**NAME : AKULA SHARATH CHANDRA**
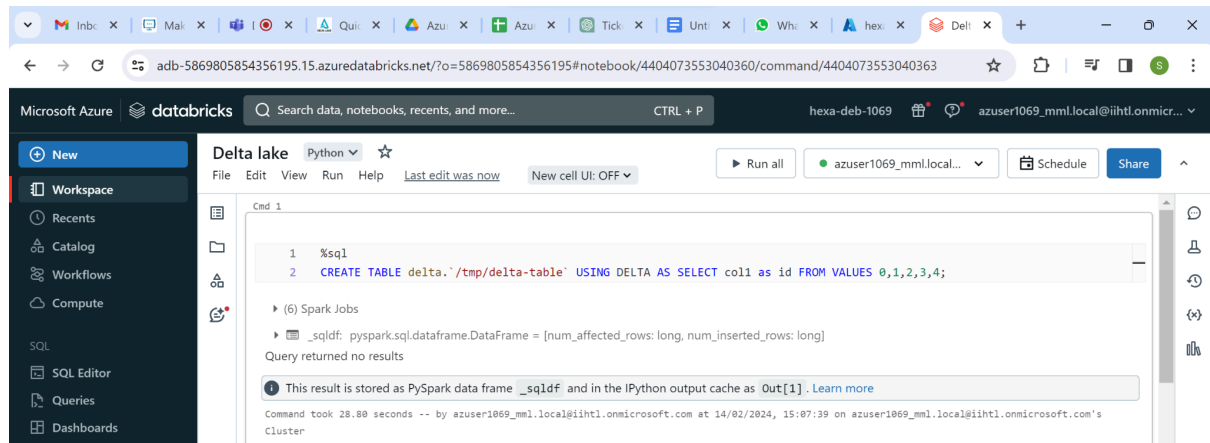
**BATCH : DATA ENGINEERING**
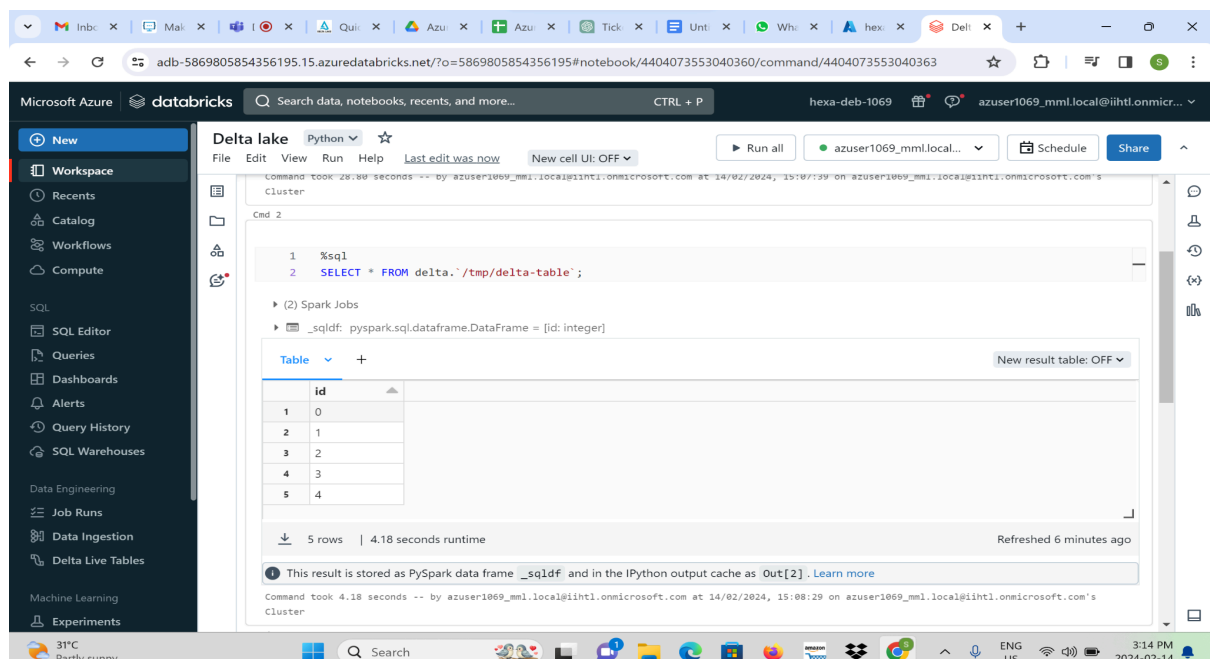
**DATE : 14-02-2024**

**TOPICS :**

**1. Creating a delta table and passing values**

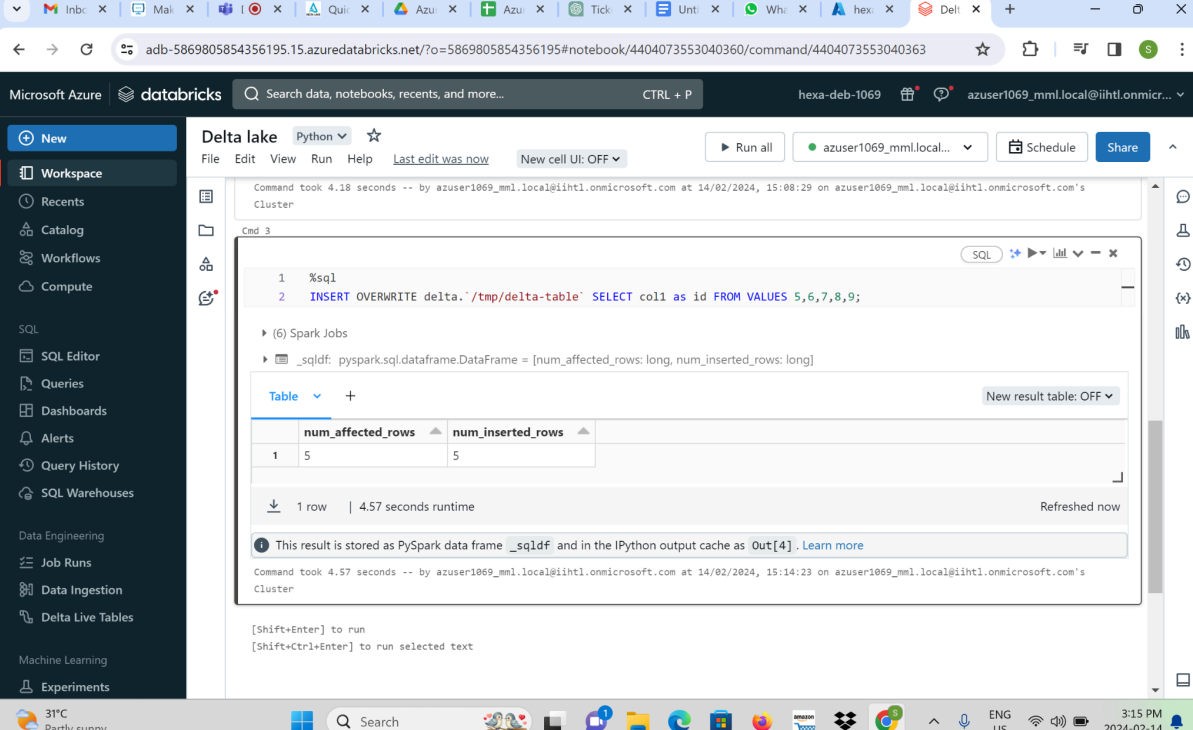**2. Creating and visualising the stream data Create a table.**

To create a Delta table, write a DataFrame out in the delta format. You can use existing Spark SQL code and change the format from parquet, csv, json, and so on, to delta.



Read data You read data in your Delta table by specifying the path to the files: "/tmp/delta-table":

Update table data Delta Lake supports several operations to modify tables using standard DataFrame APIs. This example runs a batch job to overwrite the data in the table:



Conditional update without overwrite Delta Lake provides programmatic APIs to conditional update, delete, and merge (upsert) data into tables. Here are a few examples.

**Delta lake** Python ☆

File  Edit  View  Run  Help    Last edit was 1 minute a...    New cell UI: OFF

Interrupt | ● azuser1069_mml.local... | Schedule | Share

Cmd 4

```sql
1  %sql
2  -- Update every even value by adding 100 to it
3  UPDATE delta.`/tmp/delta-table` SET id = id + 100 WHERE id % 2 == 0;
```

▶ (13) Spark Jobs

▶ ▦ _sqldf: pyspark.sql.dataframe.DataFrame = [num_affected_rows: long]

Table ∨  +                                          New result table: OFF ∨

| | num_affected_rows |
|---|---|
| 1 | 2 |

⬇ 1 row | 9.71 seconds runtime                    Refreshed 5 minutes ago

ⓘ This result is stored as PySpark data frame _sqldf and in the IPython output cache as Out[5] . Learn more

Command took 9.71 seconds -- by azuser1069_mml.local@iihtl.onmicrosoft.com at 14/02/2024, 15:31:36 on azuser1069_mml.local@iihtl.onmicrosoft.com's Cluster

Cmd 5

```sql
1  %sql
2  -- Upsert (merge) new data
3  CREATE TEMP VIEW newData AS SELECT col1 AS id FROM VALUES 1,3,5,7,9,11,13,15,17,19;
```

---

Cmd 5

```sql
1  %sql
2  -- Upsert (merge) new data
3  CREATE TEMP VIEW newData AS SELECT col1 AS id FROM VALUES 1,3,5,7,9,11,13,15,17,19;
```

OK

Command took 0.13 seconds -- by azuser1069_mml.local@iihtl.onmicrosoft.com at 14/02/2024, 15:32:43 on azuser1069_mml.local@iihtl.onmicrosoft.com's Cluster

You can query previous snapshots of your Delta table by using time travel. If you want to access the data that you overwrote, you can query a snapshot of the table before you overwrote the first set of data using the versionAsOf option.

Write a stream of data to a table You can also write to a Delta table using Structured Streaming. The Delta Lake transaction log guarantees exactly-once processing, even when there are other streams or batch queries running concurrently against the table. By default, streams run in append mode, which adds new records to the table:

Read a stream of changes from a table While the stream is writing to the Delta table, you can also read from that table as streaming source. For example, you can start another streaming query that prints all the changes made to the Delta table. You can specify which version Structured Streaming should start from by providing the startingVersion or startingTimestamp option to get changes from that point onwards. See Structured Streaming for details.

```
stream2 = spark.readStream.format("delta").load("/tmp/deltain-table").writeStream.format("console").start
()
```

f0345a3d-3e3d-4d14-8ddf-a387dab47765    *Last updated: 15 hours ago*

**Dashboard**    Raw Data

| Input vs. Processing Rate records per second | 0 rec/s Input rate | 3.2 rec/s Processing rate | Batch Duration in milliseconds | 895 ms Average | 3100 ms Latest |