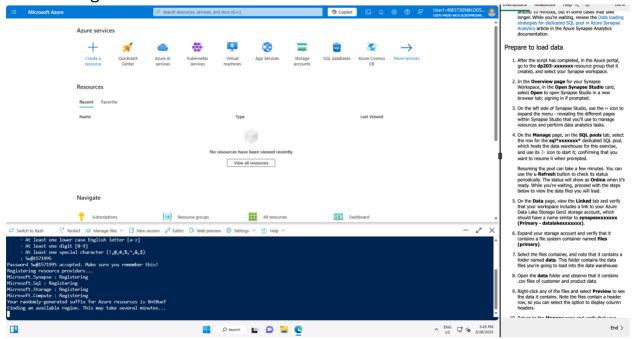# Lab 9 Load Data Warehouse
## BUAN 6390.001 – Analytics Practicum
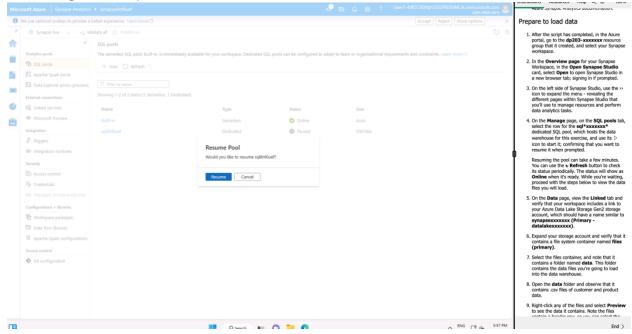
**Name: Swastik Bhatnagar**
**Net ID: Sxb220210**

## Connecting to PowerShell:



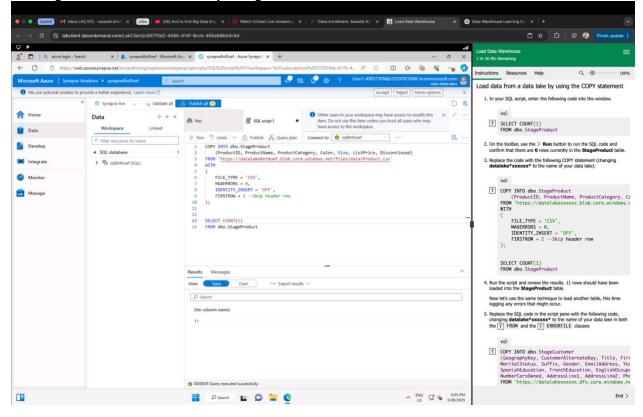## Resuming the SQL pool:
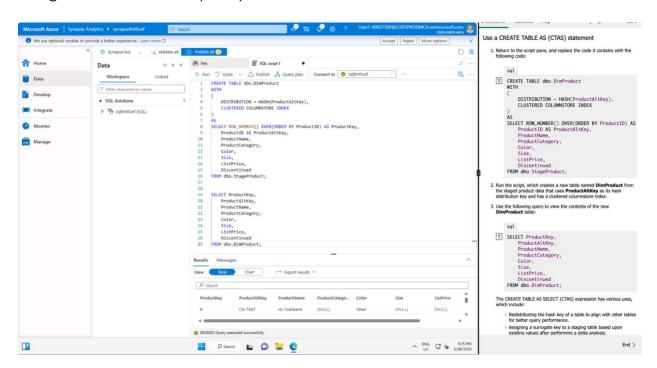
Loading the Data warehouse table:
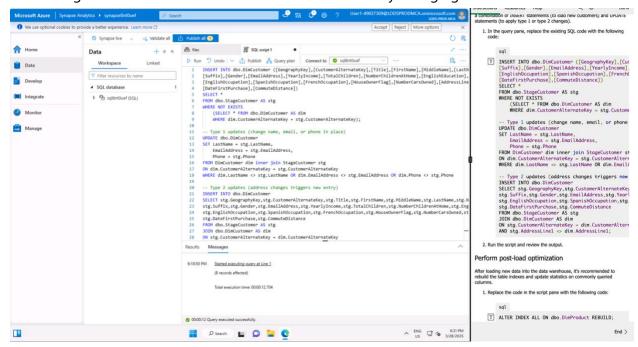
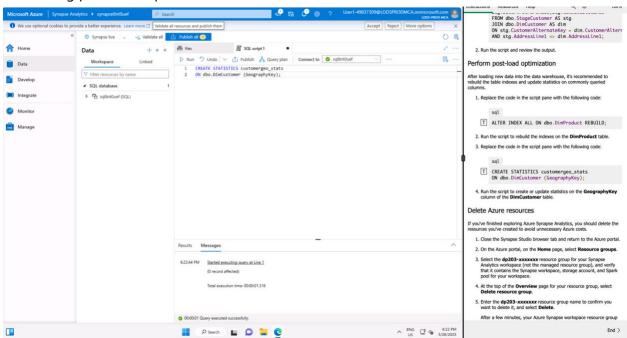Loading data from a data lake by using the COPY statement:
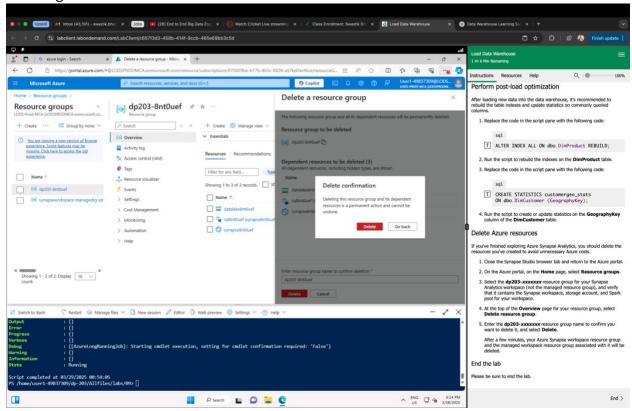


Using a CREATE TABLE AS (CTAS) statement:

Combining INSERT and UPDATE statements to load a slowly changing dimension table:



Performing post-load optimization:

Deleting Azure resources:



Conclusion:

In this lab, I learned how to resume SQL pools, load data from a data lake using the COPY statement, and utilize the CREATE TABLE AS (CTAS) statement for efficient table creation. Additionally, I practiced combining INSERT and UPDATE statements to handle slowly changing dimension tables and performed post-load optimization. This process enhanced my understanding of data loading techniques and best practices for maintaining data integrity and optimizing performance within a data warehouse environment.