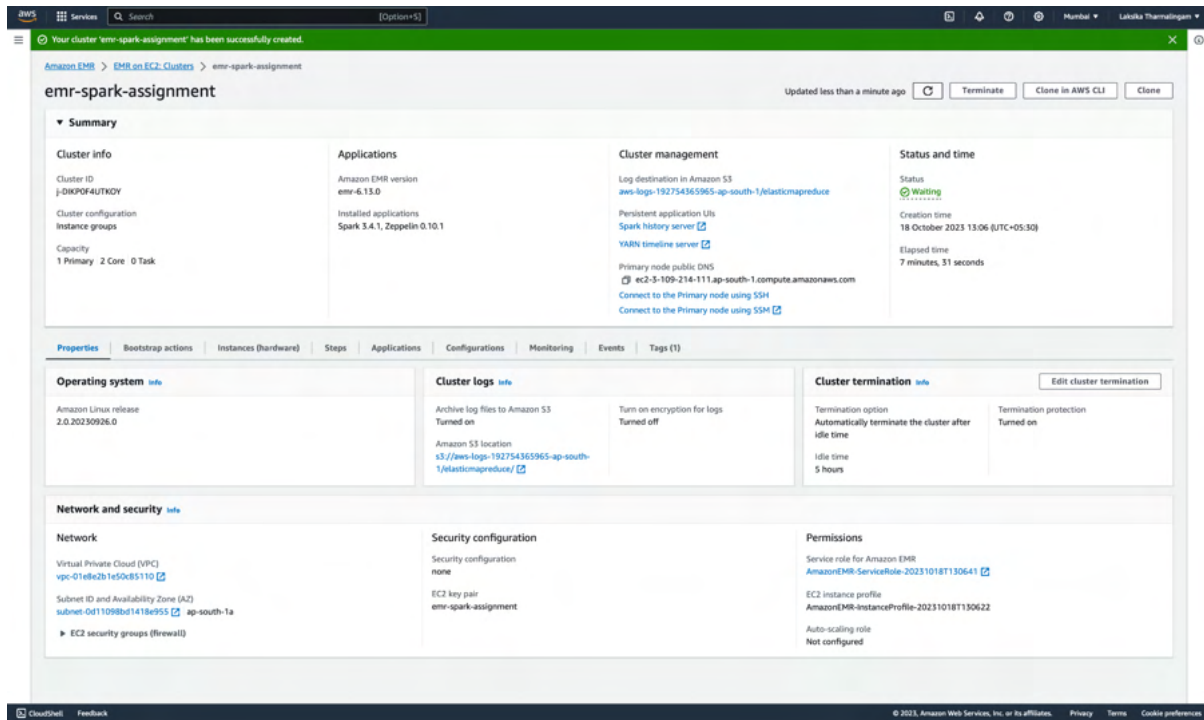# Assignment 01: Big Data Analytics using Hadoop and Spark on Amazon EMR

## 1. Task 1: Amazon EMR Setup



## 2. Task 2: Importing and Managing Data in Amazon S3
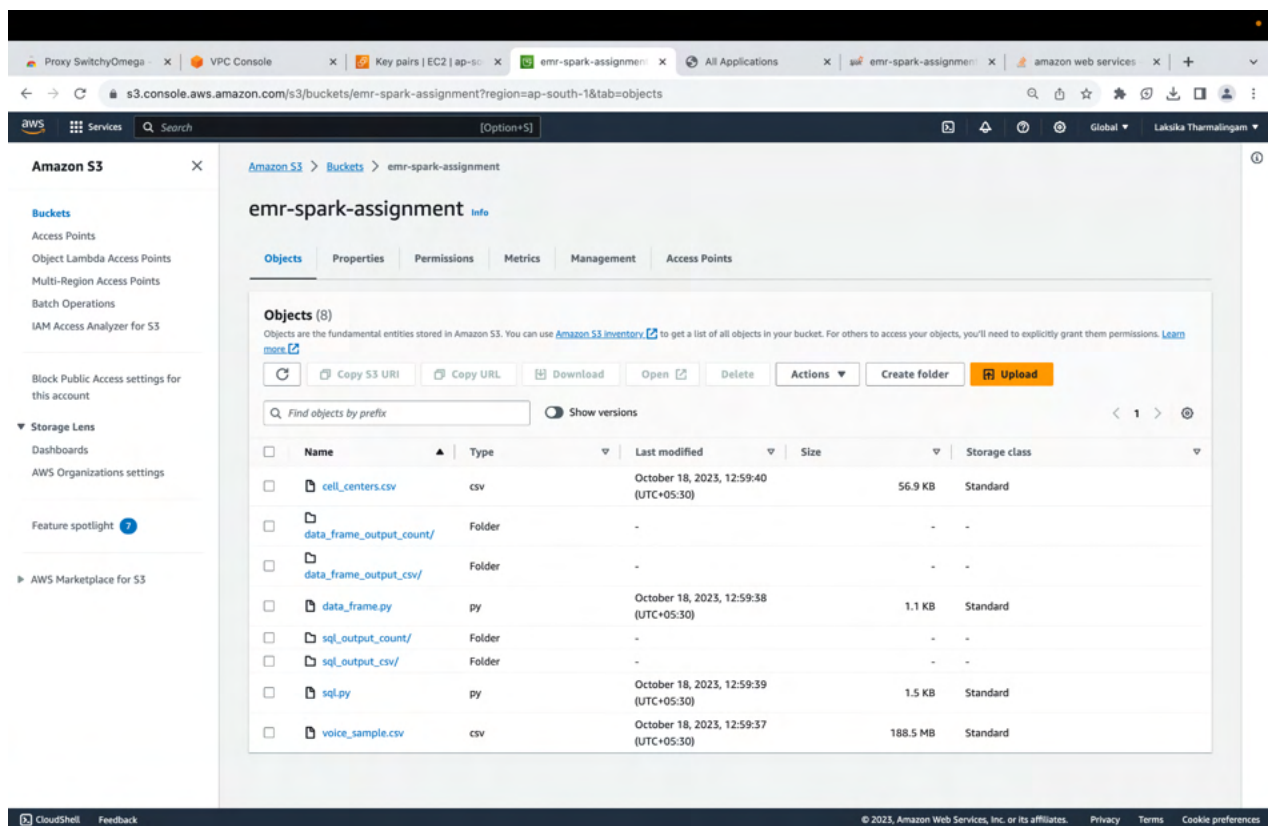
*Commands used to upload :*

```
aws s3 cp /Users/laksika/Desktop/cell_centers.csv s3://emr-spark-assignment

aws s3 cp /Users/laksika/Desktop/voice_sample.csv s3://emr-spark-assignment

aws s3 cp /Users/laksika/Desktop/sql.py s3://emr-spark-assignment

aws s3 cp /Users/laksika/Desktop/data_frame.py s3://emr-spark-assignment
```

```
(base) laksika@Laksikas-MBP Desktop % aws s3 cp /Users/laksika/Desktop/cell_centers.csv s3://emr-spark-assignment
upload: ./cell_centers.csv to s3://emr-spark-assignment/cell_centers.csv
```

*Commands used to list :*

```
aws s3 ls s3://emr-spark-assignment

aws s3 ls s3://emr-spark-assignment --recursive
```

```
[[hadoop@ip-10-0-1-228 ~]$ aws s3 ls s3://emr-spark-assignment/
                           PRE data_frame_output_count/
                           PRE data_frame_output_csv/
                           PRE sql_output_count/
                           PRE sql_output_csv/
2023-10-18 07:29:40      58305 cell_centers.csv
2023-10-18 07:29:38       1157 data_frame.py
2023-10-18 07:29:39       1525 sql.py
2023-10-18 07:29:37  197669861 voice_sample.csv
```



## 3. Task 3: Data Processing with Apache Spark via Amazon EMR

### 3.1. Spark DataFrame Approach

```
(
df1
.join(df2, ['LOCATION_ID'])
.filter(df2['PROVINCE_NAME'] == 'Western')
.groupBy('CALLER_ID').agg(countDistinct('CALL_DATE').alias('COUNT_OF_DISTINCT_CALL_DATE'))
.filter(col('COUNT_OF_DISTINCT_CALL_DATE') == df1.select('CALL_DATE').distinct().count())
.select('CALLER_ID')
)
```

## 3.2. Spark SQL Approach

```
spark.sql(
  '''
  SELECT CALLER_ID
  FROM (
    SELECT CALLER_ID, COUNT(DISTINCT CALL_DATE) AS COUNT_OF_DISTINCT_CALL_DATE
    FROM (
      SELECT CALLER_ID, CALL_DATE
      FROM view_of_voice_sample
      WHERE LOCATION_ID IN (SELECT LOCATION_ID FROM view_of_cell_centers WHERE PROVINCE_NAME = 'Western')
    )
    GROUP BY CALLER_ID
  )
  WHERE COUNT_OF_DISTINCT_CALL_DATE = (SELECT COUNT(DISTINCT CALL_DATE) FROM view_of_voice_sample)
  '''
)
```
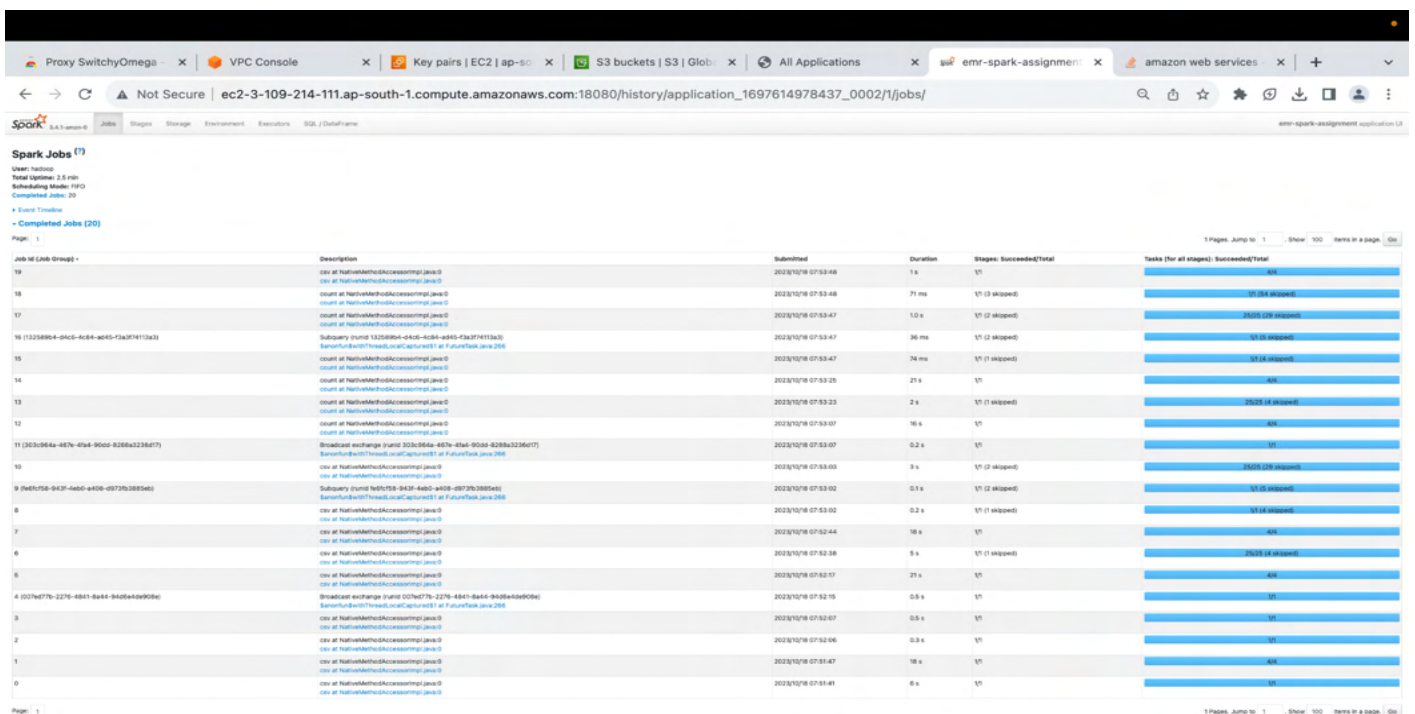
## 3.3. Analysis and Comparison



### 3.3.1. Spark SQL

Proxy SwitchyOmega | VPC Console | Key pairs | EC2 | ap-so | S3 buckets | S3 | Glob | All Applications | emr-spark-assignment | amazon web services | +

Not Secure | ec2-3-109-214-111.ap-south-1.compute.amazonaws.com:18080/history/application_1697614978437_0002/1/stages/

Spark — Jobs Stages Storage Environment Executors SQL / DataFrame — emr-spark-assignment application UI

## Stages for All Jobs

Completed Stages: 20
Skipped Stages: 15

**Completed Stages (20)**

| Stage Id | Description | Submitted | Duration | Tasks: Succeeded/Total | Input | Output | Shuffle Read | Shuffle Write |
|---|---|---|---|---|---|---|---|---|
| 34 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:53:48 | 1 s | 4/4 | | 170 B | | |
| 33 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:53:48 | 51 ms | 1/1 | | | 1290.0 B | |
| 29 | count at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:53:47 | 10 s | 25/25 | | | 11.3 MiB | 1290.0 B |
| 28 | Subquery (runId 132589b4-d4c6-4c84-ad45-f3a3f74113a3) $anonfun$withThreadLocalCaptured$1 at FutureTask.java:266 | 2023/10/18 07:53:47 | 20 ms | 1/1 | | | 50.0 B | |
| 23 | count at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:53:47 | 37 ms | 1/1 | | | 3.0 KiB | 50.0 B |
| 21 | count at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:53:25 | 21 s | 4/4 | 188.5 MiB | | | 3.0 KiB |
| 20 | count at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:53:23 | 2 s | 25/25 | | | 15.5 MiB | 11.3 MiB |
| 18 | count at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:53:07 | 16 s | 4/4 | 188.5 MiB | | | 15.5 MiB |
| 17 | Broadcast exchange (runId 3b3c966a-487e-4fa4-90cd-8288a3236d17) $anonfun$withThreadLocalCaptured$1 at FutureTask.java:266 | 2023/10/18 07:53:07 | 0.2 s | 1/1 | 58.9 KiB | | | |
| 16 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:53:03 | 3 s | 25/25 | 12.4 KiB | | 11.3 MiB | |
| 13 | Subquery (runId fe6fcf58-943f-4eb0-a406-d873fb3885eb) $anonfun$withThreadLocalCaptured$1 at FutureTask.java:266 | 2023/10/18 07:53:02 | 72 ms | 1/1 | | | 50.0 B | |
| 10 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:52:44 | 0.2 s | 1/1 | | | 3.0 KiB | 50.0 B |
| 8 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:52:44 | 18 s | 4/4 | 188.5 MiB | | | 3.0 KiB |
| 7 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:52:38 | 5 s | 25/25 | | | 15.5 MiB | 11.3 MiB |
| 5 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:52:17 | 21 s | 4/4 | 188.5 MiB | | | 15.5 MiB |
| 4 | Broadcast exchange (runId 007ed77b-2276-4841-9a44-94d6e4de908e) $anonfun$withThreadLocalCaptured$1 at FutureTask.java:266 | 2023/10/18 07:52:15 | 0.5 s | 1/1 | 58.9 KiB | | | |
| 3 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:52:07 | 0.5 s | 1/1 | 58.9 KiB | | | |
| 2 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:52:06 | 0.3 s | 1/1 | 7.4 KiB | | | |
| 1 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:51:47 | 18 s | 4/4 | 188.5 MiB | | | |
| 0 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:51:41 | 5 s | 1/1 | 1472.0 B | | | |

**Skipped Stages (15)**

| Stage Id | Description | Submitted | Duration | Tasks: Succeeded/Total | Input | Output | Shuffle Read | Shuffle Write |
|---|---|---|---|---|---|---|---|---|
| 32 | count at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/25 | | | | |
| 31 | count at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/25 | | | | |
| 30 | count at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/4 | | | | |
| 28 | count at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/25 | | | | |
| 27 | count at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/4 | | | | |
| 25 | count at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/1 | | | | |
| 24 | count at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/4 | | | | |
| 22 | count at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/4 | | | | |
| 19 | count at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/4 | | | | |
| 15 | csv at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/25 | | | | |
| 14 | csv at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/4 | | | | |
| 12 | csv at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/1 | | | | |
| 11 | csv at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/4 | | | | |
| 9 | csv at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/4 | | | | |
| 6 | csv at NativeMethodAccessorImpl.java:0 | Unknown | Unknown | 0/4 | | | | |

---

Spark — Jobs Stages Storage Environment Executors SQL / DataFrame — emr-spark-assignment application UI

## SQL / DataFrame

Completed Queries: 8

**Completed Queries (8)**

| ID | Description | Submitted | Duration | Job IDs |
|---|---|---|---|---|
| 7 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:53:48 | 2 s | [19] |
| 6 | count at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:53:07 | 41 s | [17][16][15][14][13][12][11][10] |
| 5 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:53:04 | 52 s | [9][8][6][7][5][4][3][2] |
| 4 | createOrReplaceTempView at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:52:06 | 0.3 s | |
| 3 | createOrReplaceTempView at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:52:06 | 9 ms | |
| 2 | createOrReplaceTempView at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:52:07 | 61 ms | |
| 1 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:52:06 | 0.5 s | [1] |
| 0 | csv at NativeMethodAccessorImpl.java:0 | 2023/10/18 07:51:39 | 8 s | [0] |

---

Spark — Jobs Stages Storage Environment Executors SQL / DataFrame — emr-spark-assignment application UI

## Executors

Show Additional Metrics

### Summary

| | RDD Blocks | Storage Memory | Disk Used | Cores | Active Tasks | Failed Tasks | Complete Tasks | Total Tasks | Task Time (GC Time) | Input | Shuffle Read | Shuffle Write | Excluded |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Active(2) | 0 | 0.0 B / 3.1 GiB | 0.0 B | 4 | 0 | 0 | 134 | 134 | 9.5 min (8 s) | 942.9 MiB | 53.6 MiB | 53.6 MiB | 0 |
| Dead(0) | 0 | 0.0 B / 0.0 B | 0.0 B | 0 | 0 | 0 | 0 | 0 | 0.0 ms (0.0 ms) | 0.0 B | 0.0 B | 0.0 B | 0 |
| Total(2) | 0 | 0.0 B / 3.1 GiB | 0.0 B | 4 | 0 | 0 | 134 | 134 | 9.5 min (8 s) | 942.9 MiB | 53.6 MiB | 53.6 MiB | 0 |

### Executors

Show 20 entries                                                Search:

| Executor ID | Address | Status | RDD Blocks | Storage Memory | Disk Used | Cores | Active Tasks | Failed Tasks | Complete Tasks | Total Tasks | Task Time (GC Time) | Input | Shuffle Read | Shuffle Write | Logs |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| driver | ip-10-0-11-64.ap-south-1.compute.internal:45427 | Active | 0 | 0.0 B / 1 GiB | 0.0 B | 0 | 0 | 0 | 0 | 0 | 2.5 min (0.0 ms) | 0.0 B | 0.0 B | 0.0 B | stdout stderr |
| 1 | ip-10-0-5-249.ap-south-1.compute.internal:39019 | Active | 0 | 0.0 B / 2.1 GiB | 0.0 B | 4 | 0 | 0 | 134 | 134 | 7.0 min (8 s) | 942.9 MiB | 53.6 MiB | 53.6 MiB | stdout stderr |

## 3.3.2. DataFrame API

### 3.3.3. Comparison

### Execution Time



*Spark DataFrame Approach*

The Spark DataFrame approach completed the task in a shorter execution time (2.2 min). It processed data more swiftly, resulting in quicker results.

*Spark SQL Approach*

The Spark SQL approach took slightly longer (2.5 min) to complete the task due to the complexity of some job stages with skipped steps. This resulted in a minor delay in data processing.

In this performance comparison, the Spark DataFrame approach stands out for its swifter and more efficient execution. It completed the task in a shorter time frame, ensuring quicker results. In contrast, the Spark SQL approach, while effective, resulting in a slightly longer execution time.

## Completed and Skipped Stages



Dataframe



SQL

*Completed Stages:*

Spark DataFrame: 17

Spark SQL: 20

*Skipped Stages:*

Spark DataFrame: 12

Spark SQL: 15

While Spark SQL completed more stages, it also encountered more skipped stages compared to Spark DataFrame. This suggests that Spark SQL may involve more complex execution plans or face optimization challenges, which can impact the overall efficiency and execution time of jobs. Spark DataFrame, on the other hand, completed fewer stages but encountered fewer skips, indicating a more streamlined and efficient execution. In this case, the slightly lower number of completed stages in Spark DataFrame and the lower count of skipped stages may indicate a more efficient and straightforward execution, which can be beneficial for specific workloads.

## Executors



Sql



Dataframe

*Spark DataFrame (Executor 1 ):*

Task Shuffle Read Size: 754.3 MiB

Task Shuffle Write Size: 53.6 MiB

Execution Time: 6.2 min

*Spark SQL (Executor 1):*

Task Shuffle Read Size: 942.9 MiB

Task Shuffle Write Size: 53.6 MiB

Execution Time: 7.0 min

Task Shuffle Read Size: The Spark SQL approach had a larger task shuffle read size, indicating a higher volume of data read during shuffling. This might be due to more complex SQL operations or data transformations that require more data exchange between tasks.

Task Shuffle Write Size: Both approaches had similar task shuffle write sizes (53.6 MiB), indicating that they wrote a comparable amount of data during shuffling operations.

Resource Usage and Efficiency: The larger task shuffle read size in Spark SQL suggests a potential inefficiency in data shuffling, which can impact resource usage and overall job performance.

The analysis of task shuffle read and write sizes reveals that the Spark SQL approach had a larger data read volume during shuffling, which may be due to more complex SQL operations. This highlights the importance of optimizing data shuffling processes to ensure efficient resource usage and better overall job performance.

**Completed Queries**



Sql



Dataframe

*Spark DataFrame:*

Completed Queries: 6

Duration of Longest Query: 28 seconds

*Spark SQL:*

Completed Queries: 8

Duration of Longest Query: 52 seconds

Spark DataFrame executed fewer queries (6) and had a shorter longest query run time (28 seconds), which indicates efficient processing. In contrast, Spark SQL executed more queries (8) and had a longer longest query run time (52 seconds), suggesting potential resource-intensive or complex queries. Hence, for the given task and based on the available data, the Spark DataFrame method appears to be more efficient in terms of query execution time and possibly resource utilisation.

## 3.4. Documentation

Step 1: Amazon VPC Setup

To establish a secure network environment for our big data analytics project, I initiated by accessing the Amazon VPC console. Within the console, I created a new Virtual Private Cloud (VPC) thereby ensuring the network configuration met our specifications.

Step 2: Key Pair Generation

To facilitate secure access to our Amazon EMR cluster instances, I generated an SSH key pair. This key pair, created through the EC2 service, is of critical importance for maintaining the integrity and confidentiality of our cluster.

Step 3: Amazon S3 Configuration

Data management is pivotal in big data analytics. Therefore, I set up an Amazon S3 bucket, leveraging its scalability for data storage and retrieval. As an integral part of this step, I uploaded the required data files to our designated S3 bucket, ensuring data availability for our EMR cluster.

Step 4: Amazon EMR Cluster Creation

The heart of our big data analytics operations, the Amazon EMR cluster, was configured through the Amazon EMR console. In this step, I determined cluster details, including its name, EC2 instance types(m4.large), number of instances, S3 access and the designated key pair for SSH access.

Step 5: Python Script Development

The foundation for data processing lay in the development of Python scripts for both the DataFrame and SQL methods. These scripts were crafted to execute data processing tasks with the Spark framework. The finalised Python scripts were uploaded to a directory within the Amazon S3 bucket, making them accessible to the EMR cluster for execution.

Step 6: EC2 Security Configuration

Security enhancements were integrated by modifying the security group settings for the Amazon EMR cluster instances. In particular, SSH access was explicitly granted to the cluster's EC2 instances, aligning security requirements with operational needs.

Step 7: Execution of Python Scripts

Our Python scripts, designed for data processing, were executed on the Amazon EMR cluster. I ran them as Spark applications by selecting the scripts directly from Amazon S3 through the EMR cluster's execution steps

Output: