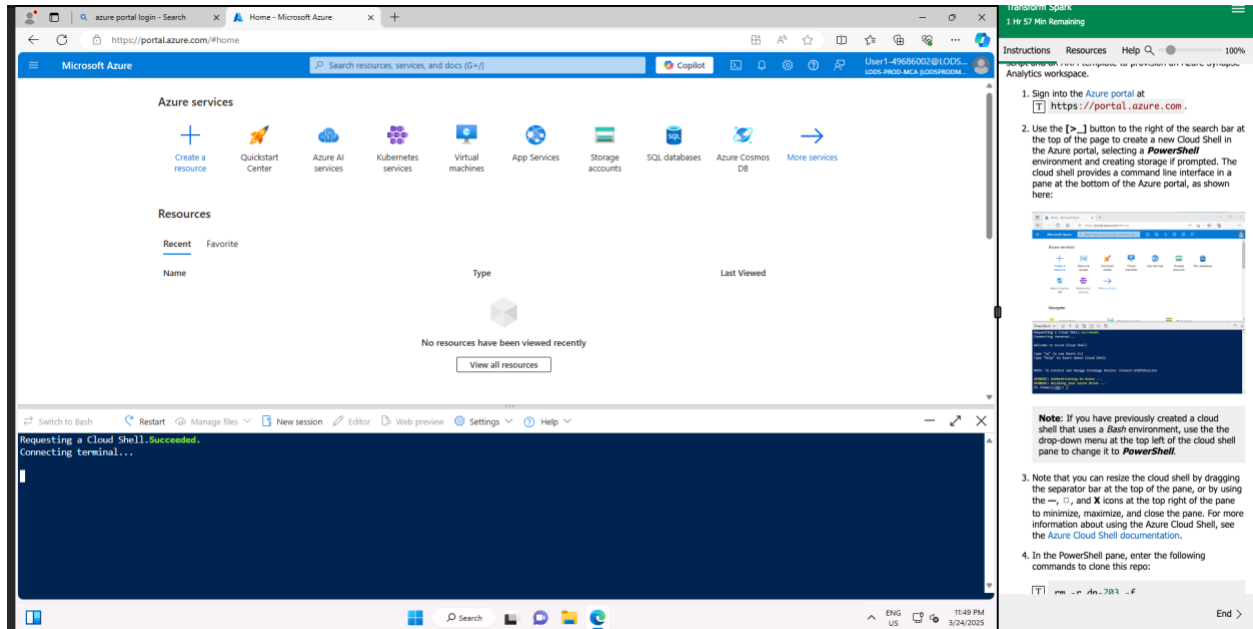# Lab 6 Transform Spark
## BUAN 6390.001 – Analytics Practicum

**Name: Swastik Bhatnagar**
**Net ID: Sxb220210**

## Connecting to Powershell:

Using Spark notebook to transform data:
Opening Synapse Studio:



Viewing the data:

Using Github to download file:



Importing the file Spark Transform.ipynb and attaching the notebook to spark pool:



Loading the Source data:

Transforming the Data Structure:

## Saving the transformed data:



## Partitioning the data by Year and month:



## Using SQL to transform data:

Dropping the external tables from the metastore without deleting the files:

Deleting the Resource group:



Conclusion:

In this lab, I learned how to use Spark for data transformation, from setting up in Synapse Studio to working with GitHub and Spark notebooks. I explored data, applied transformations, and optimized it by partitioning by year and month. Using SQL, I refined the data while also managing external tables without losing underlying files. Lastly, I practiced resource management by deleting the resource group.