)

Write Spark program to filter the response data where the daily "feels_like" temperature for the next 5-days is less than 15°C during the "day" time. Exclude the current, minutely, and hourly fields.

Below is the Python program which loads the data from the JSON file and processes the data. It filters the data where the daily "feels_like" temperature for the next 5 days is less than 15°C during the "day" time.

```python
weather_lab5.py  M  ×

Lab5 >  weather_lab5.py > ...
 1    from pyspark.sql import SparkSession
 2    from pyspark.sql.functions import col, expr
 3
 4
 5    spark = SparkSession.builder.appName("FilterWeatherData") \
 6        .config("spark.hadoop.fs.hdfs.impl", "org.apache.hadoop.fs.LocalFileSystem") \
 7        .config("spark.hadoop.mapreduce.framework.name", "local") \
 8        .config("spark.sql.warehouse.dir", "file:/tmp") \
 9        .getOrCreate()
10
11
12
13    weather_data = spark.read.json("weather.json")
14
15
16    filtered_data = weather_data.select(
17        col("lat"),
18        col("lon"),
19        col("timezone"),
20        col("timezone_offset"),
21        col("daily.dt").alias("dt"),
22        col("daily.temp.day").alias("day_temp"),
23        expr("daily.feels_like.day[0]").alias("day_feels_like"),
24        expr("daily.weather[0].description").alias("description")
25    ).filter(col("day_feels_like") < 15)
26
27
28    filtered_data.write.json("fall_weather.json")
29
30    print("Filtered data saved to fall_weather.json")
31    spark.stop()
32
```
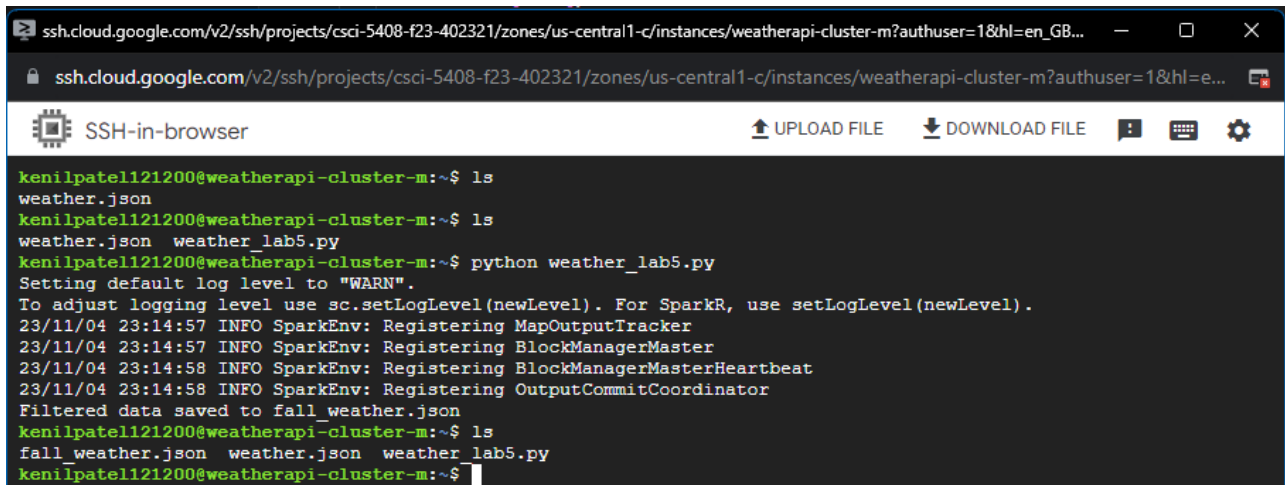
**Figure 4:** Spark Python Code

I have uploaded weather.json and weather_lab5.py file.


**Figure 5:** File Uploading.

I am running the python file and the fall_weather.json file is generated.


**Figure 6:** Running the Python file.

Displaying the fall_weather.json in command prompt.

```
kenilpatel121200@weatherapi-cluster-m:~$ cat fall_weather.json
{
  "lat": 44.6462,
  "lon": -63.5736,
  "timezone": "America/Halifax",
  "timezone_offset": -10800,
  "daily": [
    {
      "dt": 1655827200,
      "sunrise": 1655800138,
      "sunset": 1655856188,
      "moonrise": 1655786400,
      "moonset": 1655829840,
      "moon_phase": 0.75,
      "temp": {
        "day": 12.95,
        "min": 11.01,
        "max": 14.56,
        "night": 12.65,
        "eve": 14.12,
        "morn": 11.99
      },
      "feels_like": {
        "day": 12.68,
        "night": 12.22,
        "eve": 13.68,
        "morn": 11.75
      },
      "pressure": 1021,
      "humidity": 91,
      "dew_point": 11.33,
      "wind_speed": 2.66,
      "wind_deg": 1,
      "wind_gust": 4.9,
      "weather": [
        {
          "id": 500,
          "main": "Rain",
          "description": "light rain",
          "icon": "10d"
        }
      ],
      "clouds": 100,
      "pop": 0.35,
      "rain": 0.73,
      "uvi": 4.11
    }
  ]
}
kenilpatel121200@weatherapi-cluster-m:~$
```

**Figure 7:** Displaying Filtered Data.

Saved the filtered data i.e. fall_weather.json in local computer.



**Figure 8:** Filtered Data in the Local

## References

[1]     Mysql.com. [Online]. Available: https://dev.mysql.com/doc/workbench/en/wb-forward-engineering-live-server.html. [Accessed: 18-Oct-2023].

[2]     "Google Cloud Platform," Google.com. [Online]. Available: https://cloud.google.com/?hl=en. [Accessed: 18-Oct-2023].

[3]     "Apache SparkTM - Unified Engine for large-scale data analytics," Apache.org. [Online]. Available: https://spark.apache.org/. [Accessed: 04-Nov-2023].

[4]     "Python," Python.org. [Online]. Available: https://www.python.org/. [Accessed: 04-Nov-2023].