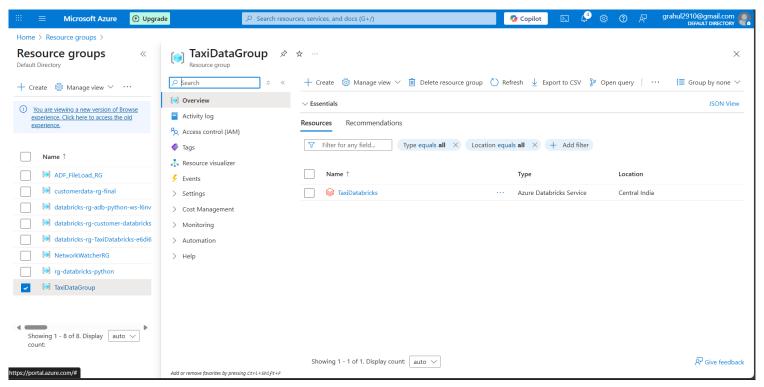# NYC Yellow Taxi Data Analysis using Azure and PySpark
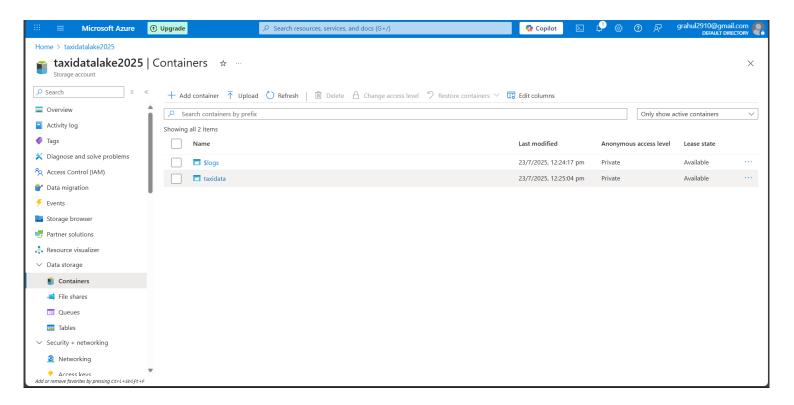
1. Created Azure Resources:
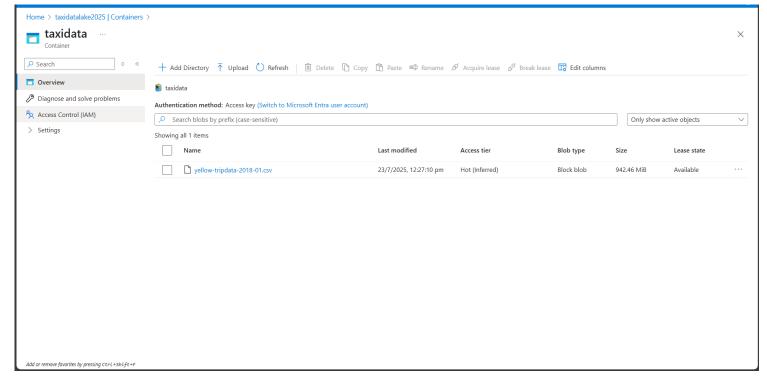
   o Resource Group
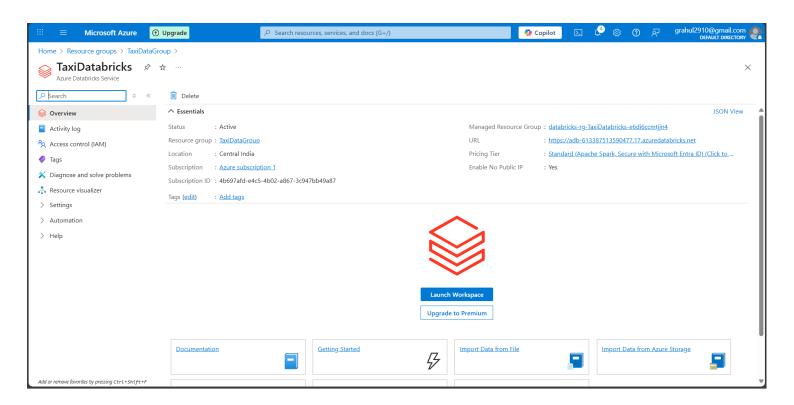


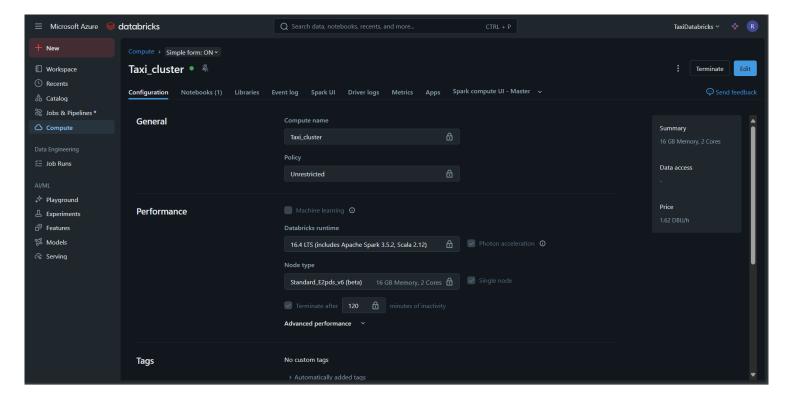   o Storage Account with ADLS Gen2 enabled
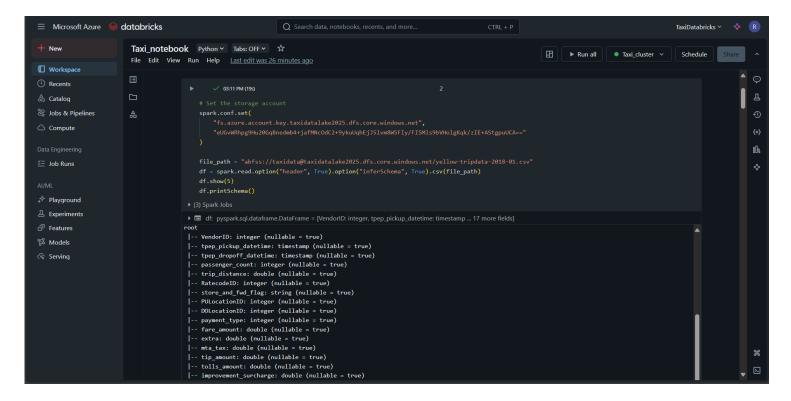
- Container named taxidata



- Azure Databricks workspace and cluster

2. Data Upload:

   o  Uploaded the dataset yellow_tripdata_2018-01.csv into the taxidata
      container in the storage account.



3. Databricks Notebook Setup:

   o  Connected Databricks to the storage using ABFSS protocol and OAuth
      credentials.

- Read the CSV file using an explicitly defined schema and converted date columns (tpep_pickup_datetime, tpep_dropoff_datetime) to timestamps.

Query 1: Top Pickup Locations

- Grouped data by PULocationID.

- Summed the number of passengers per pickup location.



Query 2: Top Dropoff Locations

- Grouped data by DOLocationID.

- Aggregated the number of passengers per dropoff location.

- Sorted to identify most common drop-off zones.

## Query 3: Revenue Generated by Each Vendor

- Created a derived column Revenue by summing:
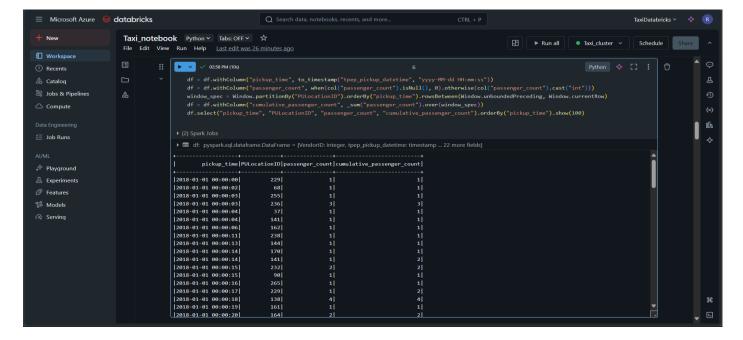
  o fare_amount + extra + mta_tax + tip_amount + tolls_amount + improvement_surcharge + congestion_surcharge

- Grouped the dataset by VendorID.

- Aggregated total revenue per vendor.



## Query 4: Moving Count of Payments by Payment Mode

- Used a Window function to calculate a rolling count of rides per payment_type, ordered by pickup timestamp.

- Enabled tracking of trends in payment methods over time.

## Query 5: Top 2 Highest Earning Vendors on a Particular Date

- Filtered rides from January 15, 2018.

- Computed Revenue per row as in Query 3.

- Aggregated total revenue, passenger count, and trip distance by vendor.

- Sorted and selected the top 2 vendors by revenue.



## Query 6: Route with Most Passengers

- Grouped by both PULocationID and DOLocationID to define a "route."

- Summed passenger_count for each route.

- Identified the route with the highest total passengers.

Query 7: Top Pickup Locations in the Last N Seconds

- Retrieved the maximum pickup timestamp from the dataset.

- Defined a time window (last 5 or 10 seconds).

- Filtered data to include only trips within that time frame.

- Aggregated and sorted pickup locations by passenger count to detect demand surges.



Output-

Home > Storage accounts > taxidatalake2025 | Containers >

## taxidata
Container

Add Directory   Upload   Refresh   |   Delete   Copy   Paste   Rename   Acquire lease   Break lease   |   Edit columns

taxidata

**Authentication method:** Access key (Switch to Microsoft Entra user account)

Search blobs by prefix (case-sensitive)

Only show active objects

Showing all 8 items

| | Name | Last modified | Access tier | Blob type | Size | Lease state | |
|---|---|---|---|---|---|---|---|
| | query1_top_pickup_locations | 25/7/2025, 5:20:14 pm | | | | | ... |
| | query2_top_dropoff_locations | 25/7/2025, 5:20:25 pm | | | | | ... |
| | query3_vendor_revenue | 25/7/2025, 5:20:38 pm | | | | | ... |
| | query4_moving_count | 25/7/2025, 5:33:17 pm | | | | | ... |
| | query5_top2_vendors | 25/7/2025, 5:33:32 pm | | | | | ... |
| | query6_top_route | 25/7/2025, 5:33:43 pm | | | | | ... |
| | query7_top_pickup_locations | 25/7/2025, 5:33:53 pm | | | | | ... |
| | yellow-tripdata-2018-01.csv | 23/7/2025, 12:27:10 pm | Hot (Inferred) | Block blob | 942.46 MiB | Available | ... |

Add or remove favorites by pressing Ctrl+Shift+F