# Reading Different File Formats from ADLS Gen2 into Delta Tables using Azure Databricks

☐ Setup Azure Resources

- Created a Resource Group in Azure Portal to logically group all resources.

- Created an Azure Data Lake Storage Gen2 (ADLS Gen2) account to store customer data files.



- Created an Azure Databricks workspace to run Spark clusters and process data.

- Configured Linked Services in Azure Databricks to securely connect to the ADLS Gen2 storage.

- Install the required Libraries



- ☐ Upload Sample Customer Data

  - Generated random customer data with 500 records in



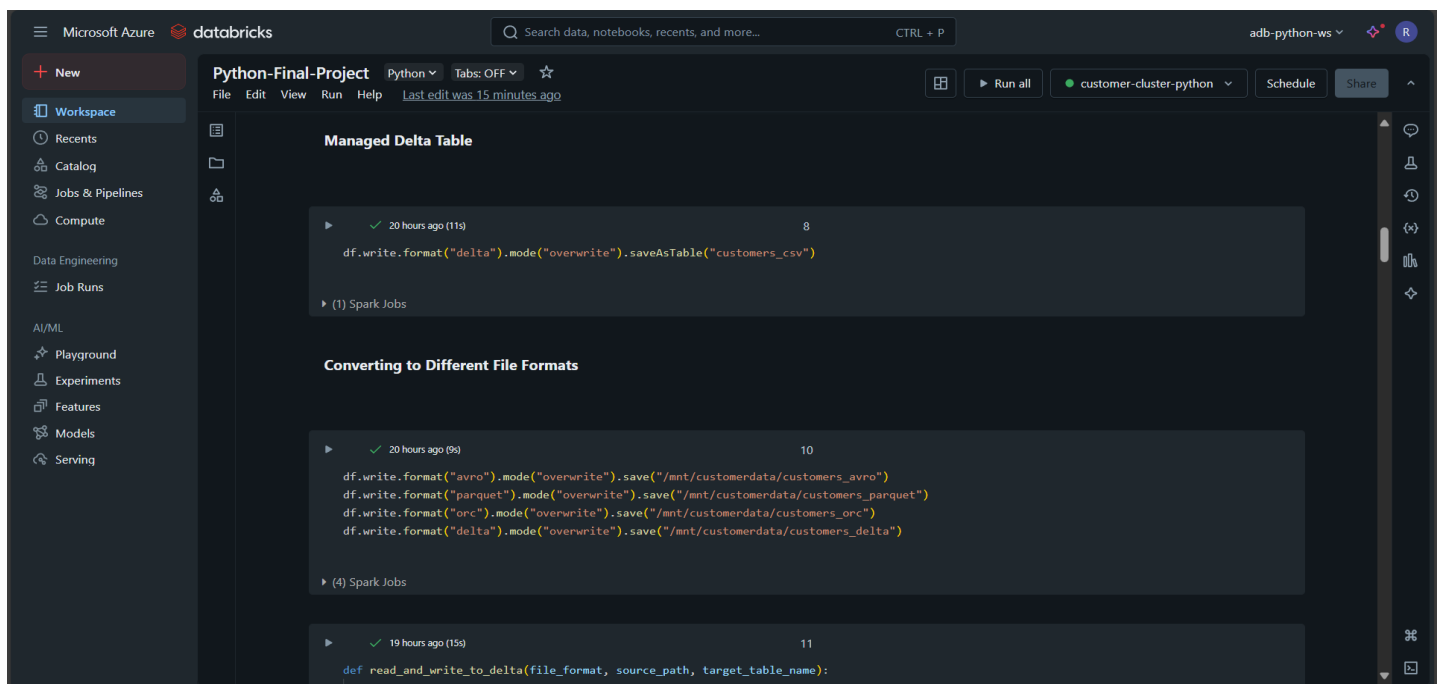- ☐ Connecting Databricks to ADLS Gen2
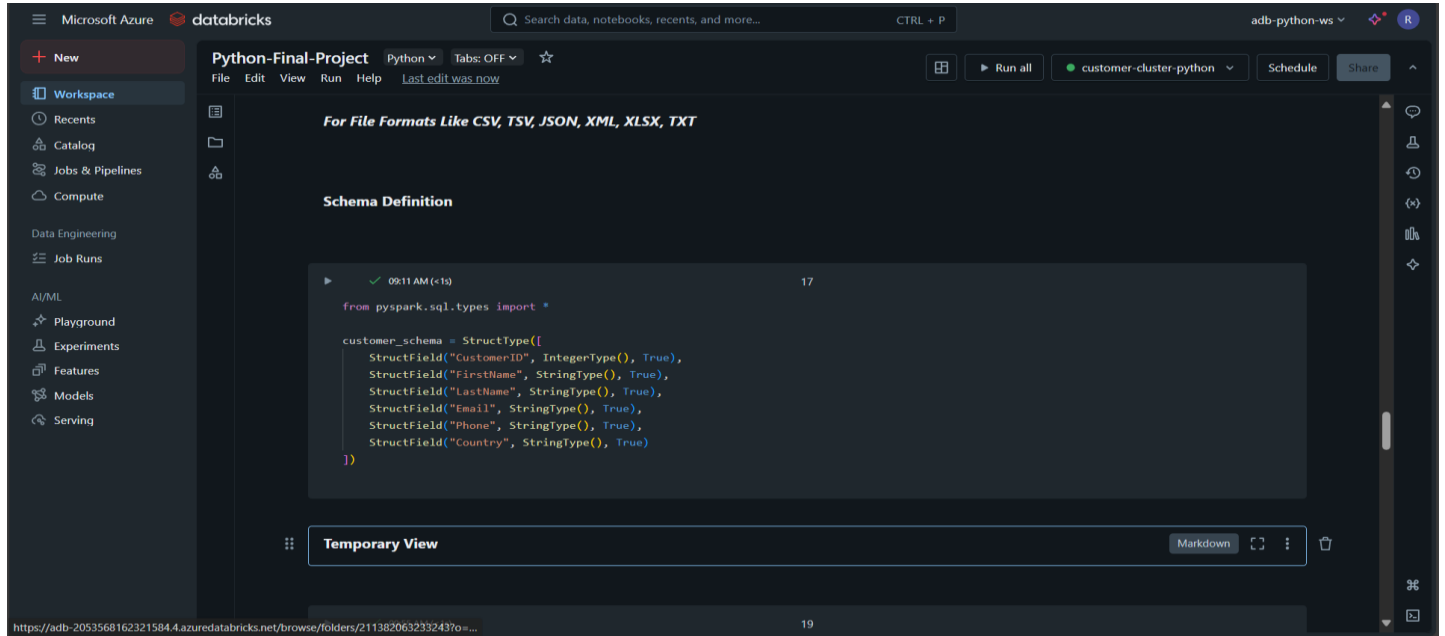
- Mount ADLS Gen2 storage in Databricks.



- Convert the sample data into various file formats including CSV, TSV, JSON, XML, XLSX, TXT, Avro, Parquet, ORC, and Delta.

☐ Creating Managed Delta Tables for Various File Formats

- For Avro, Parquet, ORC, and Delta file formats, directly created managed Delta tables using the CTAS (Create Table As Select) SQL command to read from files and save as Delta tables.

- For CSV, TSV, JSON, XML, XLSX, and TXT file formats:
    - Read files using Spark with this schema and created temporary views.
    - Created managed Delta tables from these temporary views using CTAS, ensuring consistent structure.



☐ Data Validation and Code Optimization

- Implemented a generic Python function to automate reading, view creation, Delta table creation, and validation of row counts for each file format.

- Verified that each Delta table contains exactly 500 customer records to ensure data consistency and completeness.