Name: Krushnakumar Patle

Email: krishnapatle128@gmail.com

Batch: Data Engineering Batch-1

**Delta Lake:**

Delta Lake is the optimized storage layer that provides the foundation for storing data and tables in the Databricks lakehouse. Delta Lake is open source software that extends Parquet data files with a file-based transaction log for ACID transactions and scalable metadata handling. Delta Lake is fully compatible with Apache Spark APIs, and was developed for tight integration with Structured Streaming, allowing you to easily use a single copy of data for both batch and streaming operations and providing incremental processing at scale.

To set up a Python project (for example, for unit testing), you can install Delta Lake using pip install delta-spark==3.1.0 and then configure the SparkSession with the configure_spark_with_delta_pip() utility function in Delta Lake.
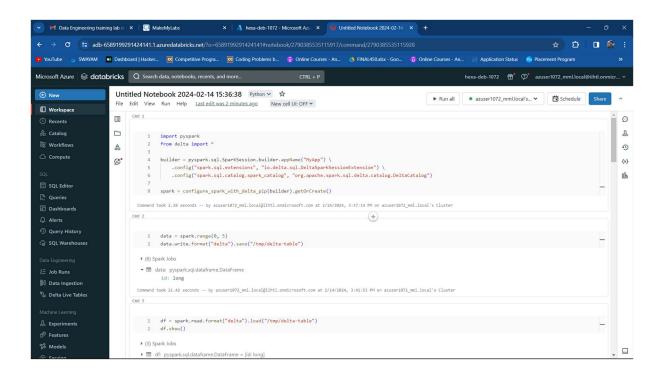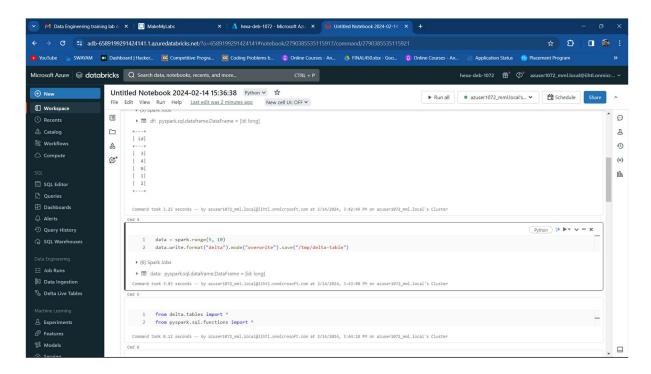
Create a table:

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.
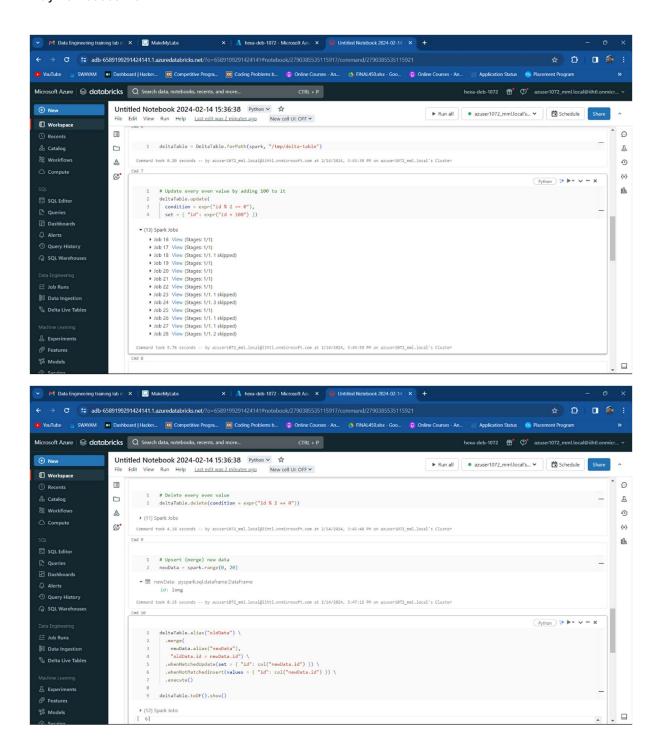
Read data:

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

# Day 18 Assessment

# Day 18 Assessment

# Day 18 Assessment