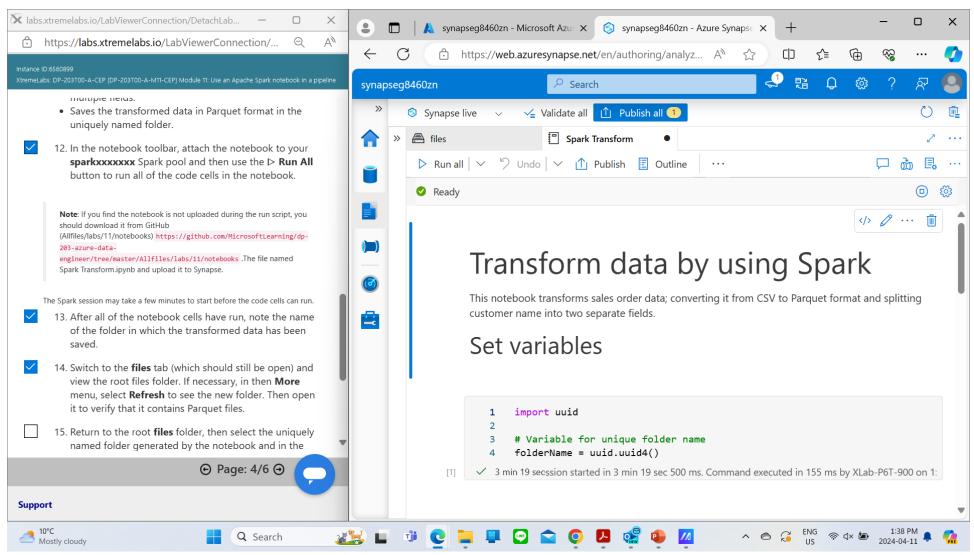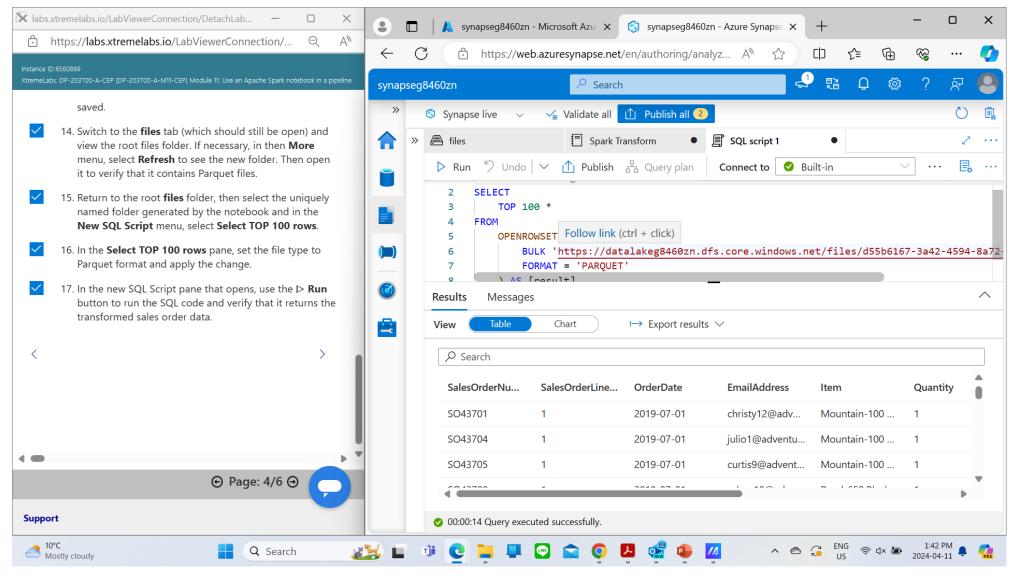# Lab 10 – Build a data pipeline in Azure Synapse Analytics

# Run a Spark notebook interactively (1)
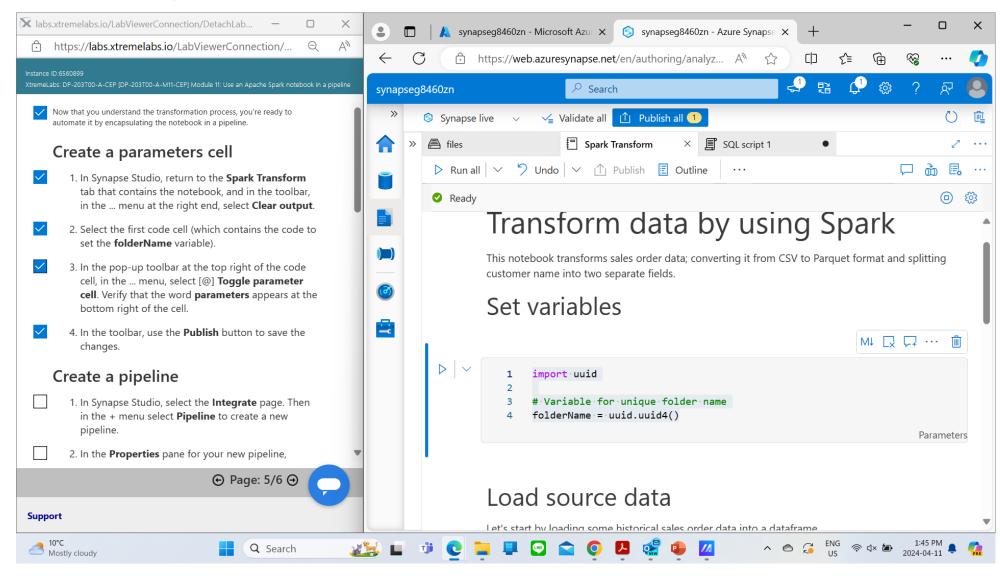
# Run a Spark notebook interactively (2)
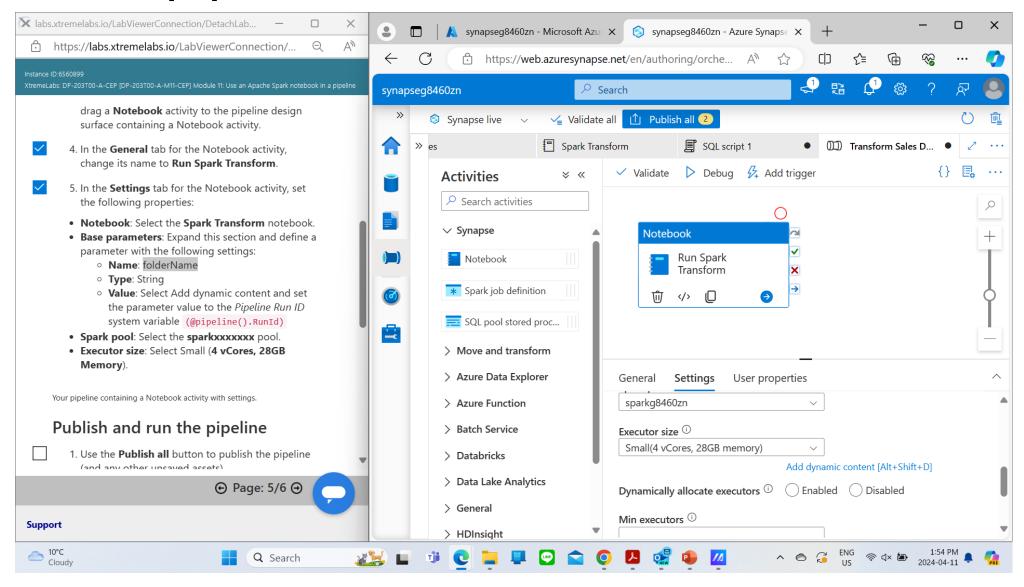
# Run the notebook in a pipeline

# Create a parameters cell

# Create a pipeline

# Publish and run the pipeline