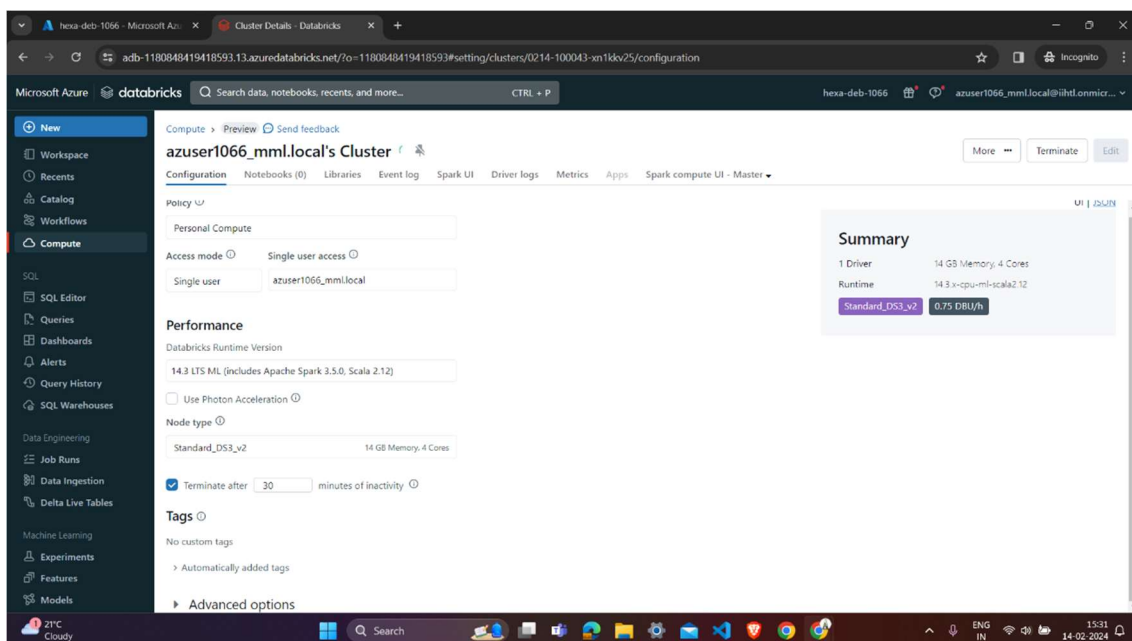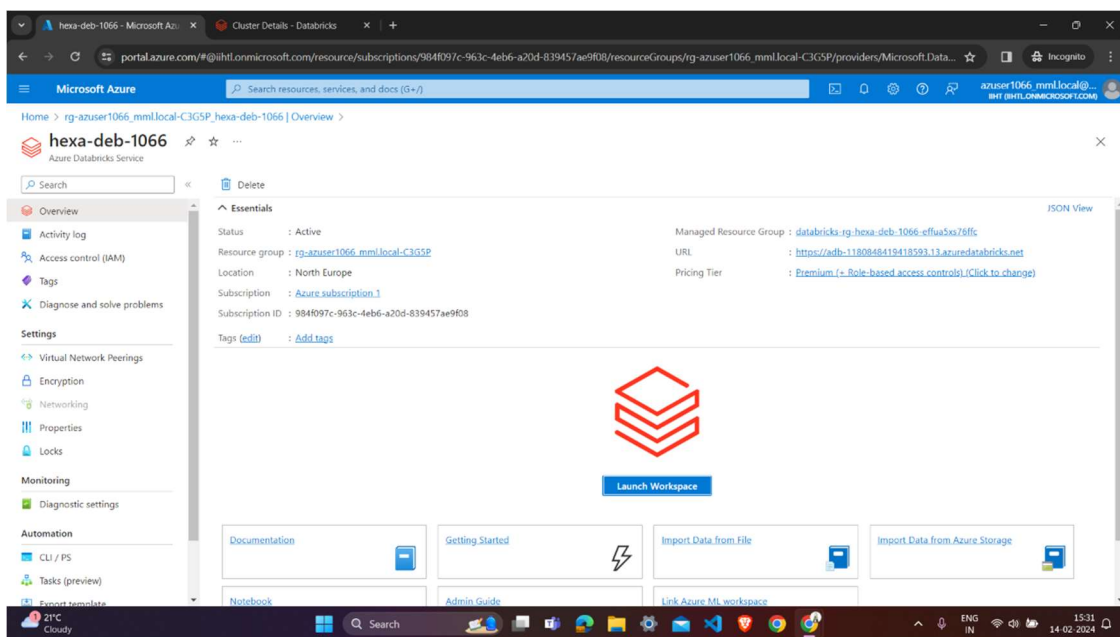**Name: Abhishek Kanoujia**

**DATA ENGINEERING BATCH 1**

**DAY 19 ASSIGNMENT**

**Introduction to Databricks Delta Lake:-** first create cluster then perform the query

# Read table:-



# Update table data:-

# Conditional Update without overwrte:-

Cmd 6

```python
1   # Update every even value by adding 100 to it
2   deltaTable.update(
3     condition = expr("id % 2 == 0"),
4     set = { "id": expr("id + 100") })
```

▶ (13) Spark Jobs

Command took 10.77 seconds -- by azuser1066_mml.local@iihtl.onmicrosoft.com at 2/14/2024, 3:39:33 PM on azuser1066_mml.local's Cluster

Cmd 7

```python
1   df = spark.read.format("delta").load("/tmp/delta-table")
2   df.show()
```

▶ (2) Spark Jobs

▶ ▦ df: pyspark.sql.dataframe.DataFrame = [id: integer]

```
+---+
| id|
+---+
|106|
|108|
|  5|
|  7|
|  9|
+---+
```

Command took 1.09 seconds -- by azuser1066_mml.local@iihtl.onmicrosoft.com at 2/14/2024, 3:40:14 PM on azuser1066_mml.local's Cluster

Cmd 8

---

Command took 1.09 seconds -- by azuser1066_mml.local@iihtl.onmicrosoft.com at 2/14/2024, 3:40:14 PM on azuser1066_mml.local's Cluster

Cmd 8

```python
1   # Delete every even value
2   deltaTable.delete(condition = expr("id % 2 == 0"))
```

▶ (10) Spark Jobs

Command took 4.91 seconds -- by azuser1066_mml.local@iihtl.onmicrosoft.com at 2/14/2024, 3:40:48 PM on azuser1066_mml.local's Cluster

Cmd 9

```python
1   df = spark.read.format("delta").load("/tmp/delta-table")
2   df.show()
```

▶ (1) Spark Jobs

▶ ▦ df: pyspark.sql.dataframe.DataFrame = [id: integer]

```
+---+
| id|
+---+
|  5|
|  7|
|  9|
+---+
```

Command took 0.39 seconds -- by azuser1066_mml.local@iihtl.onmicrosoft.com at 2/14/2024, 3:41:07 PM on azuser1066_mml.local's Cluster

Cmd 10

```python
1   # Upsert (merge) new data
```

**Untitled Notebook 2024-02-14 15:36:37** Python ∨ ☆
File Edit View Run Help  Last edit was 5 minutes ago  New cell UI: OFF ∨
Run all ● azuser1066_mml.local's... ∨  Schedule  Share

```python
1  # Upsert (merge) new data
2  newData = spark.range(0, 20)
```

▶ ▦ newData: pyspark.sql.dataframe.DataFrame = [id: long]

Command took 0.24 seconds -- by azuser1066_mml.local@iihtl.onmicrosoft.com at 2/14/2024, 3:41:21 PM on azuser1066_mml.local's Cluster

Cmd 11

```python
1  deltaTable.alias("oldData") \
2    .merge(
3      newData.alias("newData"),
4      "oldData.id = newData.id") \
5    .whenMatchedUpdate(set = { "id": col("newData.id") }) \
6    .whenNotMatchedInsert(values = { "id": col("newData.id") }) \
7    .execute()
8
9  deltaTable.toDF().show()
```

▶ (12) Spark Jobs
```
| 1|
| 2|
| 3|
| 4|
| 5|
| 6|
| 7|
| 8|
```

21°C Cloudy | Q Search | ENG IN | 15:42 14-02-2024

---

Browser window 2:

```python
8
9  deltaTable.toDF().show()
```

▶ (12) Spark Jobs
```
| 1|
| 2|
| 3|
| 4|
| 5|
| 6|
| 7|
| 8|
| 9|
| 10|
| 11|
| 12|
| 13|
| 14|
| 15|
| 16|
| 17|
| 18|
| 19|
+---+
```

Command took 8.38 seconds -- by azuser1066_mml.local@iihtl.onmicrosoft.com at 2/14/2024, 3:41:44 PM on azuser1066_mml.local's Cluster

Cmd 12

Python

21°C Cloudy | Q Search | ENG IN | 15:42 14-02-2024

**Databricks Delta Lake**   Python ∨   ☆

File   Edit   View   Run   Help   Last edit was 1 minute ago   New cell UI: OFF ∨    ▶ Run all   ● Abhishek2 ∨   Share   Publish

Cmd 1

```sql
1  %sql
2  CREATE TABLE delta.`/tmp/deltain-table` USING DELTA AS SELECT col1 as id FROM VALUES 0,1,2,3,4;
```

▸ (6) Spark Jobs

▸ ▦ _sqldf: pyspark.sql.dataframe.DataFrame = [num_affected_rows: long, num_inserted_rows: long]

Query returned no results

Command took 4.10 seconds -- by abhishek7621cse@gmail.com at 2/14/2024, 5:46:34 PM on Abhishek2

Cmd 2

SQL ▶ ∨ ⊞ ∨ — ✕

```sql
1  %sql
2  SELECT * FROM delta.`/tmp/deltain-table`;
```

▸ (2) Spark Jobs

▸ ▦ _sqldf: pyspark.sql.dataframe.DataFrame = [id: integer]

Table ∨   +

|     | id ▲ |
| --- | --- |
| 1   | 0   |
| 2   | 1   |
| 3   | 2   |
| 4   | 3   |
| 5   | 4   |

⤓ 5 rows   |   1.30 seconds runtime      Refreshed now

Command took 1.30 seconds -- by abhishek7621cse@gmail.com at 2/14/2024, 5:46:43 PM on Abhishek2

Cmd 3

```sql
1  %sql
2  INSERT OVERWRITE delta.`/tmp/deltain-table` SELECT col1 as id FROM VALUES 5,6,7,8,9;
```

▸ (6) Spark Jobs

▸ ▦ _sqldf: pyspark.sql.dataframe.DataFrame = [num_affected_rows: long, num_inserted_rows: long]

Table ∨   +

|     | num_affected_rows ▲ | num_inserted_rows ▲ |
| --- | --- | --- |
| 1   | 5   | 5   |

⤓ 1 row   |   2.52 seconds runtime      Refreshed now

ⓘ SQL cell result stored as PySpark data frame _sqldf . Learn more

Command took 2.52 seconds -- by abhishek7621cse@gmail.com at 2/14/2024, 5:46:48 PM on Abhishek2

Cmd 4

Python ▶ ∨ — ✕

```python
1  from delta.tables import *
2  from pyspark.sql.functions import *
3
4
```

Command took 0.07 seconds -- by abhishek7621cse@gmail.com at 2/14/2024, 5:46:04 PM on Abhishek2

Cmd 5

```python
1   eltaTable = DeltaTable.forPath(spark, "/tmp/delta-table")
2
3   # Update every even value by adding 100 to it
4   deltaTable.update(
5     condition = expr("id % 2 == 0"),
6     set = { "id": expr("id + 100") })
7
8   df = spark.read.format("delta").load("/tmp/delta-table")
9   df.show()
10
11
```

▸ (10) Spark Jobs

▸ ▦ df: pyspark.sql.dataframe.DataFrame = [id: integer]

```
+---+
| id|
+---+
|  7|
|  9|
| 17|
| 19|
|  1|
|  3|
|  5|
| 11|
| 13|
| 15|
+---+
```

Command took 8.59 seconds -- by abhishek7621cse@gmail.com at 2/14/2024, 5:45:41 PM on Abhishek2

Cmd 6

```
1
2    # Delete every even value
3    deltaTable.delete(condition = expr("id % 2 == 0"))
4
5    df = spark.read.format("delta").load("/tmp/delta-table")
6    df.show()
7
8
```

▸ (10) Spark Jobs

▸ 🔲 df: pyspark.sql.dataframe.DataFrame = [id: integer]

```
+---+
| id|
+---+
|  7|
|  9|
| 17|
| 19|
|  1|
|  3|
|  5|
| 11|
| 13|
| 15|
+---+
```

Command took 8.26 seconds -- by abhishek7621cse@gmail.com at 2/14/2024, 5:45:48 PM on Abhishek2

Cmd 7

```
1
2    # Upsert (merge) new data
3    newData = spark.range(0, 20)
4    df = spark.read.format("delta").load("/tmp/delta-table")
5    df.show()
6
7
```

▸ (3) Spark Jobs

▸ 🔲 newData: pyspark.sql.dataframe.DataFrame = [id: long]

▸ 🔲 df: pyspark.sql.dataframe.DataFrame = [id: integer]

```
+---+
| id|
+---+
|  7|
|  9|
| 17|
| 19|
|  1|
|  3|
|  5|
| 11|
| 13|
| 15|
+---+
```

Command took 1.59 seconds -- by abhishek7621cse@gmail.com at 2/14/2024, 5:45:50 PM on Abhishek2

Cmd 8

```
1
2    deltaTable.alias("oldData") \
3      .merge(
4        newData.alias("newData"),
5        "oldData.id = newData.id") \
6      .whenMatchedUpdate(set = { "id": col("newData.id") }) \
7      .whenNotMatchedInsert(values = { "id": col("newData.id") }) \
8      .execute()
9
10   deltaTable.toDF().show()
```

▸ (14) Spark Jobs

```
|  3|
|  4|
|  7|
|  8|
|  9|
| 17|
| 18|
| 19|
| 12|
| 13|
| 14|
|  5|
|  6|
| 15|
| 16|
|  0|
|  1|
| 10|
| 11|
+---+
```

Command took 12.58 seconds -- by abhishek7621cse@gmail.com at 2/14/2024, 5:45:52 PM on Abhishek2

Cmd 9

```
1
```