# HOME PRIJECT
# AZURE DATABRICKS DATA ENGINEERING PROJECT REPORT

OCTOBER 20, 2025

RUSHIRAJ PATHAK

# Contents

# 1. Introduction

This project demonstrates the creation of an end-to-end cloud-based data pipeline using **Microsoft Azure Databricks** integrated with **Azure Data Lake Storage Gen2**. The goal is to establish a secure data flow from ingestion to transformation and storage, while performing scalable data analytics using **PySpark** and **Delta Lake**.

The project walks through all critical stages of setup and implementation, including Azure resource provisioning, Databricks workspace configuration, data transformation using Spark, and Delta table creation for structured querying.

# 2. Objectives

- The main objectives of this project are:
- To create and configure a secure Azure environment for data engineering.
- To connect **Azure Databricks** with **Azure Data Lake Storage** using secure access keys.
- To ingest and clean raw data using **PySpark**.
- To transform, aggregate, and store data in **Delta format**.
- To register and query the processed data as a **Hive Metastore table**.

# 3. Tools and Technologies

| Tool / Service | Purpose |
|---|---|
| Microsoft Azure Portal | Resource and billing management |
| Azure Storage Account (ADLS Gen2) | Cloud-based data storage |
| Azure Databricks | Data transformation and analysis using PySpark |
| PySpark | Distributed data processing framework |
| Delta Lake | Reliable and ACID-compliant data storage layer |
| Hive Metastore | Metadata management for SQL querying |

## 4. Implementation Steps

### Step 1 – Azure Setup

- Created a free Azure account with CA**$300 free credit**.
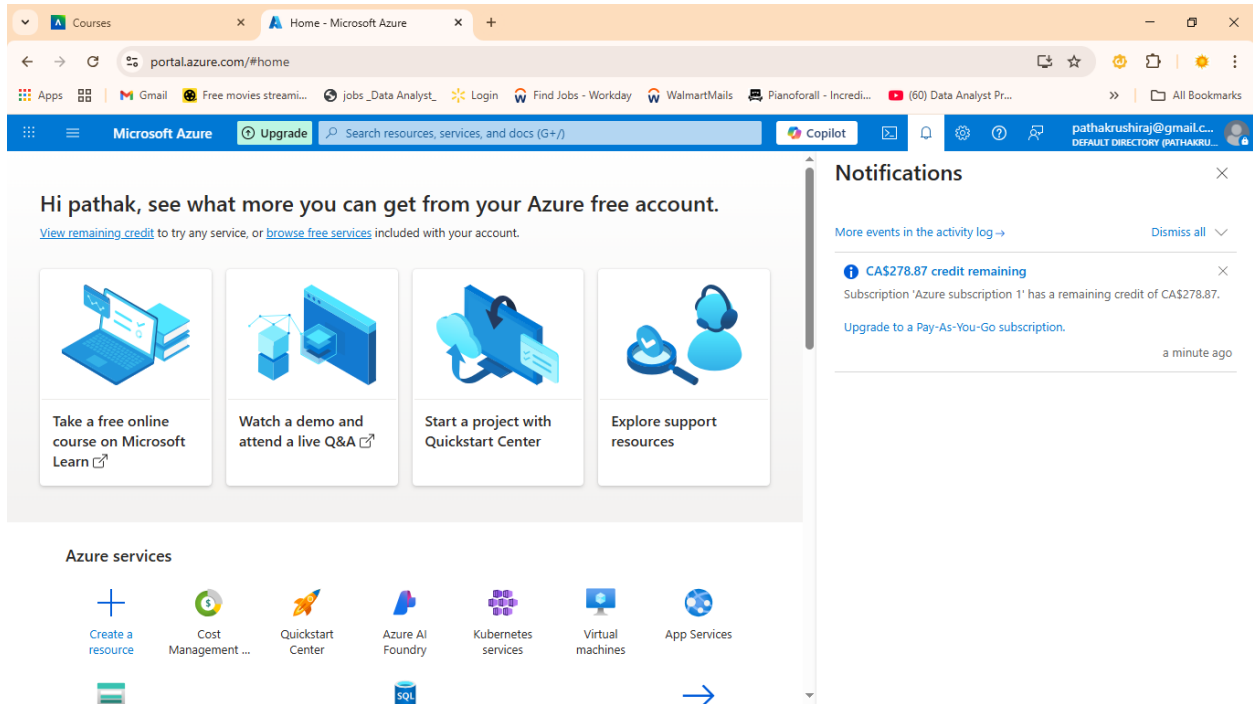- Configured **billing profiles** and verified **active subscription**.

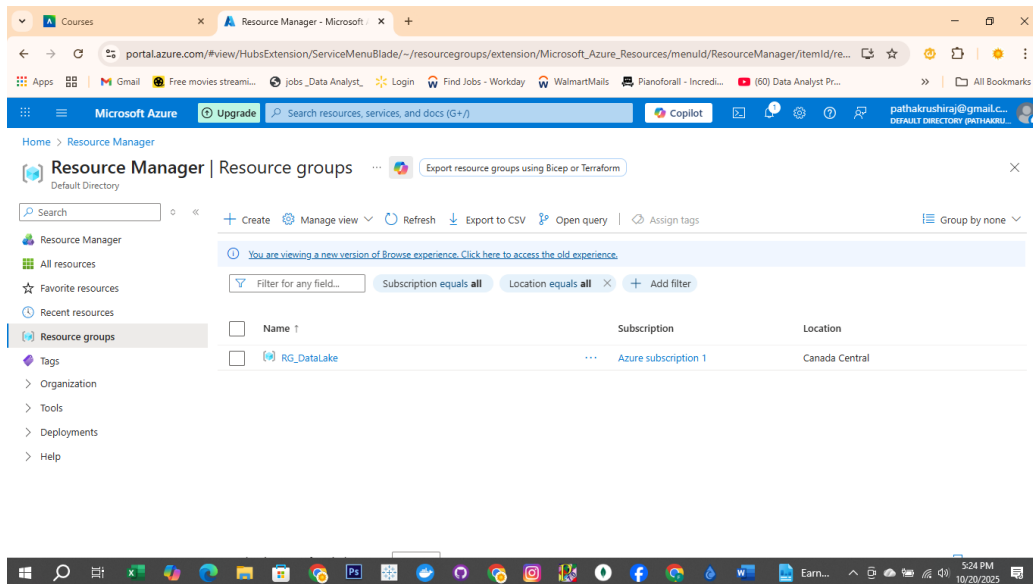

*Figure 1Azure Portal home page showing "$300 free credit"*
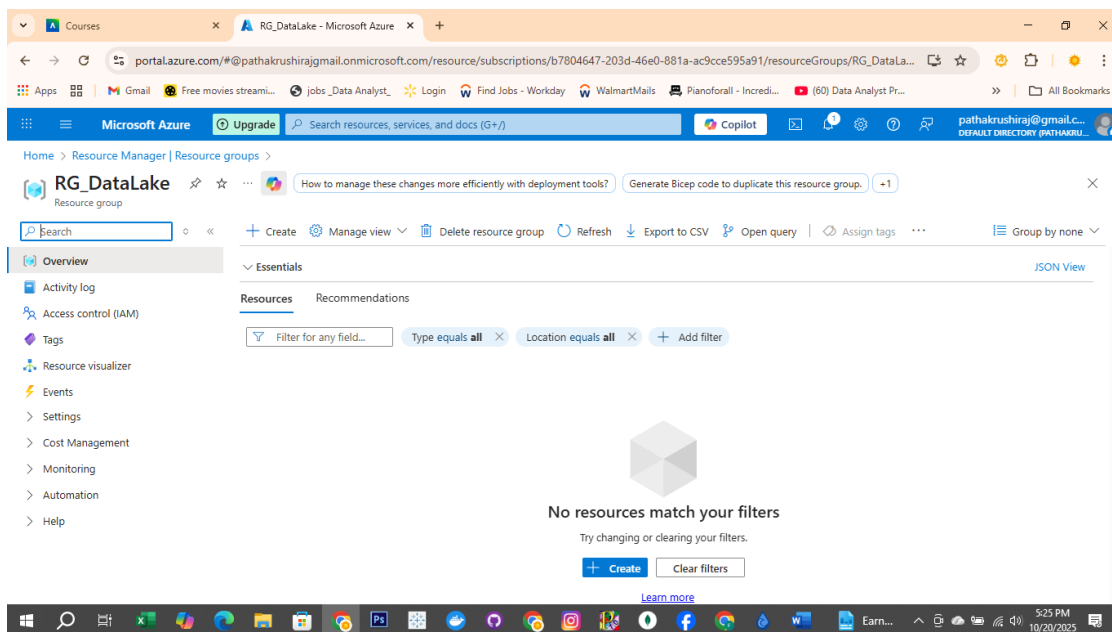
Figure 2– Resource Group Creation Form



Figure 3RG_DataLake Overview Page

# Step 2 – Storage Account Creation



*Figure 4Both Resource Groups Listed*

*Figure 5Storage Account Basics Tab Filled*



*Figure 6Advanced Tab with Hierarchical Namespace Enabled*

*Figure 7Storage Account Overview Page*



*Figure 8Containers List Page*

*Figure 9RG_Databricks Overview Page*



*Figure 10Uploaded sales.csv in Source Container*

# Step 3 – Databricks Workspace and Cluster

- Created a **Databricks workspace**.
- Configured a single-node cluster named **Cluster_SingleNode**.
- Attached the cluster to the notebook.



*Figure 11Azure Databricks Selection Screen*

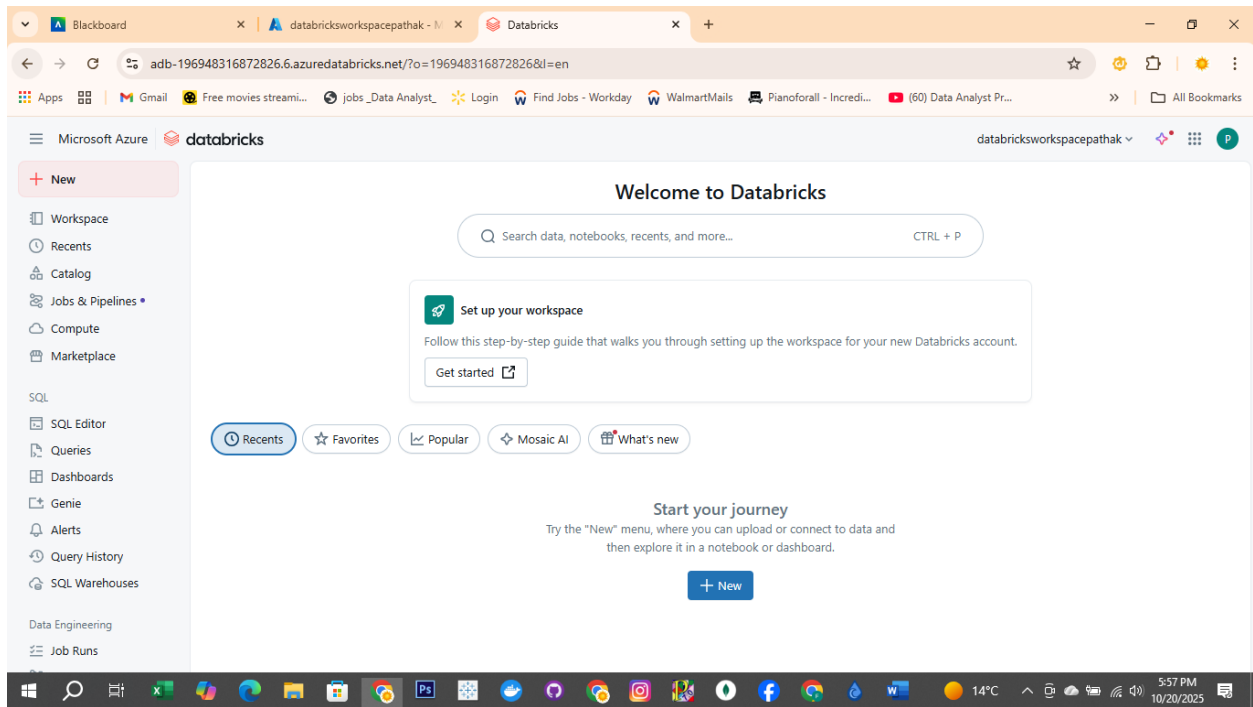*Figure 12Databricks Basics Tab Filled*



*Figure 13 Deployment Complete Page*

*Figure 14Databricks Home Page (After Launch)*



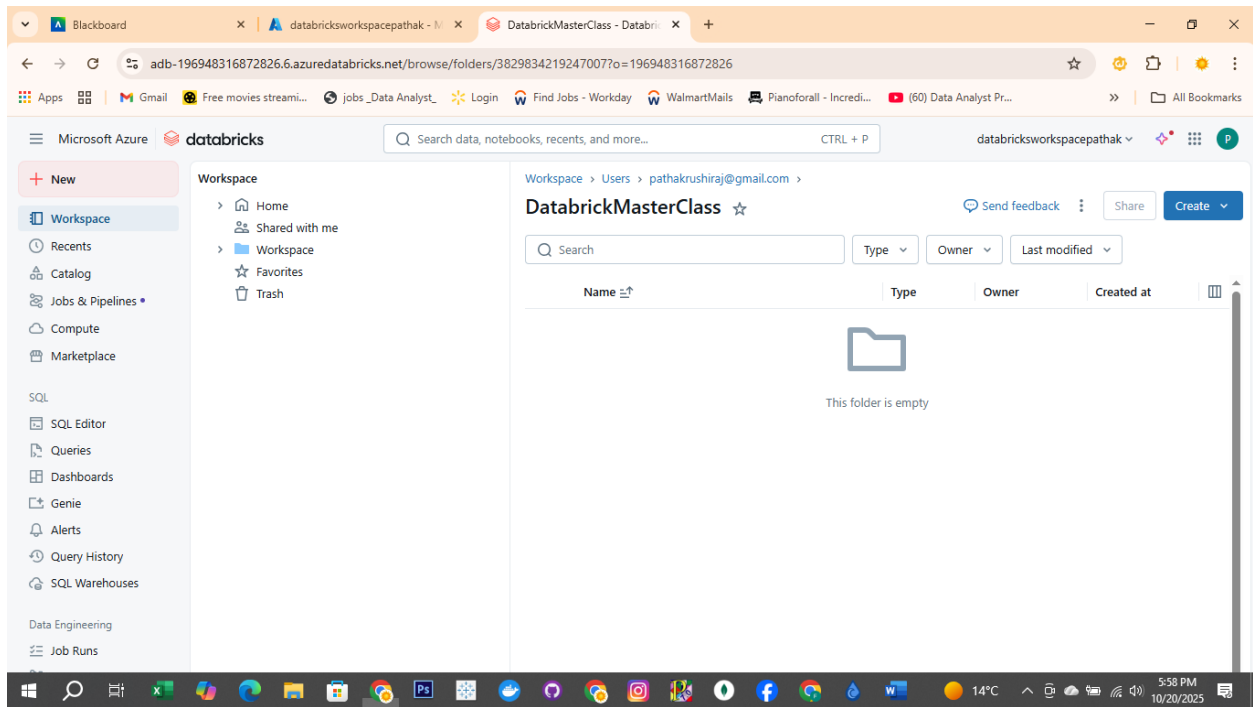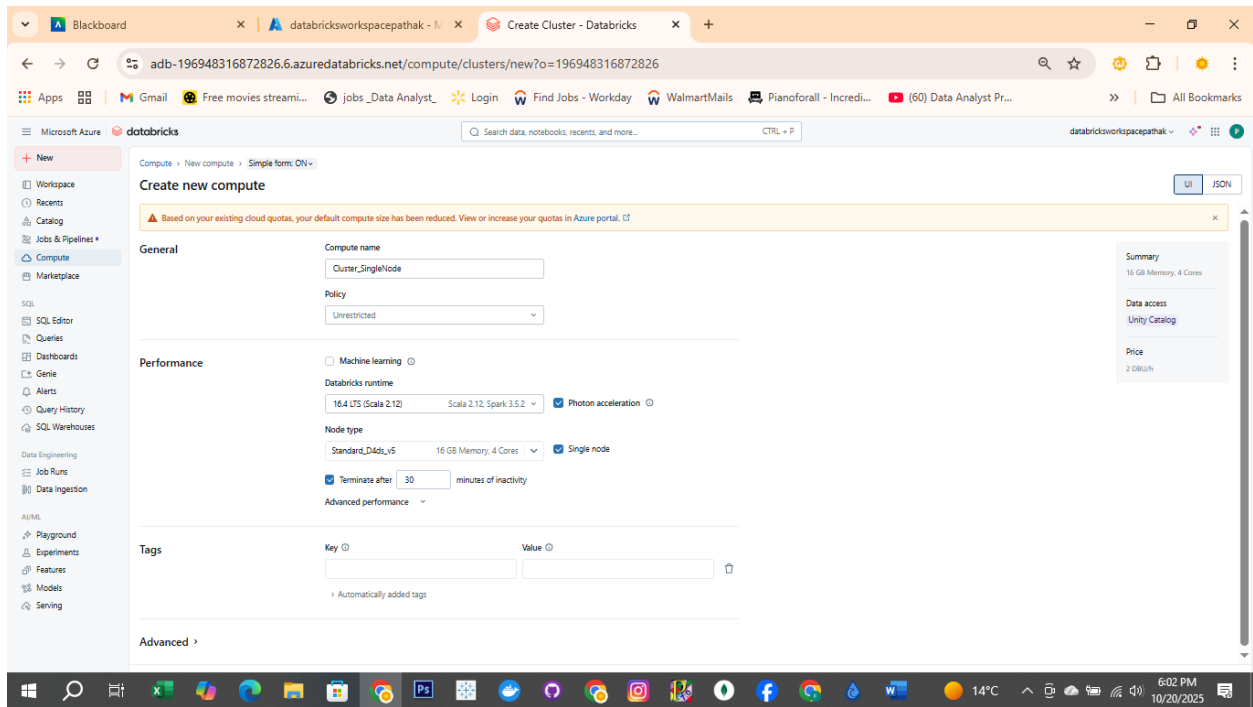*Figure 15Workspace Folder Created*

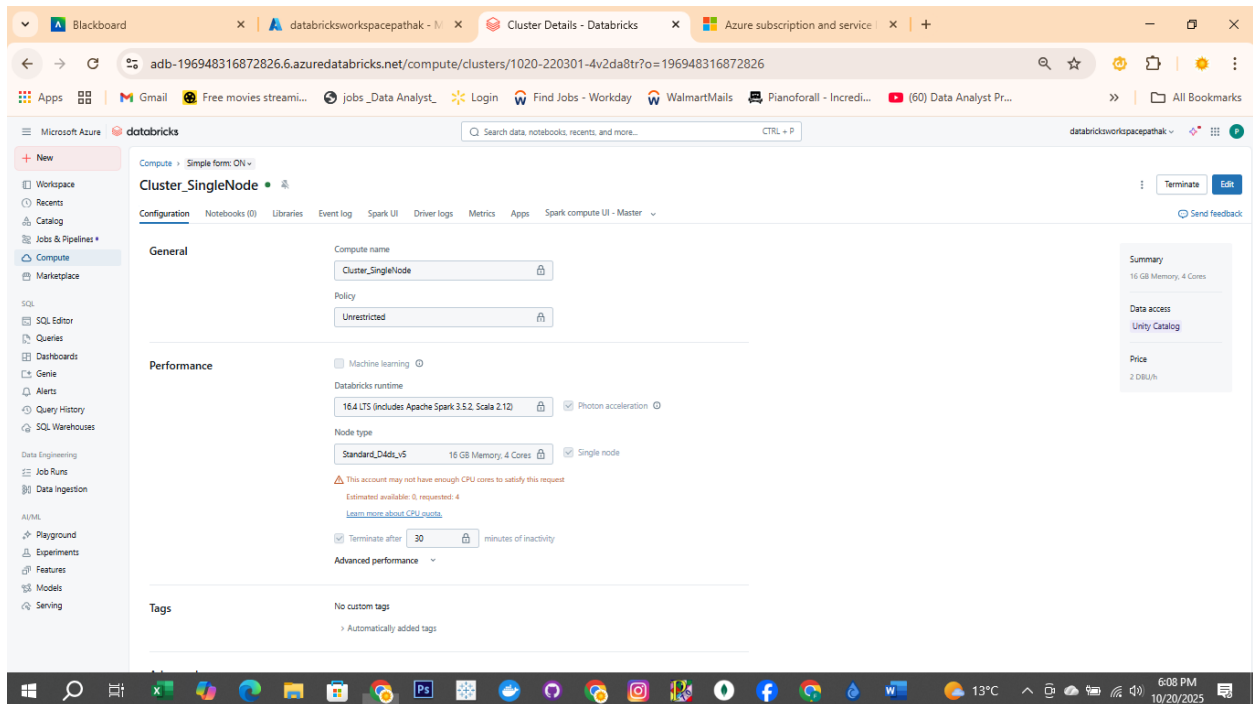*Figure 16Cluster Configuration Filled*



*Figure 17Active Cluster (Green Indicator)*

# Step 4 – Connecting Databricks to Azure Data Lake

- Used **access keys** for secure ABFSS (Azure Blob File System) connection.

spark.conf.set(

"fs.azure.account.key.datalakepathak.dfs.core.windows.net",
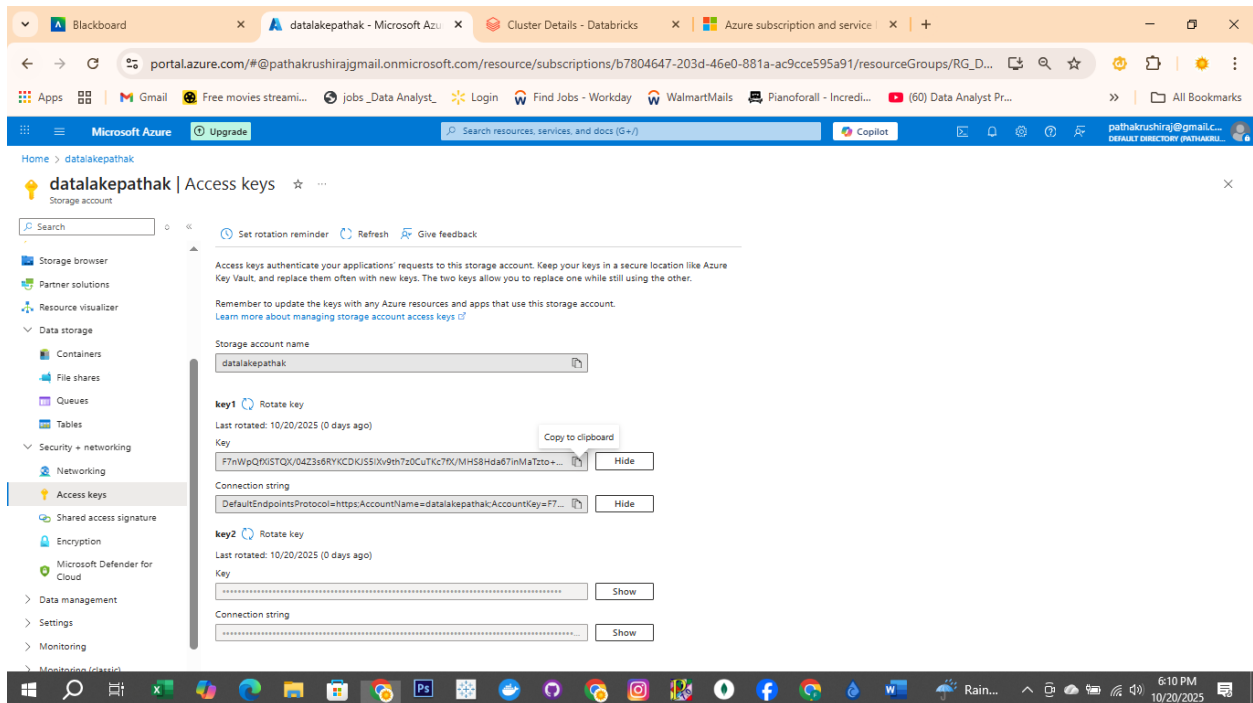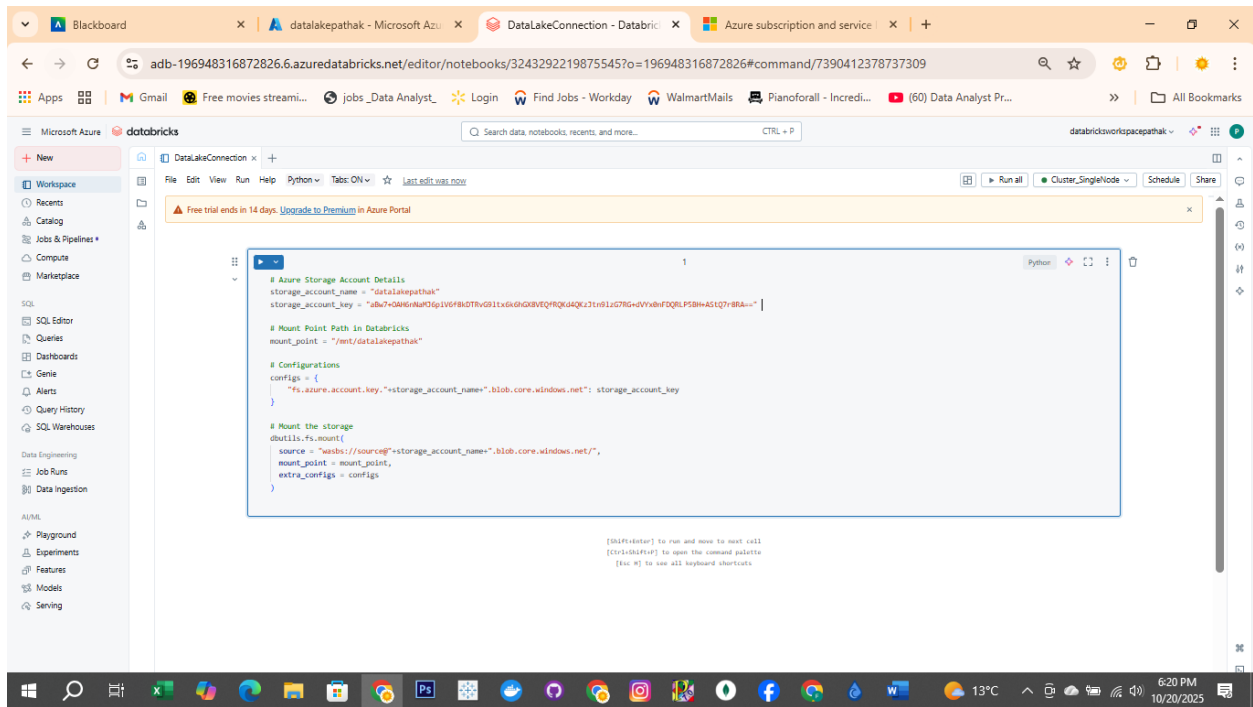
"<your-storage-access-key>"

)



*Figure 18Access Keys Page*

*Figure 19Code in Notebook*
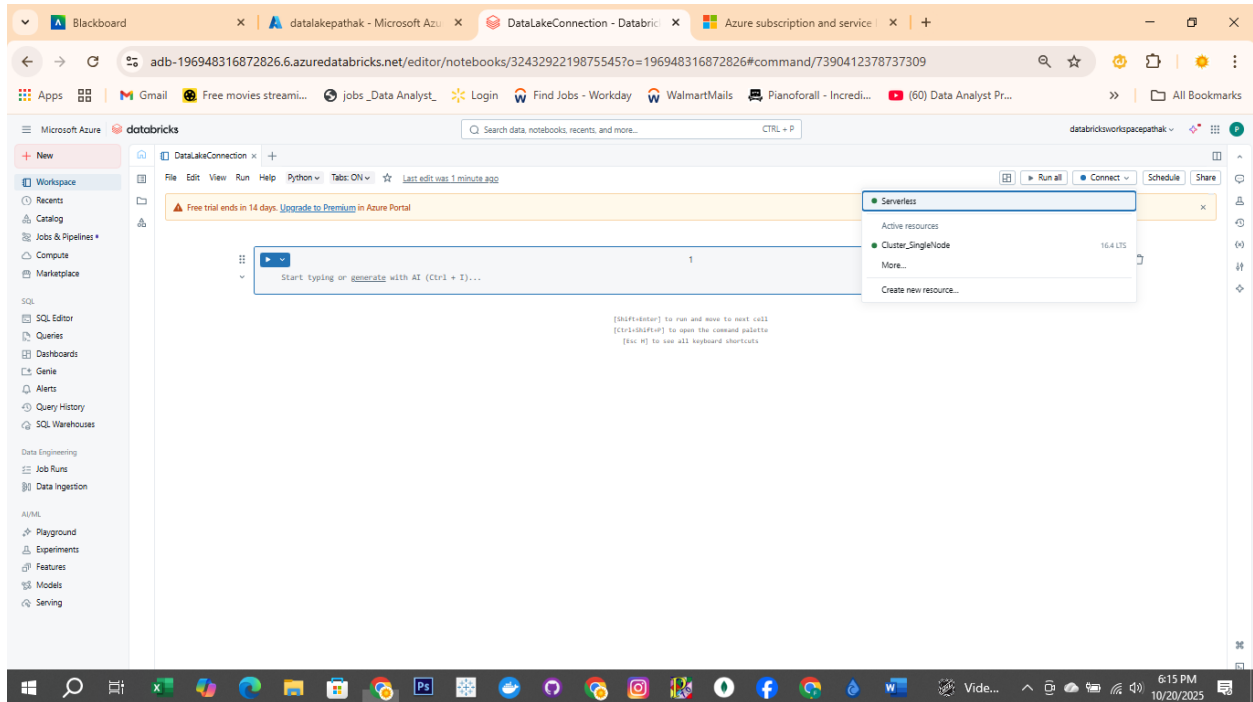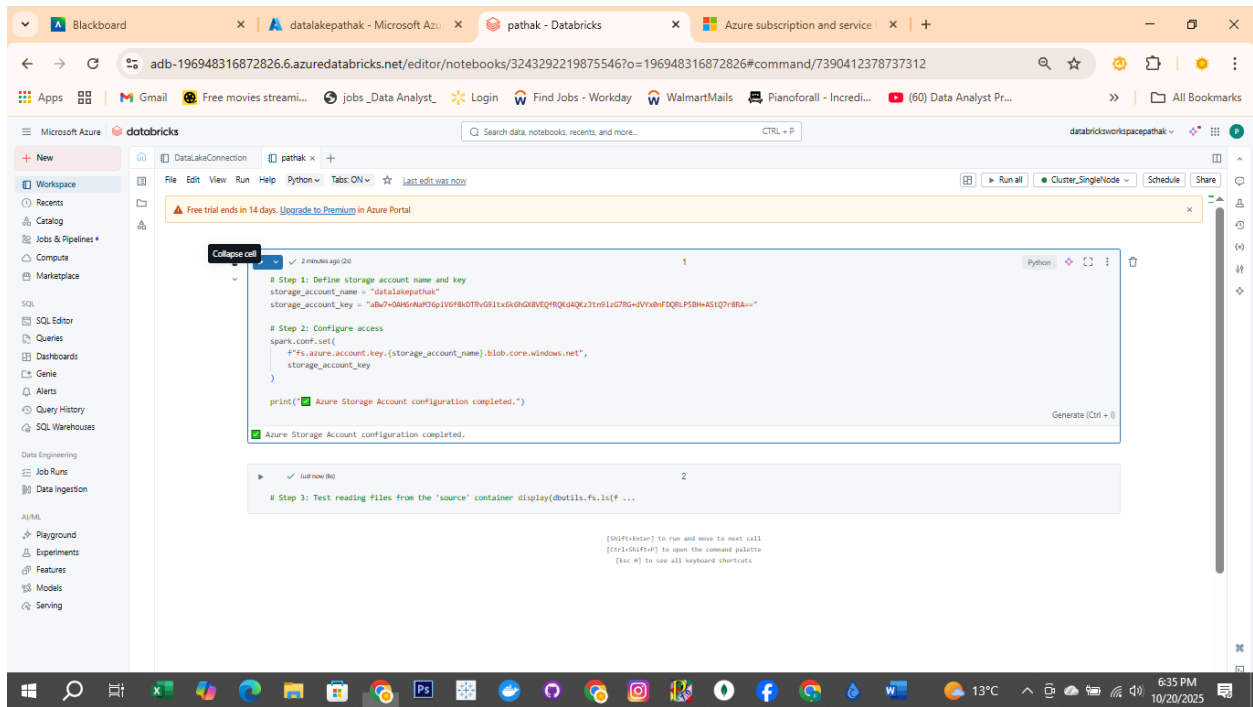


*Figure 20Notebook Creation Popup*

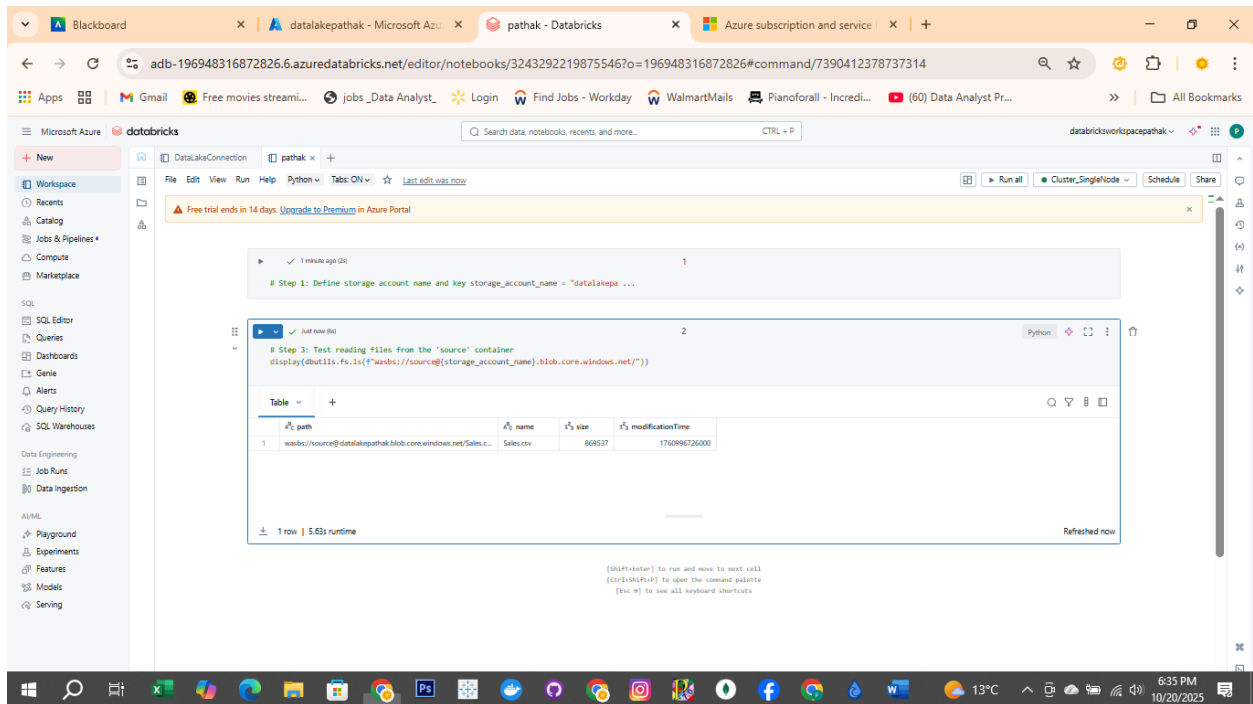*Figure 21Notebook Code Cell with Configuration*



*Figure 22Output showing sales.csv inside source container*

## Step 5 – Data Ingestion and Cleaning

- Read the Sales.csv file from the **source** container.
- Removed null values and duplicates using PySpark functions.

df_cleaned = df.dropna().dropDuplicates()



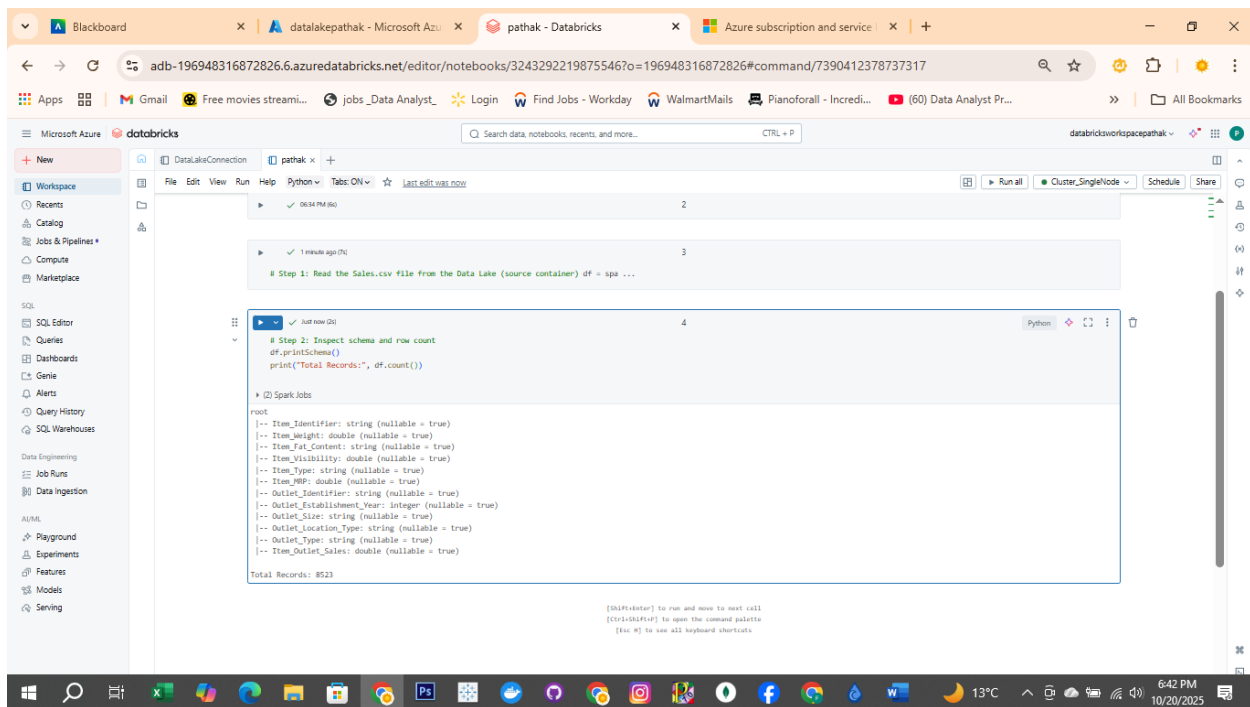*Figure 23Raw Data Loaded into Spark DataFrame*
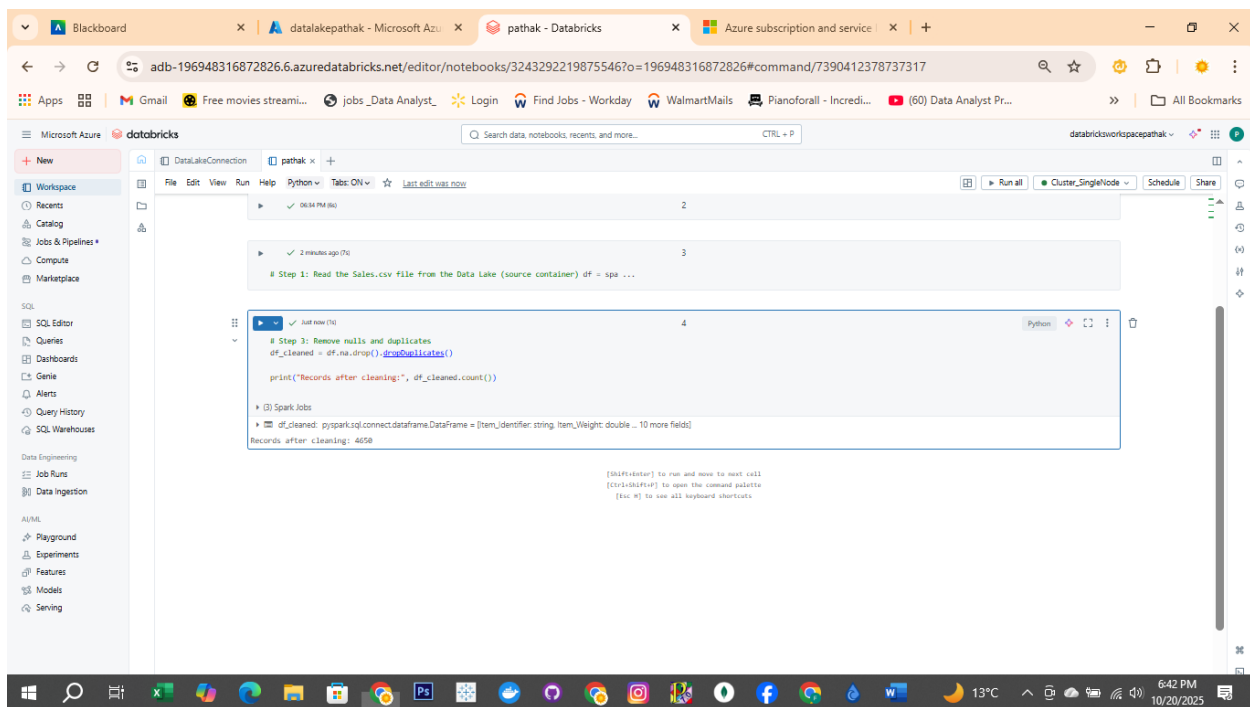
Figure 24Schema and Record Count
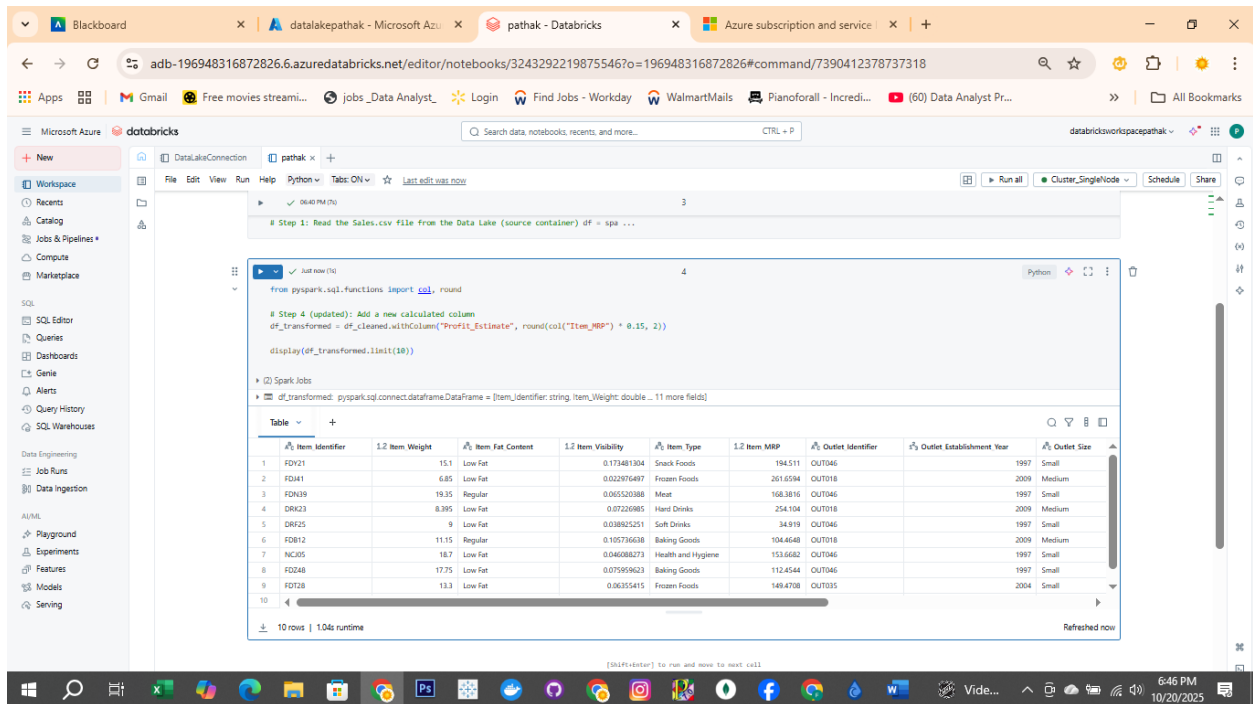


Figure 25Cleaning Step Results

## Step 6 – Data Transformation

- Grouped data by Outlet_Type and calculated total sales.

from pyspark.sql.functions import sum

df_transformed = df_cleaned.groupBy("Outlet_Type").agg(sum("Item_Outlet_Sales").alias("Total_Sales"))



*Figure 26DataFrame with Profit_Estimate Column*
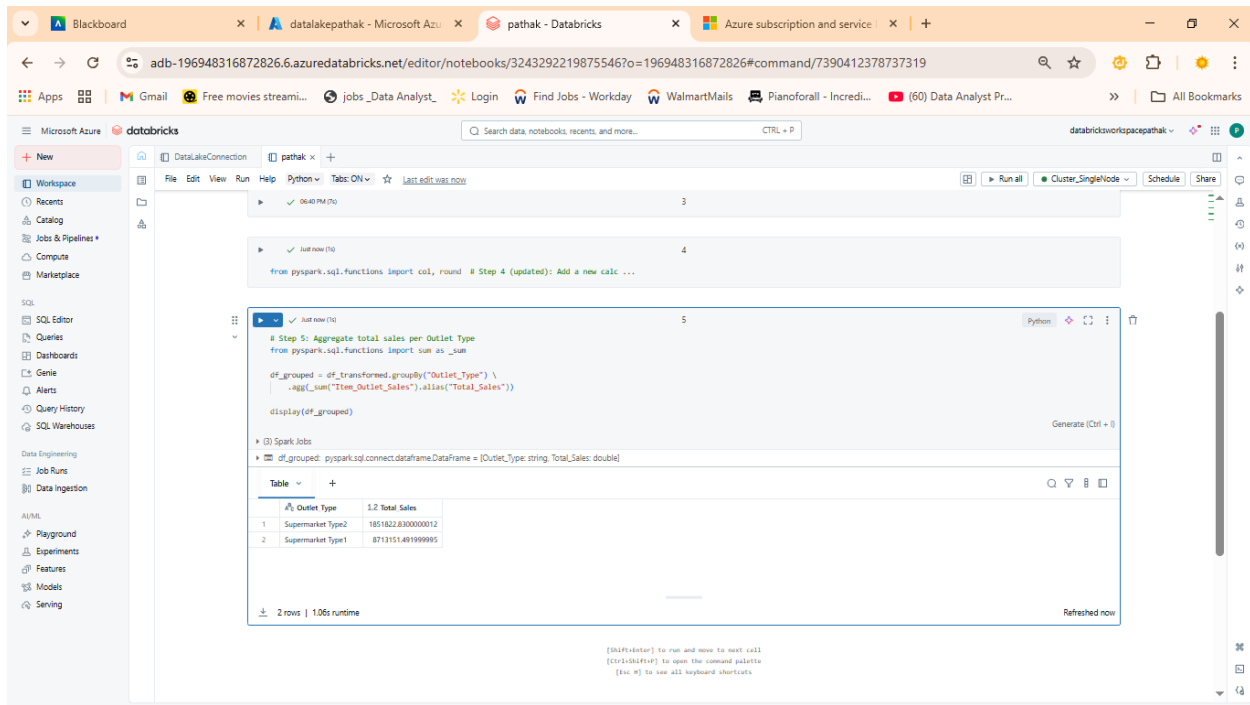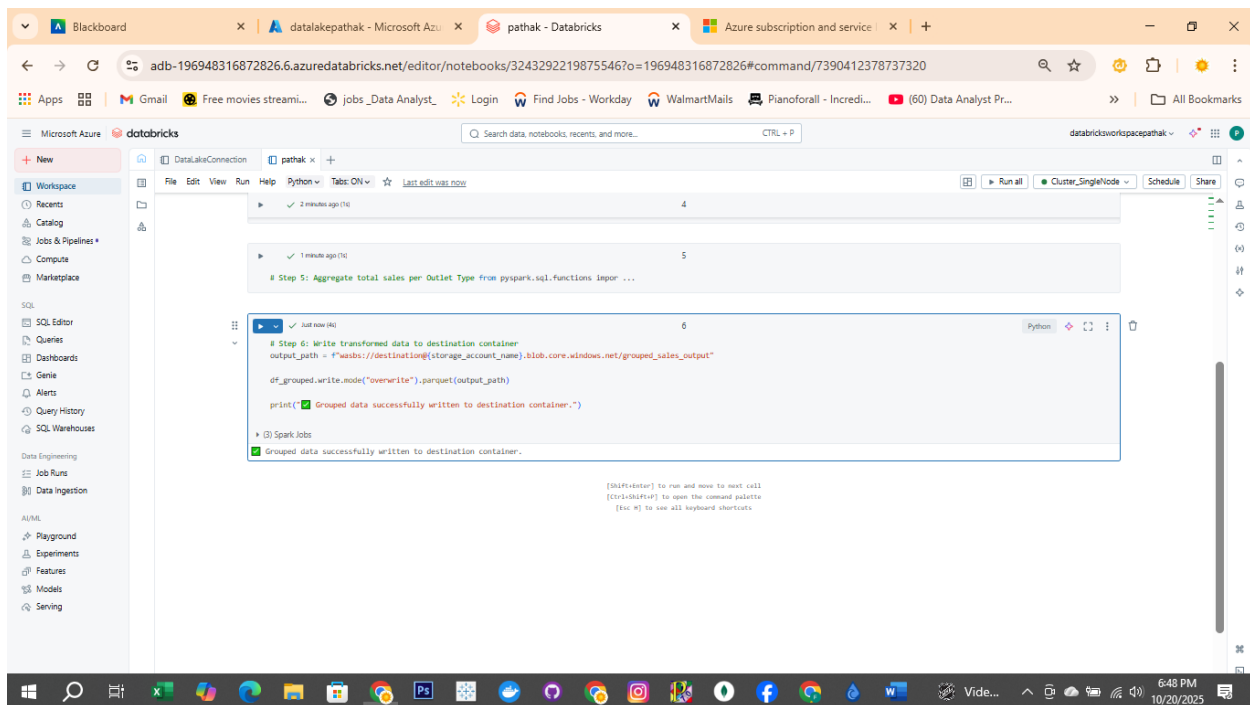
*Figure 27Grouped Total Sales by Outlet Type*



*Figure 28Output Written Confirmation*

## Step 7 – Writing Data to Data Lake

- Saved the transformed data in Parquet and Delta formats.

df_transformed.write.format("delta").mode("overwrite").save("abfss://destination@datalakepathak.dfs.core.windows.net/grouped_sales_output")



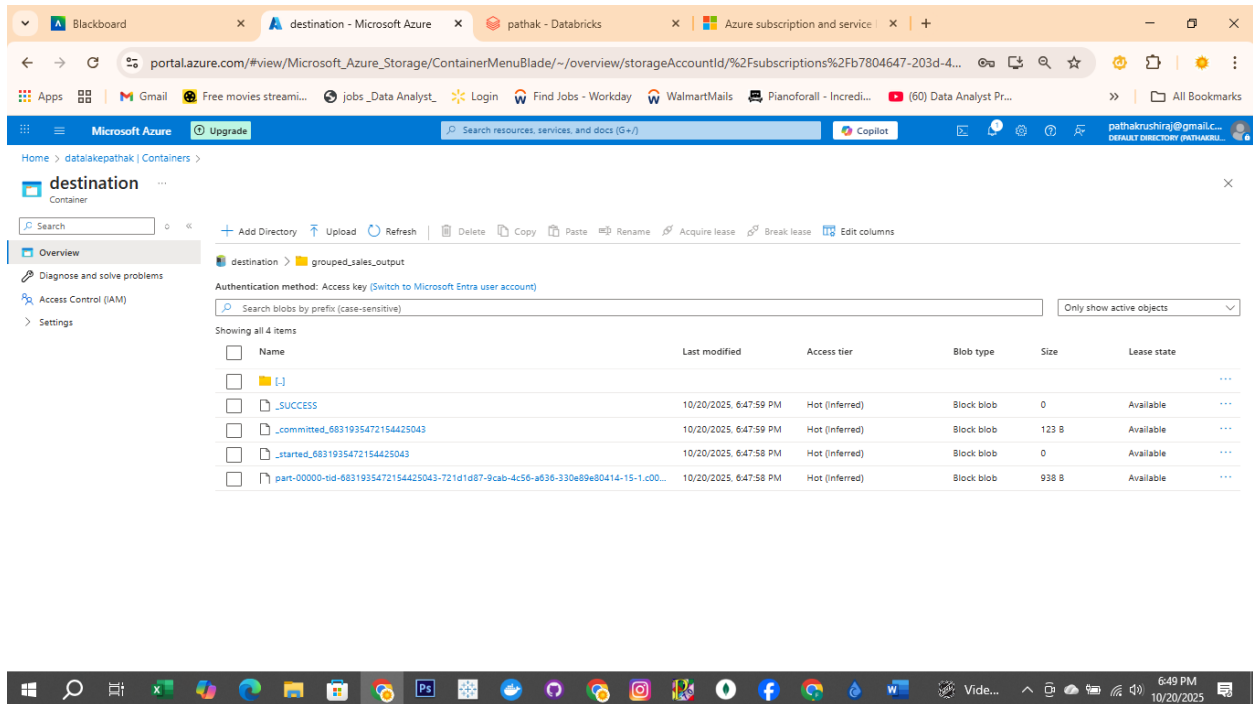*Figure 29Destination Container Showing group_output Folder*

# Step 8 – Creating Database and Delta Table

- Created a Hive Metastore database salesdb and registered a Delta table.

Spark.sql("CREATE DATABASE IF NOT EXISTS salesdb")

spark.sql("USE salesdb")

df_transformed.write.format("delta").mode("overwrite").saveAsTable("salesdb.external_sales")



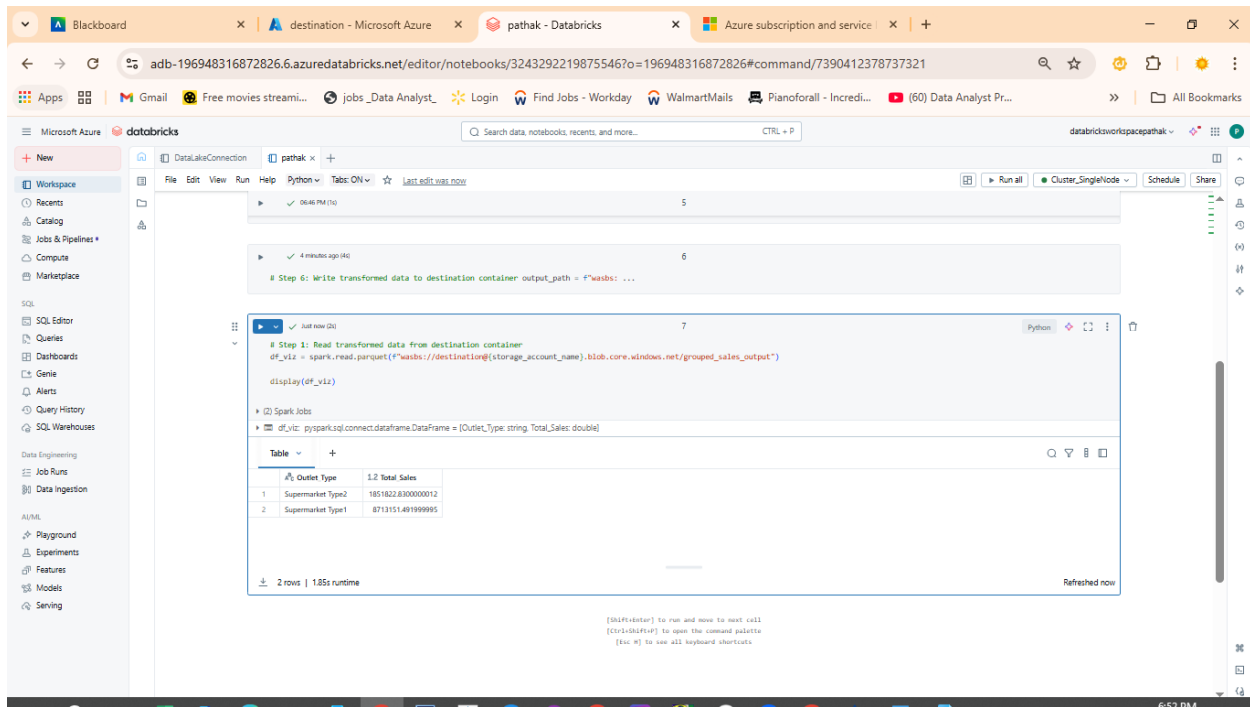*Figure 30Aggregated Data Displayed in Databricks*

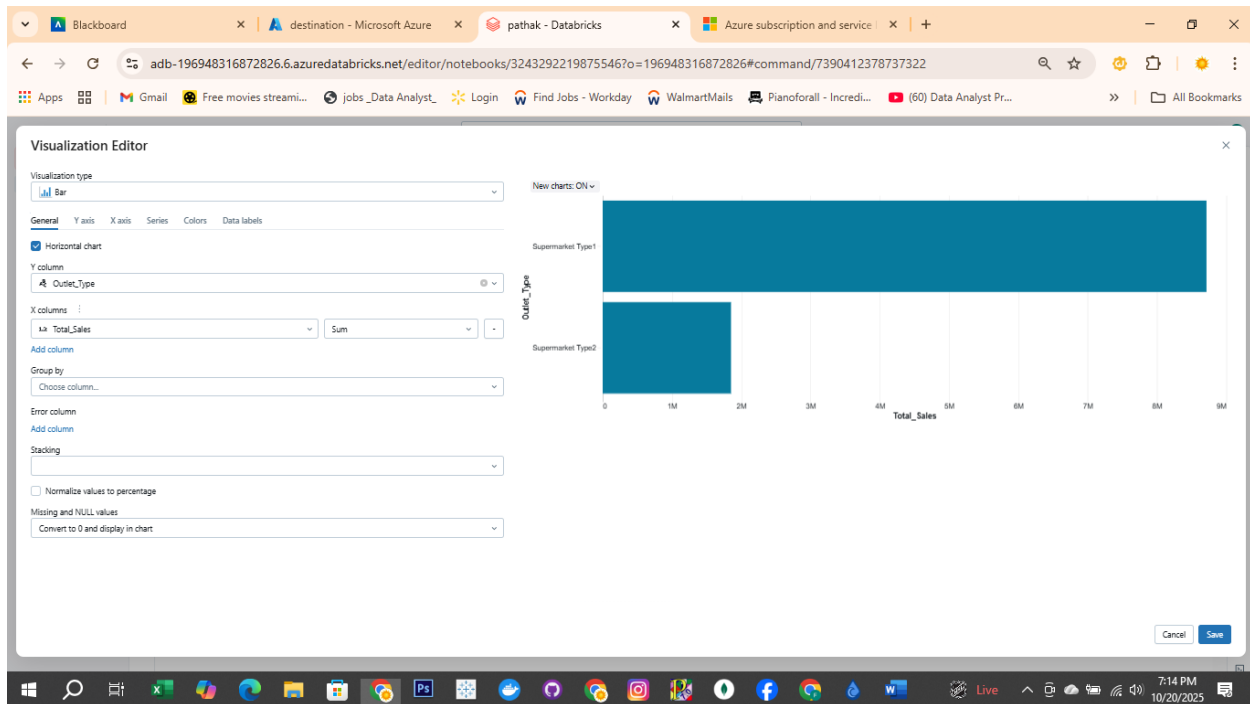*Figure 31Aggregated Data Displayed in Databricks*
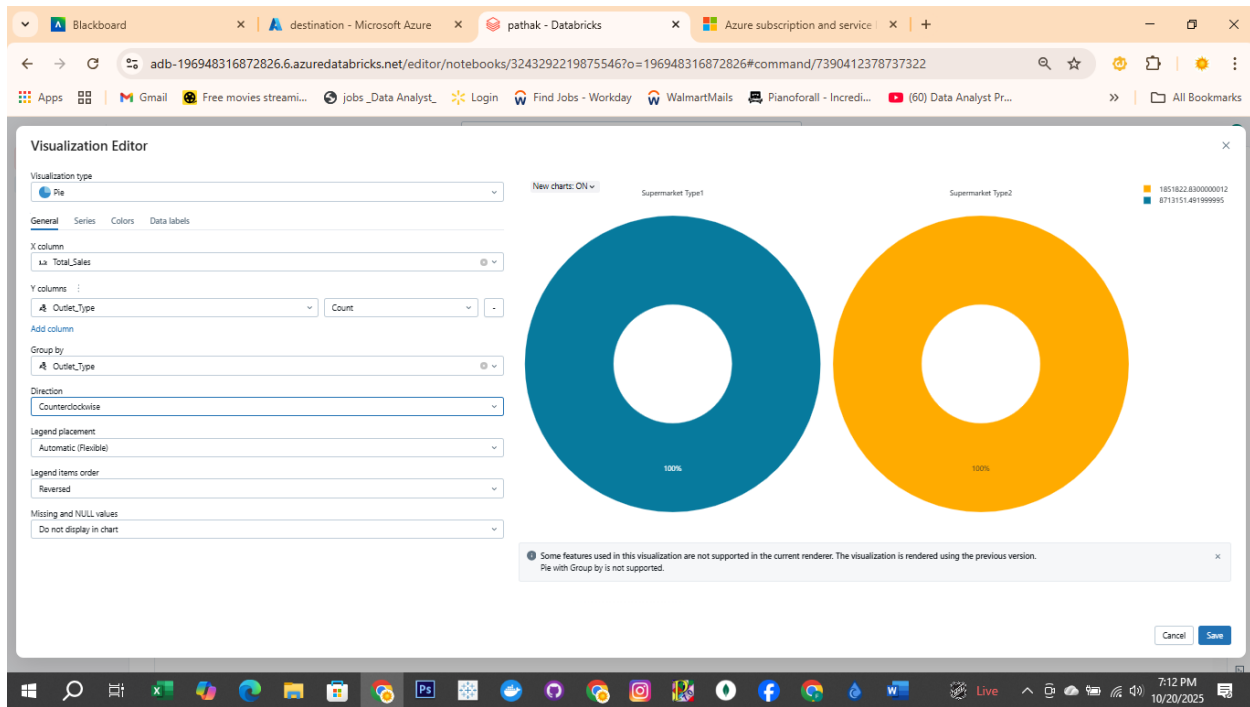


*Figure 32Bar Chart Total Sales by Outlet Type*

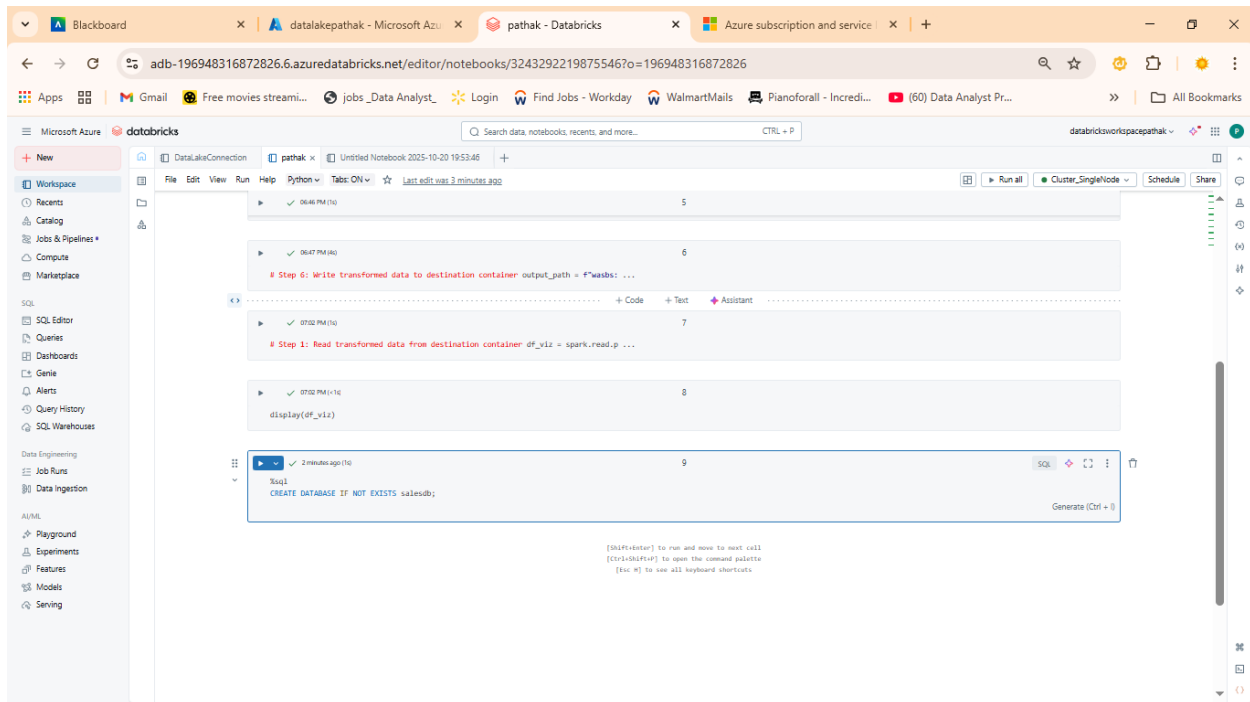*Figure 33Pie Chart Sales Share by Outlet Type*
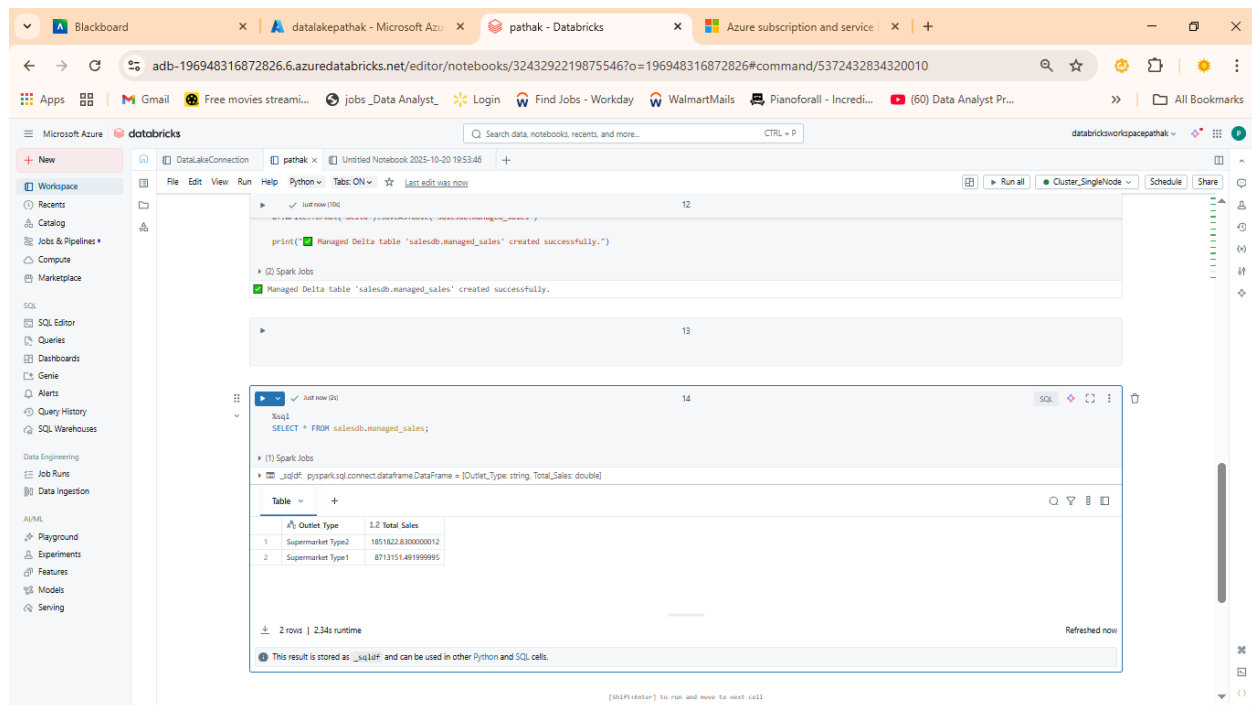


*Figure 34Create a Database*

*Figure 35Verify the Table*

## 5. Results and Observations

- Successfully established **secure integration** between Azure Databricks and ADLS Gen2.
- Verified data ingestion, transformation, and output in Delta format.
- Registered the transformed data as a queryable **Hive Metastore table**.
- Generated meaningful insights such as **total sales by outlet type**.
- *(Optional)* Visualization confirmed higher sales in Supermarket Type2.

## 6. Conclusion

This project highlights the power of **Azure Databricks** for building scalable and automated data engineering pipelines. Using PySpark, large-scale data can be processed efficiently and stored in Delta format for reliable analytics. Integrating Databricks with Azure Data Lake provides a secure, cloud-native architecture for enterprise-grade data processing.

By leveraging Databricks' flexibility and Azure's managed services, this project successfully demonstrates an end-to-end workflow — from ingestion to transformation, database registration, and querying — establishing a strong foundation for real-world big data analytics projects.