# Azure Databricks Assignment - 3
# Mitushi Vishwakarma

## Introduction to Databricks Delta Lake

Week - 4      Azure Databricks
               Day - 3

Delta format → at the centre of
data lake paradigm known as Delta Lake.

- Delta Lake is an open storage layer that brings reliability to data lakes.

- provides scalable metadata handling, unifies streaming, batch data processing

- Delta Lake runs on top of your existing data lake

- fully compatible with Apache Spark APIs.

Data lake is tool / place that stores, processes large amounts of data

Delta Lake is storage layer b/w data lake and databricks cluster

Delta Table (Program) → [Spark API] Databricks Cluster → Delta Lake Storage Layer → Data Lake

Ex. Azure data lake

How to Create a Delta Table.

Step 1 : Uploading data to DBFS
2 : data in delta format.

Create a Table :-
Write dataframe in delta format.

Read a Table :-
read data by specifying the path.
SQL    delta `/tmp/delta-table`;
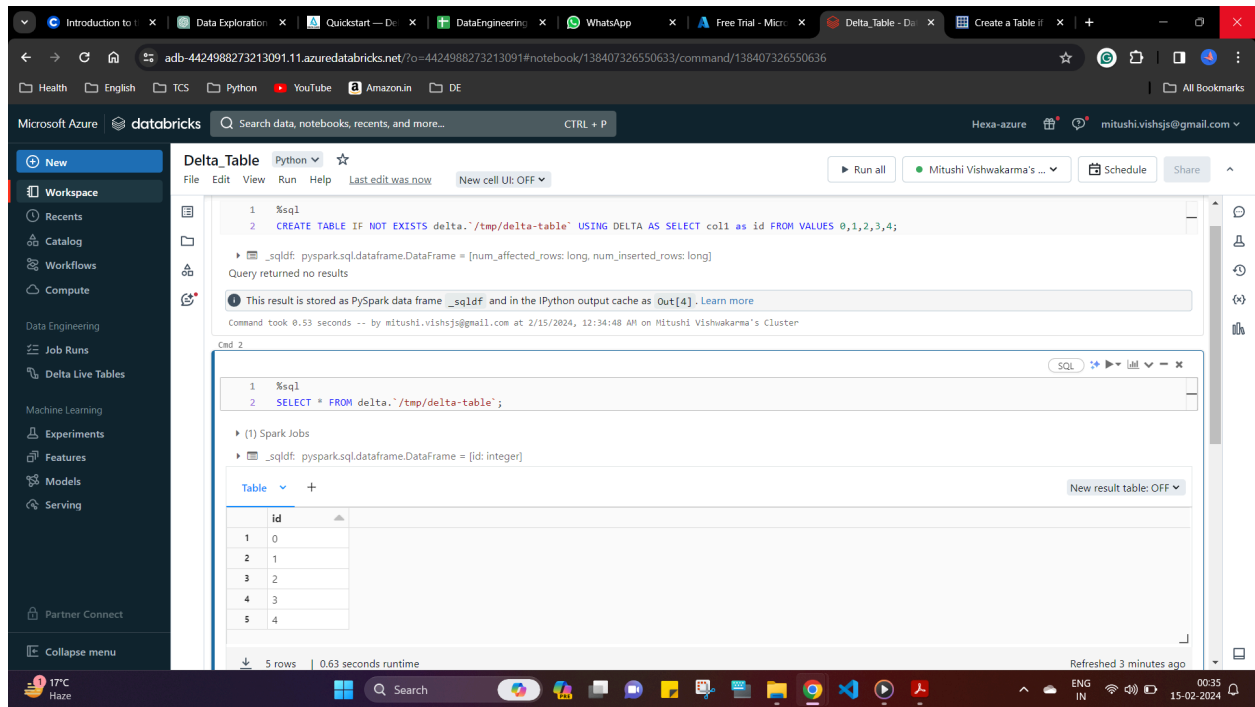   python   .format( ) . load( )

Update
   mode = overwrite

   .format() . mode() . Save( )


Reading older version.
   .option( ) . load(i

# Practice :
## Creating delta table and reading data in sql



## Creating delta table and reading data in python

# Updating data in delta table in python



```python
1   data = spark.range(5, 10)
2   data.write.format("delta").mode("overwrite").save("/tmp/delta-table1")
```

▶ (4) Spark Jobs

Command took 3.24 seconds -- by mitushi.vishsjs@gmail.com at 2/15/2024, 12:38:18 AM on Mitushi Vishwakarma's Cluster

Cmd 6

```python
1   data.show()
```

▶ (1) Spark Jobs

```
+---+
| id|
+---+
|  5|
|  6|
|  7|
|  8|
|  9|
+---+
```

Command took 0.24 seconds -- by mitushi.vishsjs@gmail.com at 2/15/2024, 12:38:24 AM on Mitushi Vishwakarma's Cluster

Cmd 7

```
1   Start typing or generate with AI (Ctrl + I)...
```

# Read older versions of data using time travel



```
|  9|
+---+
```

Command took 0.24 seconds -- by mitushi.vishsjs@gmail.com at 2/15/2024, 12:38:24 AM on Mitushi Vishwakarma's Cluster
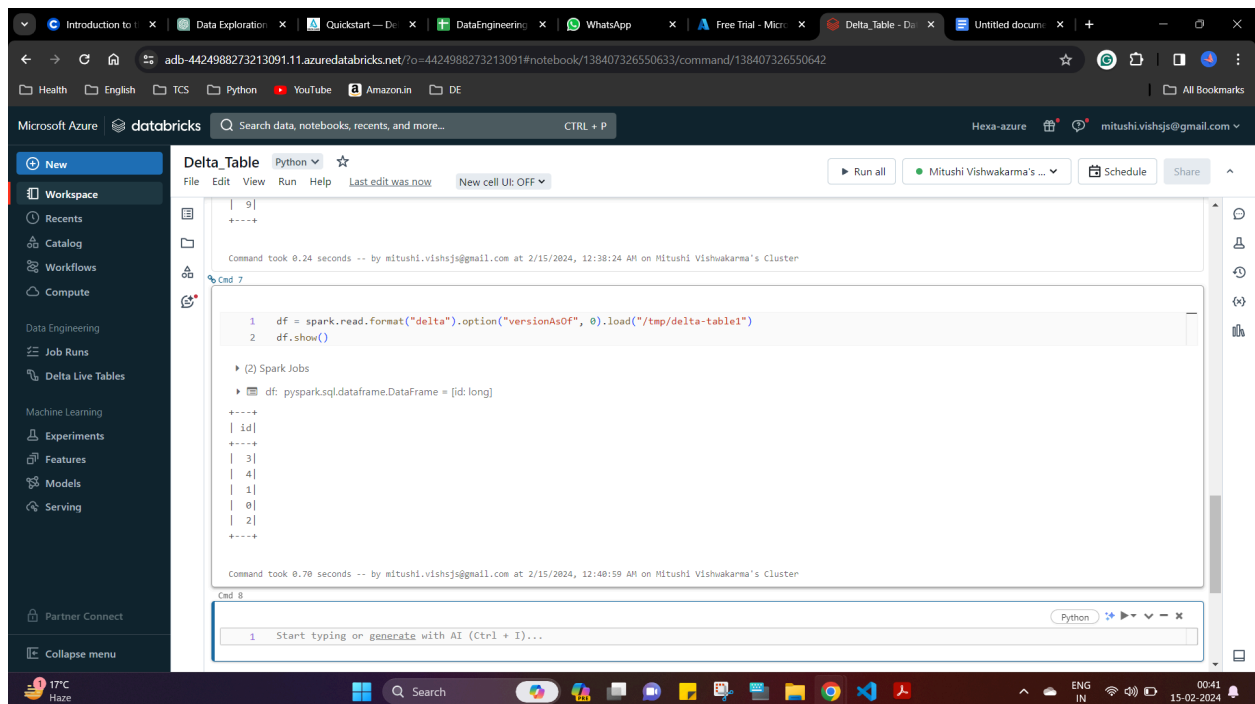
Cmd 7

```python
1   df = spark.read.format("delta").option("versionAsOf", 0).load("/tmp/delta-table1")
2   df.show()
```

▶ (2) Spark Jobs

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

```
+---+
| id|
+---+
|  3|
|  4|
|  1|
|  0|
|  2|
+---+
```

Command took 0.70 seconds -- by mitushi.vishsjs@gmail.com at 2/15/2024, 12:40:59 AM on Mitushi Vishwakarma's Cluster

Cmd 8

```
1   Start typing or generate with AI (Ctrl + I)...
```

# Write a stream of data to a table