

Project – Serverless Data Processing

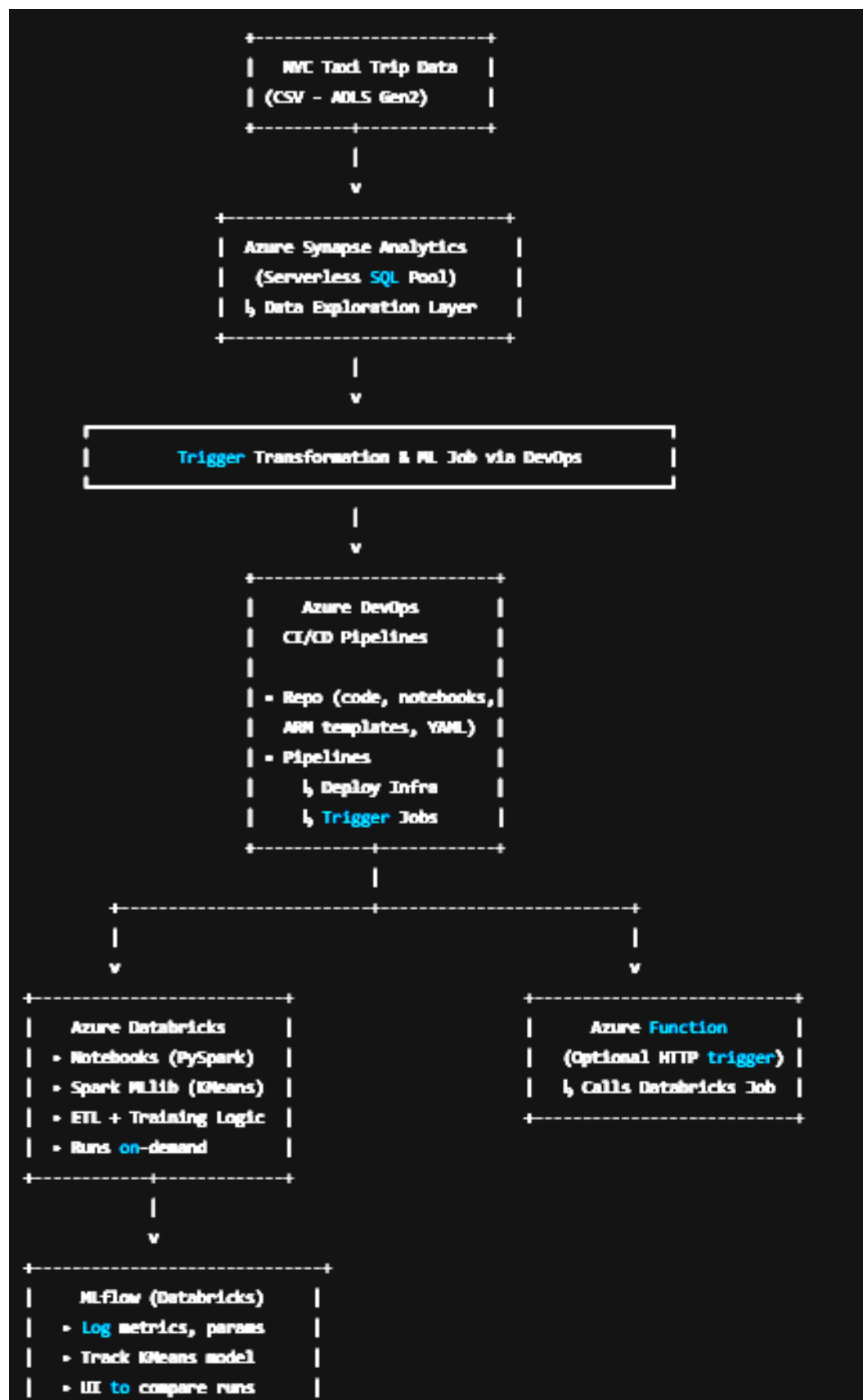
ETL Problem:

- ☒ ~~Use Azure Synapse Serverless SQL for raw data exploration directly on ADLS Gen2.~~
- ☒ ~~Trigger Azure Databricks jobs using Azure Functions for scalable transformation.~~
- ☒ ~~Code and infra managed via Azure DevOps repos and ARM templates.~~

ML Problem:

- ☒ ~~Create serverless training pipelines in Databricks with on-demand clusters.~~
- ☒ ~~Train customer segmentation model (KMeans) using Spark MLlib.~~
- ☐ Deploy using CI/CD jobs .triggered via DevOps pipelines
- ☒ ~~and monitored via MLflow.~~
- ☐ PPT
- ☐ Architecture Diagram
- ☐ YT Video
- ☐ Report Refinement
- ☐ Docs Uploading to drive
- ☐ Explaining project to team members

Architecture Diagram –



Task 1) Data Exploration Using Azure Synapse Analytics

Microsoft Azure

Search resources, services, and docs (G+/)

Copilot

[Home](#) > [Azure Synapse Analytics](#) >

Create Synapse workspace

Create a Synapse workspace to develop an enterprise analytics solution in just a few clicks.

Project details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all of your resources.

Subscription *

MML Learners

Resource group *

rg-azuser3611_mml.local-fRfkX

Create new

Managed resource group

Enter managed resource group name

Workspace details

Name your workspace, select a location, and choose a primary Data Lake Storage Gen2 file system to serve as the default location for logs and job output.

Workspace name *

dataexplusingsynapseanalytics

Region *

(Asia Pacific) Central India

Select Data Lake Storage Gen2 *

☒ From subscription ☐ Manually via URL

Workspace name must be between 1 and 50 characters long.

Workspace name must contain only lowercase letters or numbers or hyphens.

Workspace name must start with a letter or a number.

Workspace name must end with a letter or a number.

Workspace name must not contain '-ondemand' word.

Workspace name must be unique.

Review + create

< Previous

Next: Security >

Microsoft Azure

Search resources, services, and docs (G+/)

[Home](#) > [Azure Synapse Analytics](#) >

Create Synapse workspace

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all of your resources.

Subscription *

MML Learners

Resource group *

rg-azuser3611_mml.local-fRfkX

Create new

Managed resource group

Enter managed resource group name

Workspace details

Name your workspace, select a location, and choose a primary Data Lake Storage Gen2 file system to serve as the default location for logs and job output.

Workspace name *

dataexplusingsynapseanalytics

Region *

(Asia Pacific) Central India

Select Data Lake Storage Gen2 *

☒ From subscription ☐ Manually via URL

Account name *

(New) synapseadlsgen2sakshi

Create new

File system name *

(New) synapse

Create new

Review + create

< Previous

Next: Security >

Home > Azure Synapse Analytics >

Create Synapse workspace

✓

Validation succeeded

* Basics * Security Networking Tags **Review + create**

Product Details

Azure Synapse Analytics workspace
by Microsoft
[Terms of use](#) | [Privacy policy](#)

Serverless SQL est. cost/TB ⓘ
5.00 USD

Terms

By clicking Create, I (a) agree to the legal terms and privacy statement(s) associated with the Marketplace offering(s) listed above; (b) authorize Microsoft to bill my current payment method for the fees associated with the offering(s), with the same billing frequency as my Azure subscription; and (c) agree that Microsoft may share my contact, usage and transactional information with the provider(s) of the offering(s) for support, billing and other transactional activities. Microsoft does not provide rights for third-party offerings. For additional details see [Azure Marketplace Terms](#).

Basics

SubscriptionMML Learners

Resource grouprg-azuser3611_mml.local-frfkX

Create

Create

< Previous

Next >

[Download a template for automation](#)

Microsoft Azure

Search resources, services, and docs (G+/)

Copilot

azuser3611_mml.local@...
TECHADEMY LEARNING SOLUTIONS

Home >

Microsoft.Azure.SynapseAnalytics-20250714093653 | Overview

Deployment

Search

Delete

Cancel

Redeploy

Download

Refresh

Overview

Inputs

Outputs

Template

Deployment is in progress

Deployment name : Microsoft.Azure.SynapseAnalytics-20250714093653

Subscription : MML Learners

Resource group : rg-azuser3611_mml.local-frfkX

Start time : 7/14/2025, 9:50:53 AM

Correlation ID : 2ac7efb0-d530-4a06-8265-2e72f59a87f0

Deployment details

Resource	Type	Status	Operation details
synapseadlsgen2sakshi	Storage account	Accepted	Operation details

Give feedback

[Tell us about your experience with deployment](#)

Microsoft Defender for Cloud

Secure your apps and infrastructure

[Go to Microsoft Defender for Cloud >](#)

Free Microsoft tutorials

[Start learning today >](#)

Work with an expert

Azure experts are service provider partners who can help manage your assets on Azure and be your first line of support.

[Find an Azure expert >](#)

Azure Synapse Analytics

Loading

Microsoft Azure | Synapse Analytics | dataexplusingsynapseanalytics


We use optional cookies to provide a better experience. [Learn more](#)


Accept Reject More options


Synapse Analytics workspace

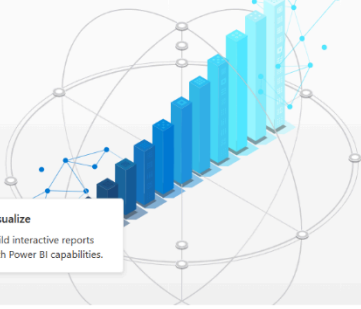
dataexplusingsynapseanalytics

New



**Ingest**
Perform a one-time or scheduled data load.

**Explore and analyze** ☒
Learn how to get insights from your data.

**Visualize**
Build interactive reports with Power BI capabilities.



Discover more

 Knowledge center  Browse partners

Recent resources

Microsoft Azure | Synapse Analytics | dataexplusingsynapseanalytics

We use optional cookies to provide a better experience. [Learn more](#)

Search

azuser3611_mml.local@techademy.com
TECHADEMY LEARNING SOLUTIONS PRIVATE LIMITED

Data

Workspace Linked

Filter resources by name

- Azure Data Lake Storage Gen2 2
- dataexplusingsynapseanalytics (Pri...
- synapse (Primary)
- (Attached Containers)

synapse

Other users in your workspace may have access to modify this workspace.

Upload New folder Select all Refresh

← → ↑ ↓ synapse

Name	Last Modified
No data available in this blob container	

Showing 0 to 0 of 0 entries

Upload Files

synapse

Destination folder /

File Upload

2020_Yellow_Taxi_Trip_Data.csv

☐ Overwrite existing files

File name	Size	Action
2020_Yellow_Taxi_Trip_Data.csv	2.23 GB	Remove

Upload Cancel

Search

🔊

📄

🔔

⚙️

?

azuser3611_mml.local@techademy.com
TECHADEMY LEARNING SOLUTIONS PRIVATE LIMITED

Other users in your workspace may have access to modify workspace.

Select all Refresh

Last Modified

tainer

Upload details

Uploading 1 files to 'synapse/'

2.04 GB of 2.23 GB

2020_Yellow_Taxi_Trip_Data.csv

2.04 GB of 2.23 GB

Search

🔊

📄

🔔

⚙️

?

azuser3611_mml.local@techademy.com
TECHADEMY LEARNING SOLUTIONS PRIVATE LIMITED

Other users in your workspace may have access to modify workspace.

New notebook New data flow New integration

Last Modified

.csv 7/14/2025, 10:32:48 AM

Upload details

Uploaded 1 files to 'synapse/'

Done

2020_Yellow_Taxi_Trip_Data.csv

Done

Upload completed

1 files uploaded to 'synapse/'. [View detail](#)

Microsoft Azure

Search resources, services, and docs (G+I)

Copilot

Home >

Microsoft.Azure.SynapseAnalytics.SparkPool-202507141121002332 | Overview

Deployment

Search

Delete Cancel Redeploy Download Refresh

Overview

Inputs

Outputs

Template

✔ Your deployment is complete

Deployment name : Microsoft.Azure.SynapseAnalytics.SparkPool-202507141121002332

Subscription : MML Learners

Resource group : rg-azuser3611_mml.local-fr&X

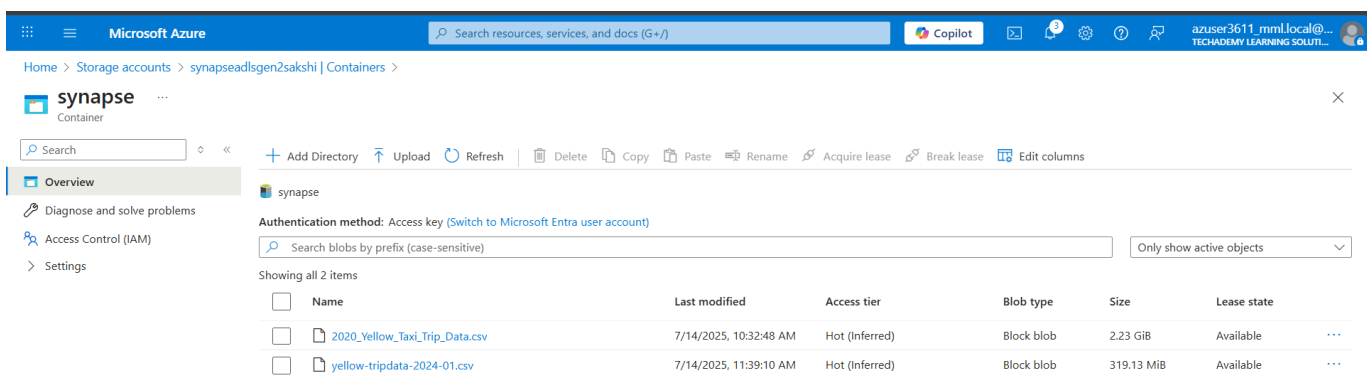
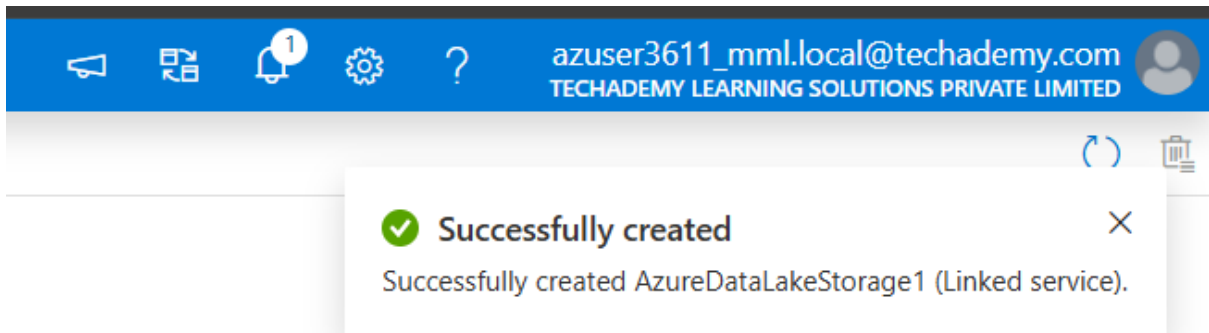
Start time : 7/14/2025, 11:21:03 AM

Correlation ID : a0bfec4e-65f4-4bbc-911e-5d66af73c0b9

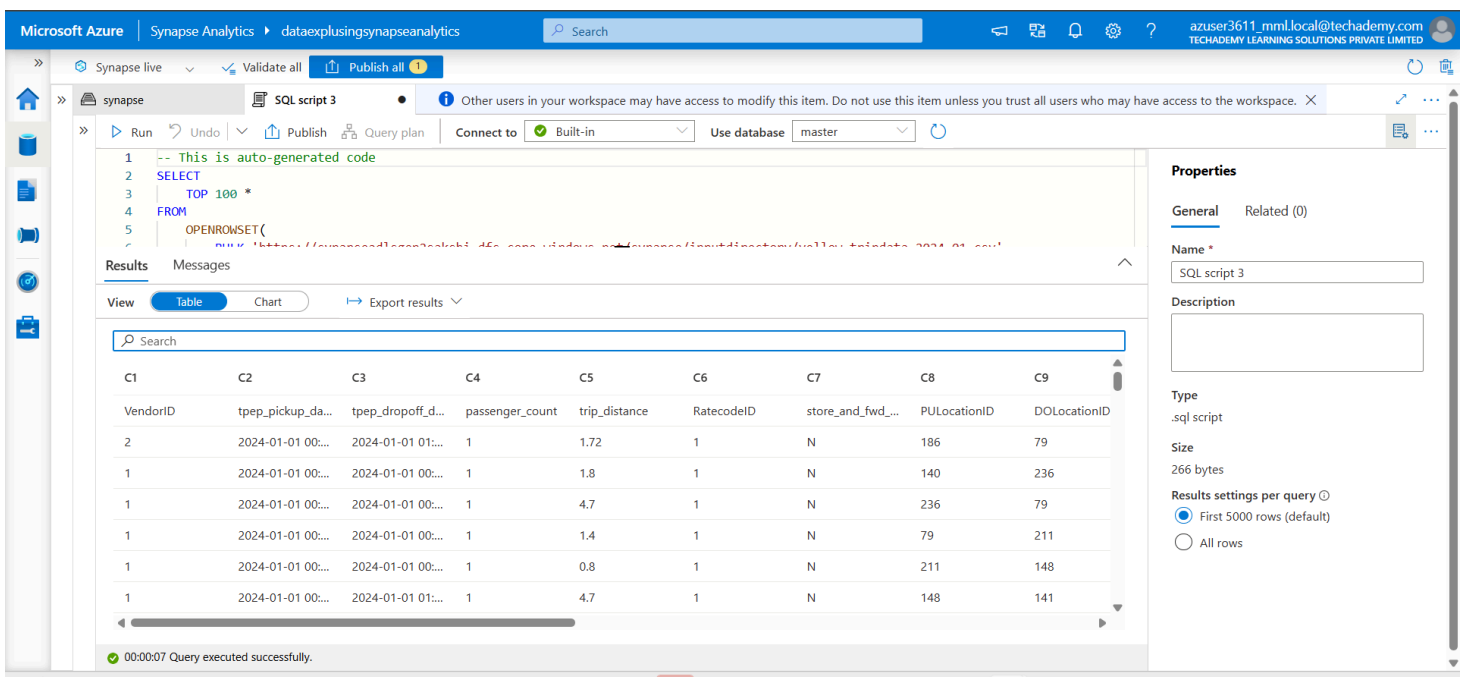
> Deployment details

< Next steps

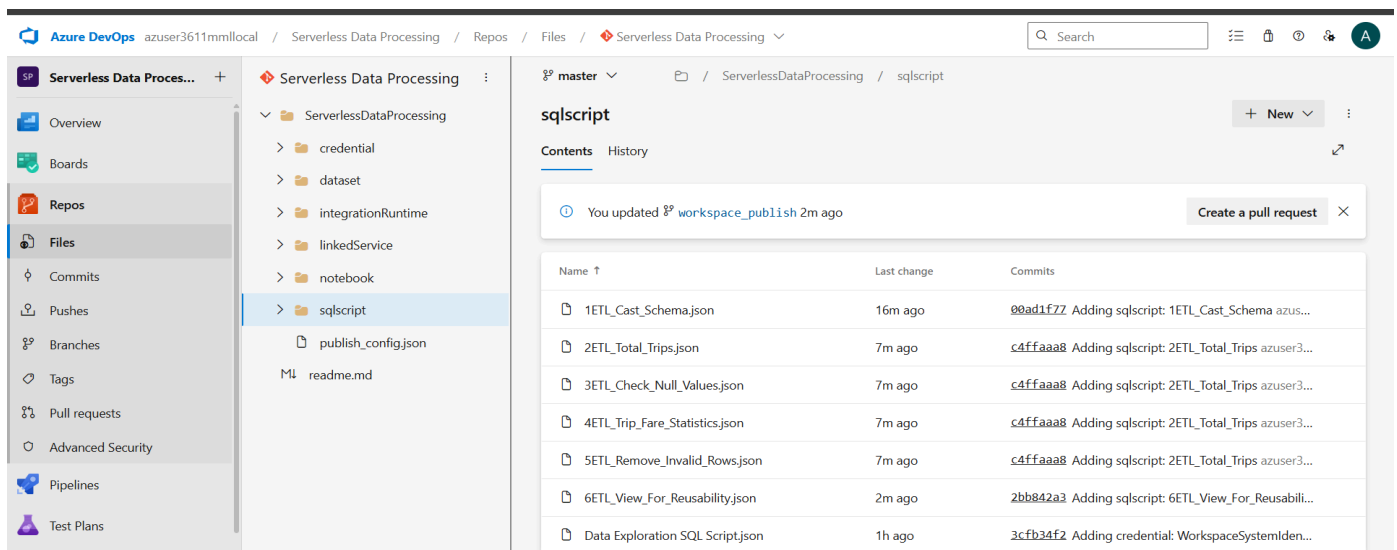
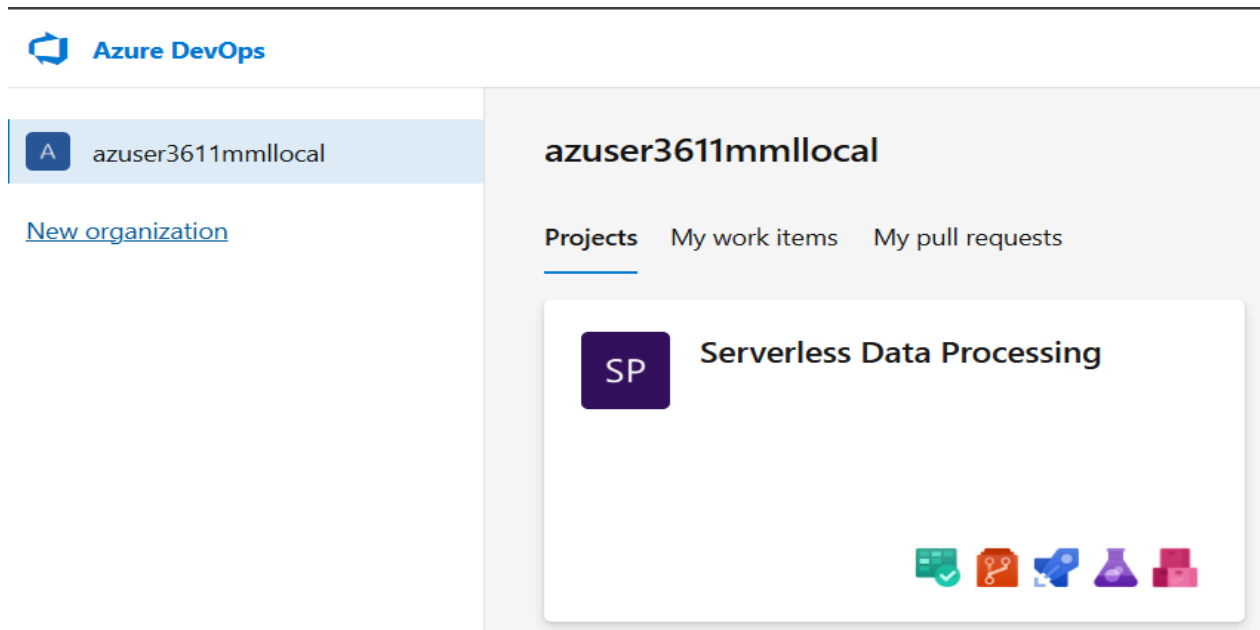
Go to resource



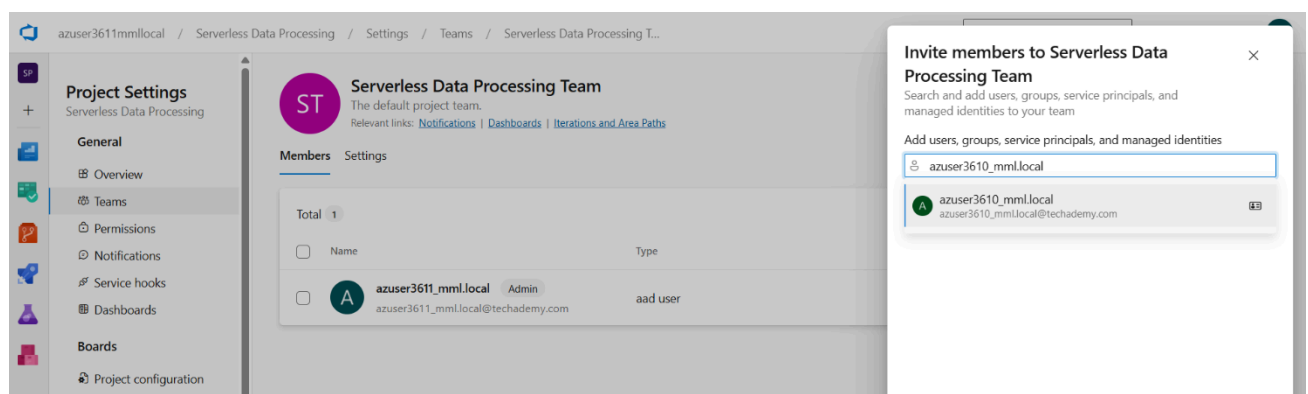
- **Fetching the first 100 rows using Synapse Serverless SQL from ADLS Gen2 dataset & ADLS Gen2 CSV file is accessible from Synapse Serverless**



Task 2) Git Configuration - Azure DevOps & Synapse Analytics



● Git Collaboration –



Task 3) Azure Databricks - Machine Learning

K-Means is a machine learning algorithm used to find groups (called clusters) in data.

Each cluster groups similar rows (data points) together.

Example:

In NYC Taxi data, we might want to group trips like:

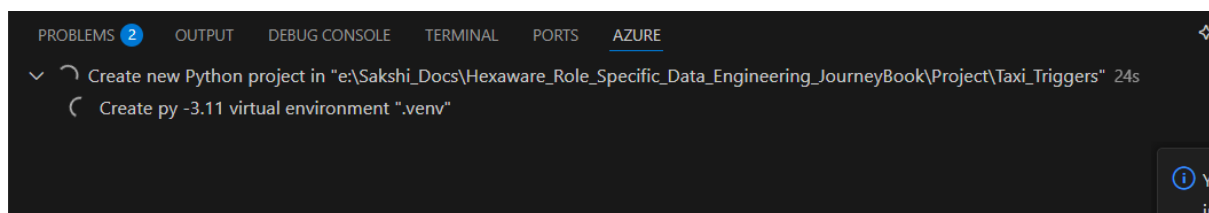
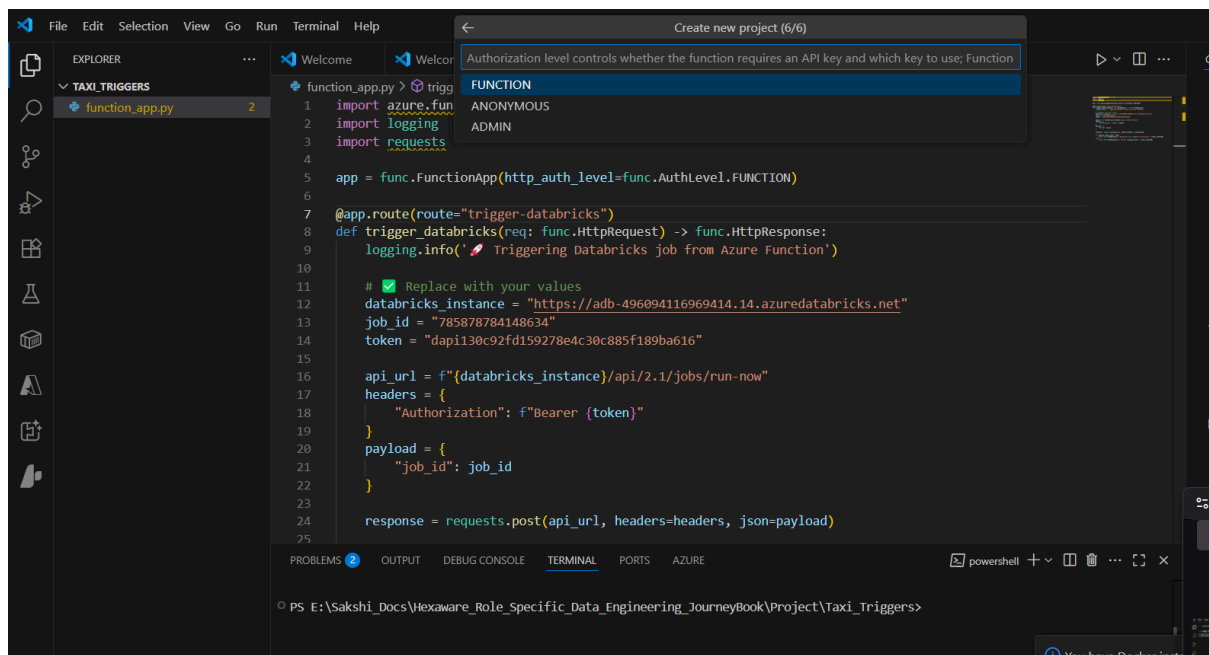
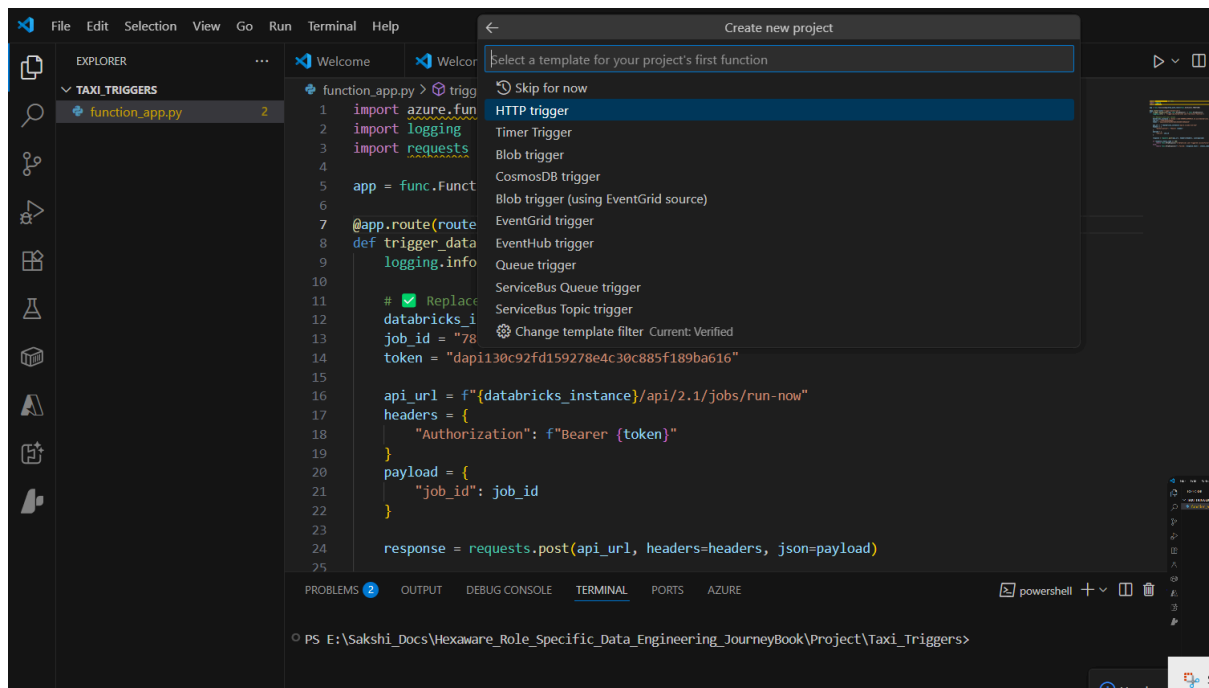
- Short trips with low fare
 - Long trips with high fare
 - Short trips with big tips
- These are called clusters.

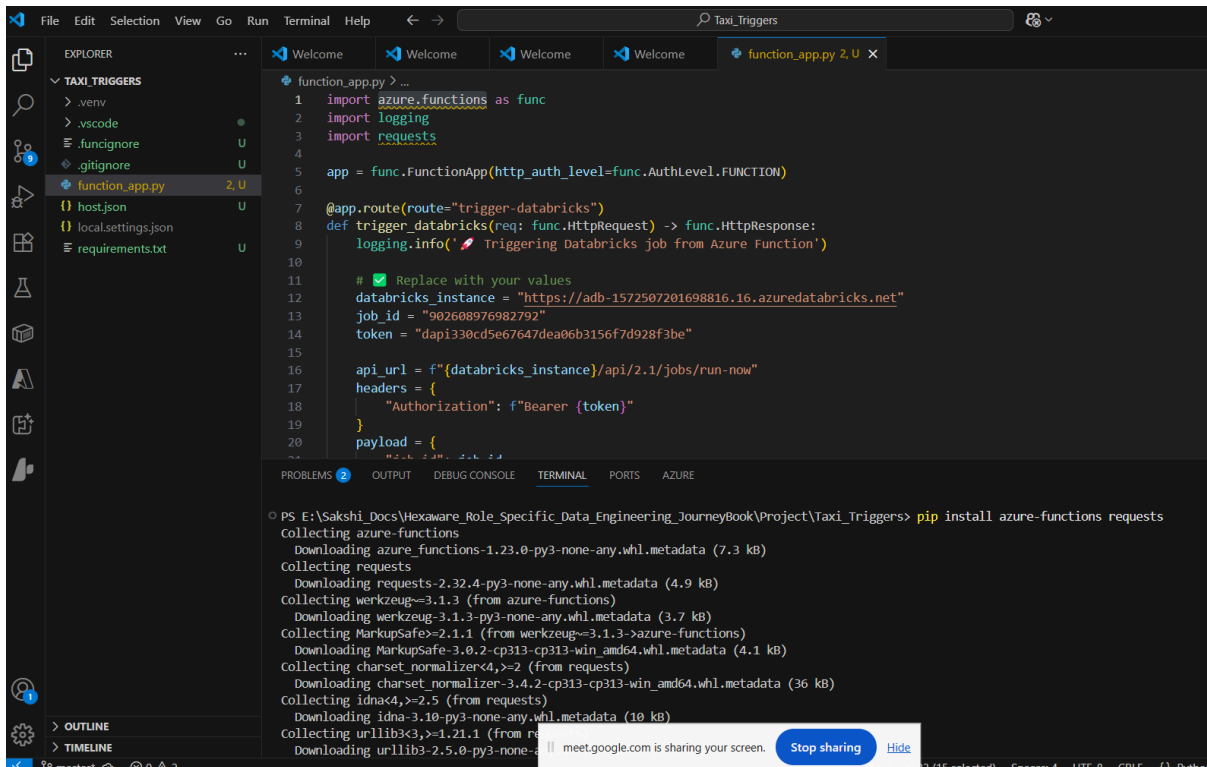
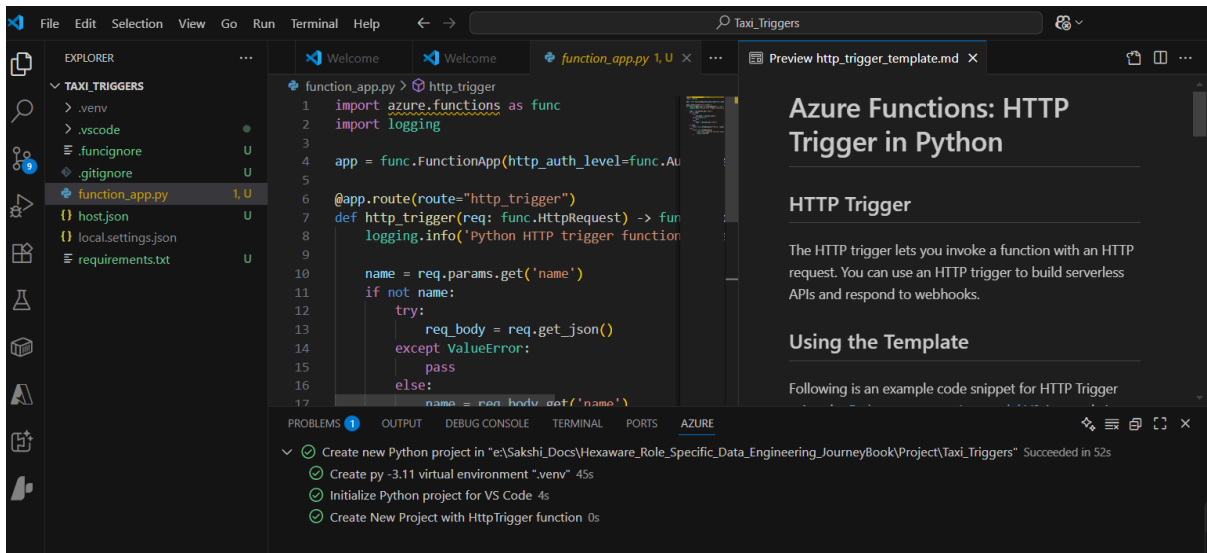
Column	Meaning
trip_distance	Distance of the taxi ride (in miles)
fare_amount	Fare for the ride (in \$)
passenger_count	No. of passengers during the trip
prediction	Which cluster the trip belongs to (0–3)

Interpreting The Results

Cluster	Approx % of Total	Likely Pattern
0	~82%	Regular trips — maybe 1–2 miles, 1 passenger, normal fare. Most common.
1	~10%	Possibly longer trips or group rides.
3	~7%	Could be airport trips or higher fare rides.

Task 4) Triggers -





The screenshot shows the Databricks Jobs & Pipelines interface. A modal window displays the details for a new job run titled "New Job Jul 15, 2025, 06:15 PM". The details include:

- Job ID: 902608976982792
- Job run ID: 315237759653411
- Launched: Manually
- Started: Jul 15, 2025, 07:09 PM
- Ended: Jul 15, 2025, 07:18 PM
- Duration: 7m 53s
- Queue duration: -
- Status: Succeeded

The background shows a table of runs with columns: Start time, Run ID, Launched, Duration, Spark, Status, Error code, and Run paramet... The first row shows a run from Jul 15, 2025, 07:09 PM with Run ID 3152377... and Status Succeeded.

“Serverless” here doesn’t mean *no servers at all*, it means:

- You don’t keep a server running all the time
- You only pay for what you use
- Resources spin up automatically when needed and shut down afterward

Task 5) MLFlow -

The screenshot shows a Databricks notebook with the following code:

```

1 import mlflow
2 import mlflow.spark
3
4 with mlflow.start_run():
5     k = 4
6     kmeans = KMeans(featuresCol="features", predictionCol="prediction", k=k, seed=42)
7     model = kmeans.fit(train_data)
8
9     predictions = model.transform(test_data)
10    silhouette = evaluator.evaluate(predictions)
11
12    mlflow.log_param("k", k)
13    mlflow.log_metric("silhouette_score", silhouette)
14
15    # Log the trained model
16    mlflow.spark.log_model(model, "kmeans_model")

```

Below the code, there is a section titled "(24) Spark Jobs" listing several jobs with their status and stage counts:

- Job 69: View (1 stage)
- Job 70: View (2 stages)
- Job 71: View (1 stage)
- Job 72: View (1 stage)
- Job 73: View (1 stage)
- Job 74: View (1 stage)
- Job 75: View (1 stage)
- Job 76: View (1 stage)
- Job 77: View (2 stages)
- Job 78: View (2 stages)

Microsoft Azure

databricks

Search data, notebooks, reents, and more...CTRL + P

New

Workspace

Recents

Catalog

Jobs & Pipelines

Compute

Data Engineering

Job Runs

AI/ML

Playground

Experiments >

NYC_Taxi_KMeans_Training

Machine learning

Runs

MLflow 3 is now available! Try out the new features and provide feedback. [Learn more](#)

Columns

Group by

metrics.silhouette_score

Run Name	Created	Dataset
wistful-dove-811	4 minutes ago	-

Search runs using a simplified version of the SQL WHERE clause. [Learn more](#)

Examples:

- metrics.rmse >= 0.8
- metrics.f1_score < 1
- params.model = 'tree'
- attributes.run_name = 'my run'
- tags.mlflow.user = 'myUser'
- metric.f1_score > 0.9 AND params.model = 'tree'
- dataset.name IN ('dataset1', 'dataset2')
- attributes.run_id IN ('a1b2c3d4', 'e5f6g7h8')
- tags.model_class LIKE 'sklearn.linear_model%'
- MIN(metrics.loss) < 1
- MAX(metrics.gpu_utilization_percentage) >= 0.5

Microsoft Azure

databricks

Search data, notebooks, reents, and more...CTRL + P

automateetlpimapiipelines

New

Workspace

Recents

Catalog

Jobs & Pipelines

Compute

Data Engineering

Job Runs

AI/ML

Playground

Experiments

Features

Models

Serving

Registered Models >

mlflow

Permissions

Use model for inference

Details

Notify me aboutAll new activity

Created Time: Jul 15, 2025, 11:55 PMLast Modified: Jul 15, 2025, 11:55 PMCreator: azuser3611_mmllocal@techademy.com

Description

Tags

Versions

All

Active 0

Compare

Version	Registered at	Created by	Stage	Pending Requests	Description	Endpoints
Version 1	Jul 15, 2025, 11:55 PM	azuser3611_mmlloc...	None	-		-

Microsoft Azure

databricks

Search data, notebooks, reents, and more...CTRL + P

automateet

New

Workspace

Recents

Catalog

Jobs & Pipelines

Compute

Data Engineering

Job Runs

AI/ML

Playground

Serving endpoints >

Create serving endpoint

Model serving is only available in Premium or Enterprise workspaces. Please contact your organization admin or Databricks support.

+ Add served entity

Route optimization

Enable route optimization

Summary

Served entities

mlflow

← ARM template deployment ⓘ

Azure Details ^

Deployment scope * ⓘ

Resource Group

Azure Resource Manager connection * ⓘ

MML Learners(2a3c6418-97b9-4d96-a2... ▾

Failed to set Azure permission 'RoleAssignmentId: 7db1cebc-7a85-45bd-a416-726e5a8d369b' for the service principal '26222bfd-76c4-4552-b044-e6c4cb900df0' on subscription ID '2a3c6418-97b9-4d96-a24b-2c2d7633d375': error code: Forbidden, inner error code: AuthorizationFailed, inner error message The client 'azuser3611_mml.local@techademy.com' with object id '81cfa09f-100c-42af-b3f5-382077f858a8' does not have authorization to perform action 'Microsoft.Authorization/roleAssignments/write' over scope '/subscriptions/2a3c6418-97b9-4d96-a24b-2c2d7633d375/providers/Microsoft.Authorization/roleAssignme7a85-45bd-a416-726e5a8d369b' or the scope is invalid. If access was recently granted, please refresh your credentials. Ensure that the user has 'Owner' or 'User Access Administrator' permissions on the Subscription.

About this task

Add




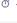


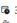
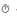

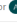
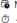



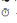
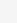
- SP Serverless Data Proces... +
- Overview
- Boards
- Repos
- Pipelines
- Pipelines
- Environments
- Library
- Test Plans
- Artifacts

Pipelines

Recent All Runs

New pipeline

Filter pipelines

Recently run pipelines		
Pipeline	Last run	
 Serverless Data Processing (4)	#20250715.2 • Update azure-pipelines.yml for Azure Pipelines Manually run by  master	 13m ago  <1s
 Serverless Data Processing (3)	#20250715.1 • Deleted SQL script 1.json Manually run by  master	 35m ago  <1s
 Serverless Data Processing (2)	#20250714.1 • Set up CI with Azure Pipelines PR automated for  azure-pipelines	 Monday  <1s
 Serverless Data Processing	#20250714.5 • Deleted SQL script 1.json Individual CI for  master	 Monday  <1s